

Classification of PDEs

consider a general 2nd Order PDE:

$$A u_{xx} + 2B u_{xy} + C u_{yy} = D$$

where $u(x,y)$ is the solution function, and A, B, C, D are smooth functions that depend on the independent variables x, y and on u, u_x, u_y (but not on the highest derivatives). The equation is nonlinear, but it is linear in its highest derivatives \rightarrow quasi-linear.

Laplace's eqn, $u_{xx} + u_{yy} = 0$, elliptic

"wave" eqn, $u_{xx} - u_{yy} = 0$, hyperbolic

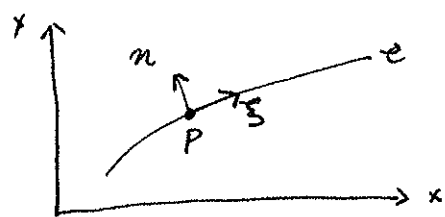
heat eqn, $u_{yy} = u_x$, parabolic

Burgers eqn, $u_{yy} = u_x + uu_y$

Tricomi eqn, $u_{xx} - xu_{yy} = 0$, changing type

The solution behavior depends on A, B, C primarily but also on D and on any boundary or initial data.

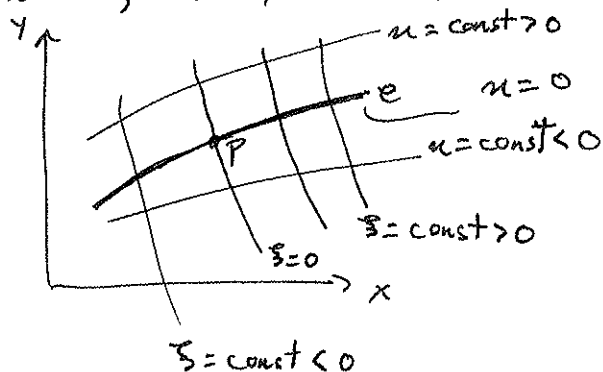
Classification: examine the local behavior of solutions about a point P given some "Cauchy" data. Ask the question, can the solution be extended beyond P .



c is a smooth curve
let $s(x,y)$ measure position along c
with $s = 0$ at the point P .
let $n(x,y)$ measure position away from c with $n = 0$ defining c

we can draw

contours, $u(x,y) = \text{const}$, $\xi(x,y) = \text{const}$.



we can define
the curve C as

$$C: u(x,y) = 0$$

assume that $\xi(x,y)$ and $u(x,y)$ are known about P .

Cauchy data: assume that u, u_y, u_x are given on C
(about P).

the task is to use the data given on C and the PDE to construct the solution in the neighborhood of P . If this can be done for any C about $P \Rightarrow$ solution is analytic, (ie has a Taylor series) and the equation is elliptic. If on the other hand there are curves C for which the solution can not be constructed off C then the solution need not be analytic and the equation is hyperbolic or parabolic.

use a Taylor series

$$u(\xi, \eta) = \underbrace{u(\xi, 0)}_{\text{given}} + \eta \underbrace{u_\eta(\xi, 0)}_{\text{given}} + \frac{\eta^2}{2} \underbrace{u_{\eta\eta}(\xi, 0)}_{\text{determine using data on } C \text{ and PDE}} + \dots$$

to compute $u_{\eta\eta}$, need $u_x, u_y, u_{xx}, u_{xy}, u_{yy}$ from chain rule

$$\rightarrow u_\eta = u_x x_\eta + u_y y_\eta$$

$$\rightarrow u_{\eta\eta} = (u_{xx} x_\eta + u_{xy} y_\eta) x_\eta + u_x x_{\eta\eta} + (u_{yx} x_\eta + u_{yy} y_\eta) y_\eta + u_y y_{\eta\eta}$$

Need to determine $u_x, u_y, u_{xx}, u_{xy}, u_{yy}$ ~~at P~~ along C .

$$u_x = u_\xi \xi_x + u_\eta \eta_x$$

$$u_y = u_\xi \xi_y + u_\eta \eta_y$$

on \mathcal{C} , we know u_ξ and u_η therefore u_x and u_y are known on \mathcal{C} from the given Cauchy problem

Now differentiate u_x, u_y along \mathcal{C}

$$\left. \begin{aligned} (u_x)_\xi &= u_{xx} \xi_\xi + u_{xy} \eta_\xi \\ (u_y)_\xi &= u_{yx} \xi_\xi + u_{yy} \eta_\xi \end{aligned} \right\} \begin{array}{l} \text{the LHS is known} \\ \text{from Cauchy data} \end{array}$$

use these two eqns with the PDE to determine u_{xx}, u_{xy} , and $u_{yy} \Rightarrow$ linear system

$$\begin{pmatrix} \xi_\xi & \eta_\xi & 0 \\ 0 & \xi_\eta & \eta_\eta \\ A & 2B & C \end{pmatrix} \begin{pmatrix} u_{xx} \\ u_{xy} \\ u_{yy} \end{pmatrix} = \begin{pmatrix} (u_x)_\xi \\ (u_y)_\xi \\ D \end{pmatrix}$$

The equation is solvable uniquely iff $\det(\text{matrix}) \neq 0$

$$\rightarrow A(\xi_\eta)^2 - 2B\xi_\xi\eta_\xi + C\xi_\xi^2 \neq 0 \quad (1)$$

If (1) holds for any smooth curve \mathcal{C} through P , then we can construct Taylor series about P and therefore the solution is analytic \Rightarrow ELLIPTIC CASE

Suppose \mathcal{C} is such that $A\xi_\eta^2 - 2B\xi_\xi\eta_\xi + C\xi_\xi^2 = 0$ then \mathcal{C} is called a characteristic.

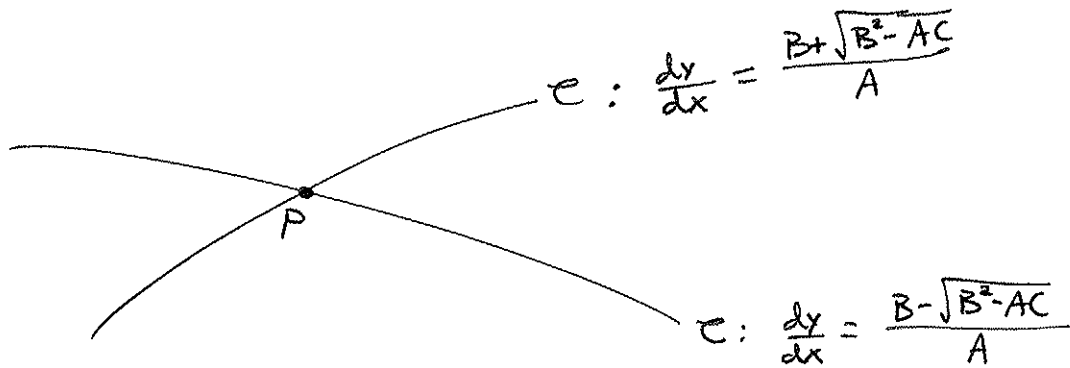
Suppose $A \neq 0$ and note that $\frac{\eta_\xi}{\xi_\xi} = \frac{dy}{dx}$ of \mathcal{C} at P

$$A\left(\frac{dy}{dx}\right)^2 - 2B\frac{dy}{dx} + C = 0$$

$$\rightarrow \frac{dy}{dx} = \frac{B \pm \sqrt{B^2 - AC}}{A}$$

→ there are 3 cases

1) $B^2 - AC \geq 0 \Rightarrow$ hyperbolic case



2) $B^2 - AC < 0 \Rightarrow$ elliptic case

no real characteristics, solutions are analytic

3) $B^2 - AC = 0 \Rightarrow$ parabolic case
coincident real characteristics

Characteristics

Consider the hyperbolic case and the linear system

$$\begin{pmatrix} x_\xi & y_\xi & 0 \\ 0 & x_\xi & y_\xi \\ A & B & C \end{pmatrix} \begin{pmatrix} u_{xx} \\ u_{xy} \\ u_{yy} \end{pmatrix} = \begin{pmatrix} u_{xy} \\ u_{y\xi} \\ D \end{pmatrix}$$

there are curves (ie characteristics) for which the determinant is zero. The system is singular so that solutions exist only if

$$\boxed{A x_\xi (u_x)_\xi + C x_\xi (u_y)_\xi = D x_\xi y_\xi}$$

we show this below \downarrow

$$(\alpha \ B \ \gamma) \begin{pmatrix} x_\xi & y_\xi & 0 \\ 0 & x_\xi & y_\xi \\ A & 2B & C \end{pmatrix} = 0$$

$$\Rightarrow (\alpha x_\xi + \gamma A, \alpha y_\xi + \beta x_\xi + 2\gamma B, \beta y_\xi + \gamma C) = 0$$

$$\text{pick } \alpha = \frac{-\gamma A}{x_\xi}, \quad \beta = \frac{-\gamma C}{y_\xi}$$

Must have $(K, B, \delta) \begin{pmatrix} u_{x\xi} \\ u_{y\xi} \\ D \end{pmatrix} = 0$

$\rightarrow d u_{x\xi} + B u_{y\xi} + \delta D = 0$

insert choices for d and B to obtain

$-\frac{\delta A}{x_\xi} (u_x)_\xi - \frac{\delta C}{y_\xi} (u_y)_\xi + \delta D = 0$

$\rightarrow A x_\xi (u_x)_\xi + C x_\xi (u_y)_\xi = D x_\xi y_\xi$ } holds along characteristics

Generally this equation is not solvable with simple analytic functions. There are some exceptions:

① wave eqn, $u_{xx} - u_{yy} = 0$, $A=1, B=0, C=-1, D=0$

$\rightarrow \frac{dy}{dx} = \frac{B \pm \sqrt{B^2 - AC}}{A} = \pm 1 \rightarrow \boxed{y = \pm x + \text{const characteristics}}$

now examine eqn that holds along characteristics

consider $y = x + \text{const}$

$\rightarrow x_\xi = x_\xi$

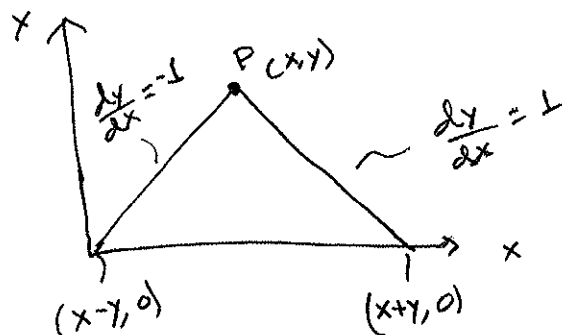
$\rightarrow x_\xi (u_x)_\xi - x_\xi (u_y)_\xi = 0 \rightarrow (u_x)_\xi - (u_y)_\xi = 0$

$\rightarrow u_x - u_y = f(x) = \text{constant along each characteristic}$
 $y = x + \text{const}$

likewise, $y = -x + \text{const} \Rightarrow u_x + u_y = \text{const.}$

use characteristics and Riemann variables to construct u

eg, $u_{xx} - u_{yy} = 0$, $|x| < \infty$, $y > 0$, $u(x, 0) = f(x)$, $u_y(x, 0) = g(x)$



$$(u_x + u_y)|_p = f'(x+y) + g(x+y)$$

$$(u_x - u_y)|_p = f'(x-y) - g(x-y)$$

add and subtract

$$u_x(x, y) = \frac{1}{2} f'(x-y) + \frac{1}{2} f'(x+y) + \frac{1}{2} g(x+y) - \frac{1}{2} g(x-y)$$

$$u_y(x, y) = \frac{1}{2} f'(x+y) - \frac{1}{2} f'(x-y) + \frac{1}{2} g(x+y) + \frac{1}{2} g(x-y)$$

integrate with respect to x

$$u(x, y) = \frac{1}{2} f(x+y) + \frac{1}{2} f(x-y) + \frac{1}{2} \int_{x-y}^{x+y} g(s) ds + h(y)$$

differentiate with respect to y , compare to $u_y(x, y)$

from above, and conclude that $h'(y) = 0$

$\rightarrow h(y) = 0$ and $u(x, 0)$ shows that $h(y) = 0$

to obtain d'Alembert's solution

Previously

considered the eqn $Au_{xx} + 2Bu_{xy} + Cu_{yy} = D$

where A, B, C, D depend on u, u_x, u_y .

considered local behavior near a given curve C where u and its normal derivative is specified.

hyperbolic case, $B^2 - AC > 0$, 2 real characteristics

elliptic case, $B^2 - AC < 0$, no real characteristics

parabolic case, $B^2 - AC = 0$, 1 real characteristic

We would like to show that equation types are invariant under coordinate transformation.

Introduce new independent variables (s, t) , ie $s = s(x, y)$, $t = t(x, y)$ with Jacobian non-zero (so the mapping isn't singular) and bounded.

$$J \equiv \frac{\partial(s, t)}{\partial(x, y)} = s_x t_y - s_y t_x \neq 0, \text{ (and bounded)}$$

The chain rule gives

$$\textcircled{2} \quad \alpha u_{ss} + 2\beta u_{st} + \gamma u_{tt} = \delta$$

where $\alpha = A s_x^2 + 2B s_x s_y + C s_y^2$

$$\beta = A s_x t_x + B(s_x t_y + s_y t_x) + C s_y t_y$$

$$\gamma = A t_x^2 + 2B t_x t_y + C t_y^2$$

You can show that

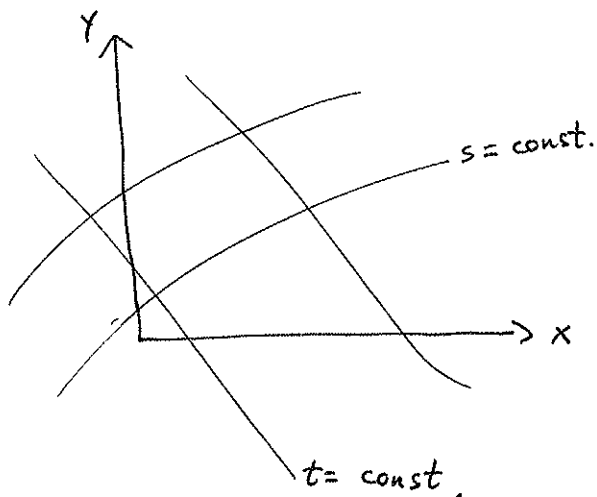
$$\beta^2 - \alpha\gamma = (B^2 - AC) \left(\frac{\partial(s, t)}{\partial(x, y)} \right)^2 \neq 0$$

\Rightarrow the sign of the discriminant is invariant \Rightarrow the equation type is unchanged

Suggests that we could choose a coordinate system to bring (2) to a certain canonical form.

Hyperbolic Case

Let $s(x,y) = \text{const.}$ be one family of characteristics and $t(x,y) = \text{const.}$ be the other family.



\Downarrow

$$A y_t^2 - 2 B x_t y_t + C x_t^2 = 0$$

$$A y_s^2 - 2 B x_s y_s + C x_s^2 = 0$$

eqn on $s = \text{const.}$
using t as
parameter

eqn for
characteristic
 $t = \text{const.}$, using
 s as parameter

use

use

$$\text{use } s_x = \frac{1}{J} y_t, \quad s_y = -\frac{1}{J} x_t, \quad t_x = -\frac{1}{J} y_s, \quad t_y = \frac{1}{J} x_s$$

$$\text{Recall, } d = A s_x^2 + 2 B s_x s_y + C s_y^2$$

$$d = A \left(\frac{1}{J} y_t \right)^2 + 2 B \left(\frac{1}{J} y_t \right) \left(-\frac{1}{J} x_t \right) + C \left(-\frac{1}{J} x_t \right)^2$$

$$d = \frac{1}{J^2} \left(\underbrace{A y_t^2 - 2 B x_t y_t + C x_t^2}_{\text{vanishes}} \right)$$

$$d = 0$$

Likewise, can show that $\delta = 0$

\Rightarrow Equation (2) reduces to the following hyperbolic canonical eqn

$$u_{st} = \delta$$

In the homogeneous problem where $\delta = 0$

$$u_{st} = 0 \rightarrow u = F(s) + G(t)$$

Another form: let $\hat{s} = s - t$
 $\hat{t} = s + t$

then obtain $u_{\hat{s}\hat{s}} - u_{\hat{t}\hat{t}} = 0$, the other ~~characteristic~~ ^{canonical} form

Elliptic Case

no real characteristics, however use real and imaginary parts of the complex characteristics

suppose $(s + it) = \text{const.}$ is a complex characteristic

suppose that the characteristic is parametrized by a real variable ξ , ie

$$s(x(\xi), y(\xi)) + it(x(\xi), y(\xi)) = \text{const.}$$

Differentiate:

$$(s_x + it_x)x_\xi + (s_y + it_y)y_\xi = 0$$

$$\rightarrow y_\xi = -\frac{(s_x + it_x)}{(s_y + it_y)} x_\xi$$

this complex characteristic is defined by

$$A y_\xi^2 - 2B x_\xi y_\xi + C x_\xi^2 = 0$$

Eliminate y_ξ :

$$A \left[-\frac{(s_x + it_x)}{(s_y + it_y)} x_\xi \right]^2 - 2B x_\xi \left[-\frac{(s_x + it_x)}{(s_y + it_y)} x_\xi \right] + C x_\xi^2 = 0$$

After some algebra, the real part becomes:

$$A s_x^2 + 2B s_x s_y + C s_y^2 = A t_x^2 + 2B t_x t_y + C t_y^2$$

\parallel \parallel
 χ χ

$$\Rightarrow \boxed{d = \chi, \text{ in this coordinate system}}$$

the imaginary part tells us that $B=0$ (in this coordinate system)

Therefore in this coordinate system, we have

$$\alpha u_{ss} + \gamma u_{tt} = \delta$$

with $\alpha=\gamma$, we have $u_{ss} + u_{tt} = \frac{\delta}{\alpha}$, Poisson's formula

The homogeneous eqn gives us $u_{ss} + u_{tt} = 0$, Laplace's eqn

Parabolic Case

Assume that $A \neq 0$ and take as our coordinates

$t = \text{const}$ is the characteristic

$$s = x$$

Can show that $\alpha = A \neq 0$, $B = \gamma = 0 \Rightarrow$

$$\boxed{u_{ss} = \bar{\delta} = \frac{\delta}{A}}$$

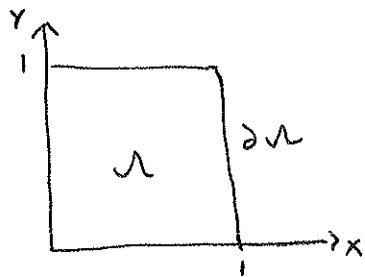
canonical form

Interesting case, $u_{ss} = u_t$, heat eqn

Boundary Conditions

Elliptic Case, the canonical form is Laplace's eqn.

Consider unit square as domain, and solve Laplace's eqn.



For a well-posed problem, need one condition on the whole perimeter. There are several forms:

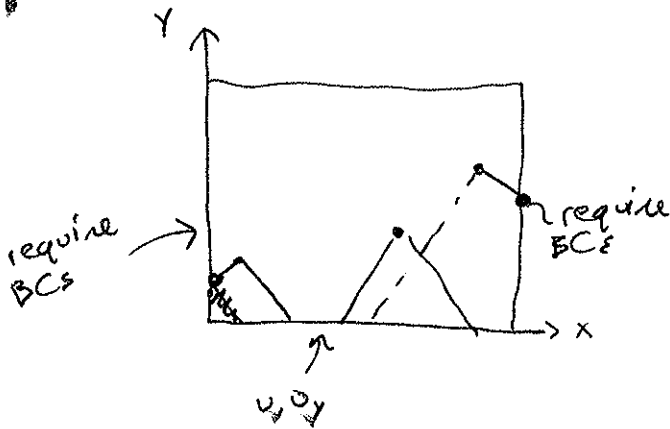
a) $u(x,y)$ given on $\partial\Omega$, (Dirichlet)

b) $\frac{\partial u}{\partial n} = \nabla u \cdot \hat{n}$ given on $\partial\Omega$, (Neumann)

c) $\alpha u + \beta \frac{\partial u}{\partial n}$ given on $\partial\Omega$, (Robin)

Hyperbolic Case (Boundary Conditions)

canonical form, $u_{xx} - u_{yy} = 0$



characteristics, $x+y = \text{const}$
 $x-y = \text{const}$

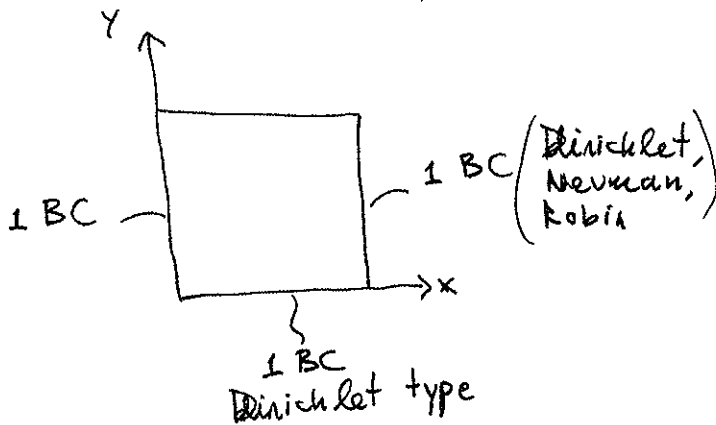
require two conditions on $y=0$
typically specify u, u_y given

domain of dependence / region of influence



Parabolic Case (BCs)

canonical form, $u_y = u_{xx}$, heat eqn



First Order Eqns

$$Au_x + Bu_y = C$$

where A, B, C may depend on x, y, u (at most) \Rightarrow quasi-linear

Always have one family of characteristics \Rightarrow hyperbolic.

Suppose ξ is a parameter along the characteristic.

$$\text{Set } \frac{dx}{d\xi} = A(x, y, u), \quad \frac{dy}{d\xi} = B(x, y, u)$$

These two differential eqns define a path, where ξ measures position along the path. On the characteristic

$$u(x, y) = u(x(\xi), y(\xi))$$

$$\rightarrow \frac{du}{d\xi} = u_x \frac{dx}{d\xi} + u_y \frac{dy}{d\xi} = Au_x + Bu_y = C$$

$$\rightarrow \text{Along the characteristic, } \frac{du}{d\xi} = C(x, y, u)$$

Reduced the PDE to 3 ODEs.

First Order Systems

$$\underline{A} u_x + \underline{B} u_y = \underline{C}$$

The system may be elliptic, hyperbolic, parabolic, or a mix.
We'll talk about this later.

Finite Difference Methods For PDEs

Begin with a simple example involving heat equation.

$$u_t = \gamma u_{xx}, \quad 0 \leq x \leq 1, \quad t > 0, \quad \gamma = \text{constant}$$

$$u(x, 0) = f(x), \quad u(0, t) = a(t), \quad u(1, t) = b(t)$$

Assume $a(0) = f(0)$ and $b(0) = f(1)$ to enforce continuity, though not important due to diffusive nature of heat eqn.

Introduce a mesh

$$x_j = j\Delta x, \quad \Delta x = \frac{1}{N}$$

$$0 = x_0 < x_1 < x_2 < \dots < x_N = 1$$

} spatial discretization.

define $v_j(t) \approx u(x_j, t)$ ($v_j(t)$ approximates exact solution u at $x = x_j$ and time t)

Approximate derivatives

the usual thing for the heat eqn is to take a centered appx.

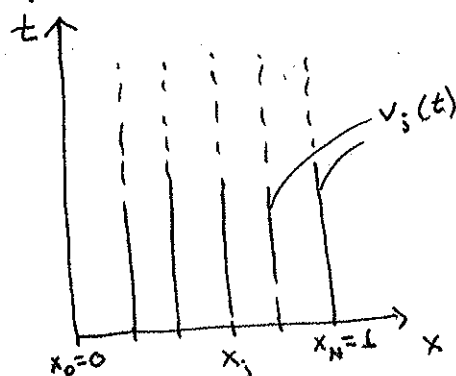
$$u_{xx}(x_j, t) \approx \frac{1}{\Delta x^2} (v_{j-1}(t) - 2v_j(t) + v_{j+1}(t)) \quad - \text{2nd order centered appx}$$

$$u_t(x_j, t) \approx v_j'(t)$$

$$\Rightarrow v_j'(t) = \frac{\gamma}{\Delta x^2} (v_{j-1}(t) - 2v_j(t) + v_{j+1}(t)), \quad j = 1, \dots, N-1$$

$$v_0(t) = a(t), \quad v_N(t) = b(t), \quad v_j(0) = f(x_j), \quad j = 1, \dots, N-1$$

we have changed PDE into a bunch of ODEs
apply a method to solve ODEs numerically



"Method of lines" - discretize eqns in space first to obtain an ODE in time then use ODE solver to integrate ODEs

9/4/03

Solve the ODEs using Forward Euler;
ie replace $v'(t)$ with Forward difference

$$v'(t) \approx \frac{v_j(t+\Delta t) - v_j(t)}{\Delta t}, \text{ where } \Delta t \text{ is chosen}$$

$$\rightarrow \frac{v_j(t+\Delta t) - v_j(t)}{\Delta t} \approx v \left[\frac{v_{j-1}(t) - 2v_j(t) + v_{j+1}(t)}{\Delta x^2} \right]$$

Suppose Δt is constant and $t_n = n\Delta t$ then

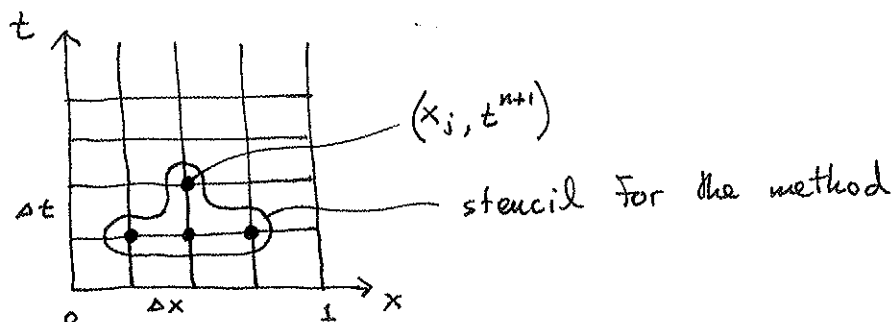
define $v_j^n \approx v_j(t_n)$

$$\Rightarrow \frac{v_j^{n+1} - v_j^n}{\Delta t} = v \left[\frac{v_{j-1}(t) - 2v_j(t) + v_{j+1}(t)}{\Delta x^2} \right]$$

Finite
difference
method

$$v_0^n = a(t_n), \quad v_N^n = b(t_n), \quad v_j^0 = f(x_j)$$

we now have
a full grid,
(mesh)



you can solve the Finite difference method explicitly
here is the algorithm:

1) specify $N, v, \Delta t$

2) $v_j^0 = f(x_j)$ for $j = 0, \dots, N$

3) For $n = 0, 1, \dots, n_{\text{final}}$

4) $v_j^{n+1} = v_j^n + \frac{\Delta t v}{\Delta x^2} (v_{j-1}^n - 2v_j^n + v_{j+1}^n)$ for $j = 1, \dots, N-1$

5) $v_0^{n+1} = a(t_{n+1}), \quad v_N^{n+1} = b(t_{n+1})$

Example:

$$u_t = \nu u_{xx}, \quad 0 < x < 1, \quad t > 0$$

$$u(x, 0) = x + \sin \pi x + \sin 3\pi x$$

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

Solve using $v_j^{n+1} = v_j^n + r(v_{j-1}^n - 2v_j^n + v_{j+1}^n)$, $r = \frac{\nu \Delta t}{\Delta x^2}$

Exact solution is

$$u(x, t) = x + e^{-\nu \pi^2 t} \sin \pi x + e^{-9\nu \pi^2 t} \sin 3\pi x$$

Solve numerically For $N = 20, 40, 80$ (spatial discretizations)
and for various choices of Δt .

Issues

- 1) accuracy \rightarrow error
- 2) instabilities? \rightarrow choice for Δt
- 3) cost?

For $N = 20$, $r = 0.4$, $\max(\text{error}) = 9.84 \times 10^{-3}$

$N = 40$, $r = 0.4$, $\max(\text{error}) = 2.43 \times 10^{-3}$

$N = 80$, $r = 0.4$, $\max(\text{error}) = 6.05 \times 10^{-4}$

notice $\Delta x = \frac{1}{N}$ is halved and the error decreases by a factor of 4, which suggests that the order decreases by Δx^2

$N = 40$, $r = 0.6$, \Rightarrow instability
 \rightarrow limited by how big Δt can be

Consistency, Order of Accuracy

Previous example demonstrated that maximum error $\sim O(\Delta x^2)$ for $r = \frac{v \Delta t}{\Delta x^2}$ Fixed less than some limiting value. Why?

$$\frac{v_j^{n+1} - v_j^n}{\Delta t} = v \left(\frac{v_{j-1}^n - 2v_j^n + v_{j+1}^n}{\Delta x^2} \right)$$

Introduce difference operators: (linear operators)

$$\delta_{+x} v_j^n = v_{j+1}^n - v_j^n, \text{ Forward difference operator}$$

$$\delta_{-x} v_j^n = v_j^n - v_{j-1}^n, \text{ backward difference operator}$$

$$\delta_{0x} v_j^n = v_{j+1}^n - 2v_j^n + v_{j-1}^n, \text{ centered difference operator}$$

$$\begin{aligned} \delta_x^2 v_j^n &= \delta_{+x} \delta_{-x} v_j^n \\ &= v_{j+1}^n - 2v_j^n + v_{j-1}^n, \end{aligned} \quad \begin{array}{l} \text{2nd order centered} \\ \text{difference operator} \end{array}$$

can also define in time

$$\delta_{+t} v_j^n = v_j^{n+1} - v_j^n$$

and so on

Using this notation, our finite difference method becomes

$$\boxed{\frac{1}{\Delta t} \delta_{+t} v_j^n = \frac{v}{\Delta x^2} \delta_x^2 v_j^n}$$

Introduce the error grid function

$$e_j^n = v_j^n - u(x_j, t_n)$$

$$\boxed{e_j^n = v_j^n - u_j^n}, \text{ For notational convenience}$$

$$\text{set } v_j^n = u_j^n + e_j^n$$

substitute this into Finite difference method

$$\rightarrow \frac{1}{\Delta t} \delta_{+t} u_j^n + \frac{1}{\Delta t} \delta_{+t} e_j^n = \frac{\gamma}{\Delta x^2} \delta_x^2 u_j^n + \frac{\gamma}{\Delta x^2} \delta_x^2 e_j^n$$

$$\rightarrow \left[\frac{1}{\Delta t} \delta_{+t} e_j^n - \frac{\gamma}{\Delta x^2} \delta_x^2 e_j^n \right] = - \left[\frac{1}{\Delta t} \delta_{+t} u_j^n - \frac{\gamma}{\Delta x^2} \delta_x^2 u_j^n \right] \quad \text{error equation}$$

$$e_j^0 = v_j^0 - u(x_j, 0) = 0$$

$$e_0^n = v_0^n - u(0, t_n) = 0$$

$$e_N^n = v_N^n - u(x_N, t_n) = 0$$

◆ Examine the RHS of the error equation

$$\begin{aligned} \frac{1}{\Delta t} \delta_{+t} u_j^n &= \frac{u(x_j, t_n + \Delta t) - u(x_j, t_n)}{\Delta t} \\ &= \frac{1}{\Delta t} \left[\Delta t u_t(x_j, t_n) + \frac{\Delta t^2}{2} u_{tt}(x_j, t_n) + \dots \right] \\ &= u_t(x_j, t_n) + \frac{\Delta t}{2} u_{tt}(x_j, t_n) + O(\Delta t^2) \end{aligned}$$

$$\begin{aligned} \frac{1}{\Delta x^2} \delta_x^2 u_j^n &= \frac{1}{\Delta x^2} \left[u(x_j - \Delta x, t_n) - 2u(x_j, t_n) + u(x_j + \Delta x, t_n) \right] \\ &= \frac{1}{\Delta x^2} \left[\begin{aligned} &\cancel{u(x_j, t_n)} - \cancel{\Delta x u_x} + \frac{\Delta x^2}{2} u_{xx} - \frac{\Delta x^3}{6} u_{xxx} + \frac{\Delta x^4}{24} u_{xxxx} + \dots \\ &\cancel{-2u} + \cancel{u} + \cancel{\Delta x u_x} + \frac{\Delta x^2}{2} u_{xx} + \frac{\Delta x^3}{6} u_{xxx} + \frac{\Delta x^4}{24} u_{xxxx} + \dots \end{aligned} \right] \\ &= u_{xx}(x_j, t_n) + \frac{\Delta x^2}{12} u_{xxxx}(x_j, t_n) + O(\Delta x^4) \end{aligned}$$

The error equation becomes

$$\frac{1}{\Delta t} \delta_{+t} e_j^n - \frac{\nu}{\Delta x^2} \delta_x^2 e_j^n = - \left[\begin{aligned} &u_t + \frac{\Delta t}{2} u_{tt} + O(\Delta t^2) \\ &- \nu \left(u_{xx} + \frac{\Delta x^2}{12} u_{xxxx} + O(\Delta x^4) \right) \end{aligned} \right]_{(x_j, t_n)}$$

since $u(x_j, t_n)$ solves the differential equation, what remains is the negative of the truncation error

$$\frac{1}{\Delta t} \delta_{+t} e_j^n - \frac{\nu}{\Delta x^2} \delta_x^2 e_j^n = - \underbrace{\left[\frac{\Delta t}{2} u_{tt} - \nu \frac{\Delta x^2}{12} u_{xxxx} + O(\Delta t^2, \Delta x^4) \right]}_{\text{local truncation error}}$$

local truncation error: amount by which the exact solution fails to satisfy the difference approximation

In our example, $r = \frac{\nu \Delta t}{\Delta x^2} = \text{const} \rightarrow \Delta t = \frac{r \Delta x^2}{\nu}$

Then the right hand side becomes

$$- \left[\frac{r \Delta x^2}{2\nu} u_{tt} - \frac{\nu \Delta x^2}{12} u_{xxxx} + \dots \right]$$

$$= - \frac{\Delta x^2}{2\nu} \left[r u_{tt} - \frac{\nu^2}{6} u_{xxxx} + \dots \right]$$

$$= O(\Delta x^2)$$

So if the solution is well-behaved then the size of e_j^n would be $O(\Delta x^2)$. This suggests that

$$\max_{j,n} |e_j^n| = O(\Delta x^2)$$

For our difference equation, the local truncation error is

$$\tau_j^n = \frac{1}{\Delta t} \delta_{+t} u_j^n - \frac{\nu}{\Delta x^2} \delta_x^2 u_j^n$$

A difference method is said to be consistent if

$$\max_{j,n} |\tau_j^n| \rightarrow 0 \text{ as } \Delta x, \Delta t \rightarrow 0$$

So for our example

$$\tau_j^n = \frac{\Delta t}{2} u_t - \frac{\nu \Delta x^2}{12} u_{xxxx} + \dots$$

$$= O(\Delta t, \Delta x^2) \rightarrow 0 \text{ as } \Delta t, \Delta x \rightarrow 0$$

The order of accuracy refers to the rate at which the truncation error τ_j^n goes to zero.

So for our example, the order of accuracy is first order in time and second order in space.

For our example, we had that

$$e_j^n \approx \Delta x^2 E(x_j, t_n)$$

\uparrow amplitude \uparrow "shape of plot of error"

Go back to the error equation

$$\frac{1}{\Delta t} \delta_{+t} e_j^n - \frac{\nu}{\Delta x^2} \delta_x^2 e_j^n = \Delta x^2 F_j^n + \dots, \text{ with ICs \& BCs}$$

$$\rightarrow \frac{1}{\Delta t} \delta_{+t} \left(\frac{e_j^n}{\Delta x^2} \right) - \frac{\nu}{\Delta x^2} \delta_x^2 \left(\frac{e_j^n}{\Delta x^2} \right) = F_j^n + \dots$$

as $\Delta x \rightarrow 0$ for $r = \text{fixed}$

$$E_t - \nu E_{xx} = F + \text{homogeneous ICs, BCs}$$

Computational Cost

number of Floating point operations required to compute v_j^n to a time $t_{\text{final}} = O(1)$

number of Flops needed per time step is $O(N)$

number of time steps is $O(1/\Delta t)$

if r is fixed then $O(1/\Delta t) = O(\Delta x^2) = O(1/N^2)$

$\rightarrow \text{cost} = O(N^3)$

Best you could hope for is $O(N^2)$ i.e. one flop per grid point

A Neumann Problem For the Heat Equation

9/8/03

$$u_t = \nu u_{xx}, \quad 0 < x < 1, \quad t > 0$$

$$\text{IC: } u(x, 0) = f(x)$$

$$\text{BC: } u_x(0, t) = 0, \quad u_x(1, t) = 0 \quad - \text{homogeneous Neumann BCs}$$

Let us pursue a finite difference approximation of the problem.

Approximate the solution, $v_j(t) \approx u(x_j, t)$

Require a finite difference eqn for $v_j(t)$:

$$\rightarrow v_j'(t) = \frac{\nu}{\Delta x^2} \delta_x^2 v_j(t), \quad \text{recall, } \delta_x^2 v_j(t) = v_{j-1} - 2v_j + v_{j+1}$$

This equation is defined for some range of the subscript j .

In the Dirichlet case, did not have to solve at $j=0, N$

The difficulty in the Neumann case is that $u(x, t)$ is not specified on the boundary.

For example at $j=1$,

$$v_1'(t) = \frac{\gamma}{\Delta x^2} (v_2 - 2v_1 + v_0)$$

In the Dirichlet case, v_0 was given, but this is not the situation in the Neumann case.

We need to solve

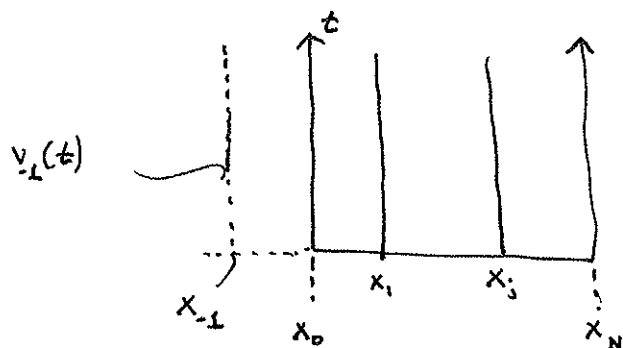
$$v_j'(t) = \frac{\gamma}{\Delta x^2} \delta_x^2 v_j(t), \quad j=0, \dots, N$$

But now at $j=0$,

$$v_0'(t) = \frac{\gamma}{\Delta x^2} (v_1 - 2v_0 + v_{-1})$$

So we now have the negative subscript!

The Fix is to introduce a "ghost" line to handle the Neumann B.C.



consider the boundary condition

$$v_x(0, t) = 0$$

approximate the derivative using a finite difference equation

$$v_x \Big|_{x=0} \approx \frac{1}{2\Delta x} \delta_{0x} v_0(t)$$

- centered difference, $\frac{v_1(t) - v_{-1}(t)}{2\Delta x}$

the discrete BC is

$$\frac{1}{2\Delta x} \delta_{0x} v_0(t) = 0$$

\rightarrow

$$\delta_{0x} v_0(t) = 0$$

$$\rightarrow \boxed{v_{-1}(t) = v_1(t)}$$

Therefore at $j=0$:

$$v_0'(t) = \frac{\nu}{\Delta x^2} (v_1 - 2v_0 + v_{-1}) = \frac{\nu}{\Delta x^2} (2v_1 - 2v_0)$$

Can choose to make this substitution into the FDE, in the program or just do the loop and use extra equation - just a matter of personal choice.

Similarly near $x=1$,

$$u_x(1,t) = 0 \rightarrow \frac{1}{2\Delta x} \delta_{0x} v_N = 0 \rightarrow \boxed{v_{N+1} = v_{N-1}}$$

The result of all this is a set of ODEs, solving for $v_j(t)$, $j=0, \dots, N$:

$$v_0' = \frac{\nu}{\Delta x^2} (2v_1 - 2v_0)$$

$$v_j' = \frac{\nu}{\Delta x^2} (v_{j+1} - 2v_j + v_{j-1}), \quad j=1, \dots, N-1$$

$$v_N' = \frac{\nu}{\Delta x^2} (2v_{N-1} - 2v_N)$$

with ICs $v_j(0) = f(x_j)$, $j=0, \dots, N$

Now integrate the ODEs using method of your choice.

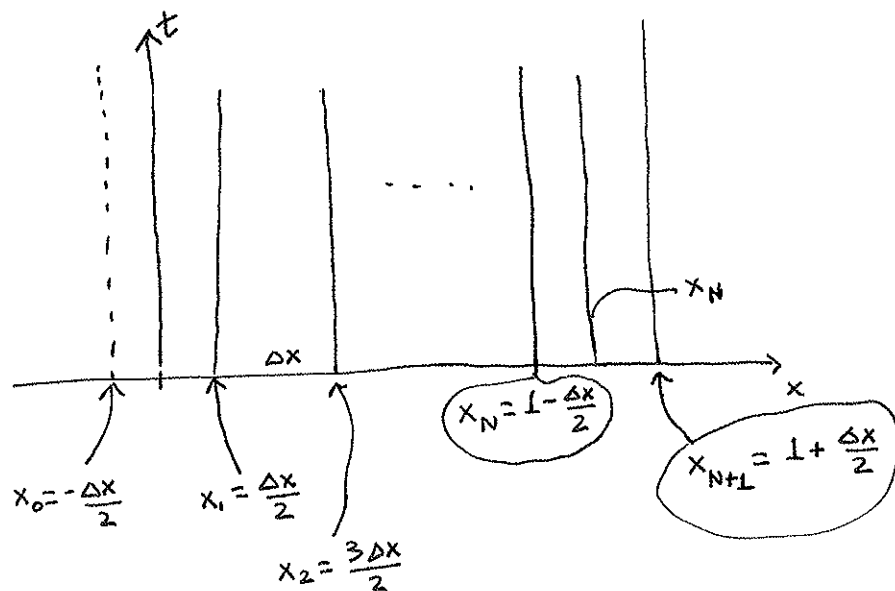
Alternate approach

$$u_t = \nu u_{xx}, \quad 0 < x < 1, \quad t > 0$$

$$u(x,0) = f(x), \quad u_x(0,t) = u_x(1,t) = 0$$

Let us discuss the use of a staggered grid:

$$x_j = (j - \frac{1}{2}) \Delta x, \quad \Delta x = \frac{1}{N}$$



Use $v_j'(t) = \frac{\gamma}{\Delta x^2} \delta_{0x} v_j$, $j = 1, \dots, N$ (interior grid lines)

At $x=0$, need (would like) a centered approximation

$$v_x(0, t) = 0 \rightarrow \frac{1}{\Delta x} \delta_{-x} v_1(t) = 0 \quad \text{recall, } \delta_{-x} = \frac{v_1(t) - v_0(t)}{\Delta x}$$

$$\rightarrow \boxed{v_0(t) = v_1(t)}$$

$$\rightarrow \boxed{v_1'(t) = \frac{\gamma}{\Delta x^2} (v_2 - 2v_1 + v_0) = \frac{\gamma}{\Delta x^2} (v_2 - v_1)} \quad \left(\text{notice the factor of 2 is gone} \right)$$

and similarly on the other side, $\boxed{v_{N+1} = v_N}$

$$\rightarrow \boxed{v_N'(t) = \frac{\gamma}{\Delta x^2} (v_{N-1} - v_N)}$$

Also Integral Conservation For Neumann Problem

Integrate the heat equation with respect to x from $x=0$ to $x=1$.

$$\int_0^1 v_t dx = \gamma \int_0^1 v_{xx} dx$$

$\rightarrow 0$, due to Neumann BCs

$$\frac{d}{dt} \int_0^1 v dx = \left[v_x \right]_0^1$$

\rightarrow

$$\boxed{\begin{aligned} \frac{d}{dt} \int_0^1 v dx &= 0 \\ \rightarrow \int_0^1 v dx &= \text{const}, \end{aligned}}$$

Question: In what sense is this physical rule preserved in the discrete approximation?

Let us define $I(t) = \int_0^1 u(x, t) dx$

and define the discrete approximation (trapezoidal rule)

$$I_h(t) = \frac{\Delta x}{2} v_0(t) + \Delta x \sum_{j=1}^{N-1} v_j(t) + \frac{\Delta x}{2} v_N(t)$$

(this is for the non-staggered approximation)

Let us show that $I_h = \text{constant}$. Calculate $I_h'(t)$

$$I_h'(t) = \frac{\Delta x}{2} v_0'(t) + \Delta x \sum_{j=1}^{N-1} v_j'(t) + \frac{\Delta x}{2} v_N'(t)$$

$$I_h'(t) = \frac{\Delta x}{2} \cdot \frac{v}{\Delta x^2} (2v_1 - 2v_0) + \Delta x \sum_{j=1}^{N-1} \frac{v}{\Delta x^2} (v_{j+1} - 2v_j + v_{j-1}) + \frac{\Delta x}{2} (2v_{N-1} - 2v_N) \frac{v}{\Delta x^2}$$

$$I_h'(t) = 0$$

→

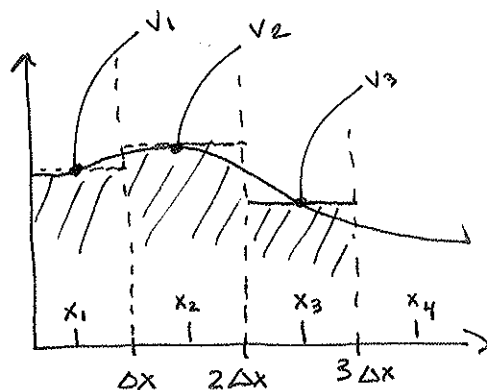
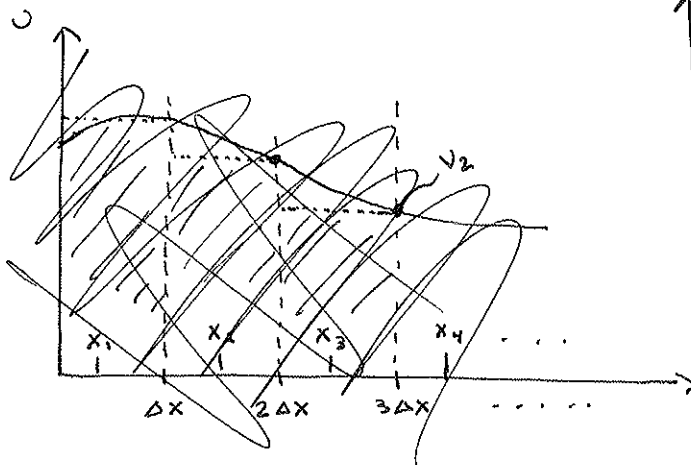
$$I_h(t) = \text{constant}$$

→

$$I_h(t) = I_h(0)$$

The staggered configuration also preserves an integral.

Define $\hat{I}_h(t) = \Delta x \sum_{j=1}^N v_j(t)$, Midpoint rule



Generating Discrete Approximations

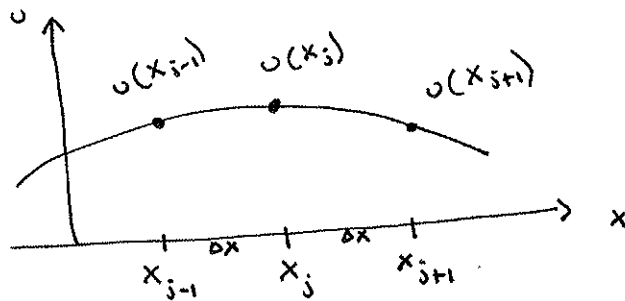
Basic problem is to find discrete approximations to derivatives that appear in PDE.

eg For the second derivative in the heat eqn, we had

$$u_{xx}(x_j, t) \approx \frac{u(x_{j-1}, t) - 2u(x_j, t) + u(x_{j+1}, t))}{\Delta x^2}$$

$$x_{j-1} = x_j - \Delta x, \quad x_{j+1} = x_j + \Delta x$$

Suppose we want to approximate $u_x(x_j)$



truncation
error

$$\text{let } u_x(x_j) = \underbrace{a u(x_{j-1}) + b u(x_j) + c u(x_{j+1}))}_{\substack{\text{Finite difference appx} \\ \text{of } u_x(x_j) \text{ for some choice} \\ \text{of coefficients } a, b, c}} + \tau(x_j)$$

How do you choose a, b, c ?

- 1) we want $\tau(x_j) \rightarrow 0$ as $\max\{\Delta x_j, \Delta x_{j+1}\} \rightarrow 0$
(this is the notion of consistency)
- 2) we would like $\tau(x_j) \rightarrow 0$ as fast as possible, ie this is the notion of maximizing order of an approximation
- 3) you may want the approximation to preserve some features of the PDE

$$u_x(x_j) = a u(x_j - \Delta x_j) + b u(x_j) + c u(x_j + \Delta x_{j+1}) + \tau(x_j)$$

- For simplicity, set $\Delta x_j = \Delta x_{j+1} = \Delta x$
- Also note that a, b, c are independent of x_j
we can set $x_j = 0$ without loss of generality

$$\rightarrow u_x(0) = a u(-\Delta x) + b u(0) + c u(\Delta x) + \tau$$

Taylor series approach

$$u_x(0) = a \left(u(0) - \Delta x u_x(0) + \frac{\Delta x^2}{2} u_{xx}(0) + \dots \right) + b u(0) + c \left(u(0) + \Delta x u_x(0) + \frac{\Delta x^2}{2} u_{xx}(0) + \dots \right) + \tau$$

equate the coeffs of derivatives

$$\left. \begin{array}{l} u(0): \quad a + b + c = 0 \\ u_x(0): \quad -\Delta x a + c \Delta x = 1 \\ u_{xx}(0): \quad \frac{\Delta x^2}{2} a + \frac{\Delta x^2}{2} c = 0 \end{array} \right\} \text{linear system for coefficients } a, b, c$$

$$\rightarrow \boxed{a = \frac{-1}{2\Delta x}, \quad b = 0, \quad c = \frac{1}{2\Delta x}}$$

The unmatched terms are absorbed into τ

$$\begin{aligned} \rightarrow u_x(0) = & \frac{-1}{2\Delta x} \left(u(0) - \Delta x u_x(0) + \frac{\Delta x^2}{2} u_{xx}(0) - \frac{\Delta x^3}{6} u_{xxx}(0) + \frac{\Delta x^4}{24} u_{xxxx}(0) + \dots \right) \\ & + \frac{1}{2\Delta x} \left(u(0) + \Delta x u_x(0) + \frac{\Delta x^2}{2} u_{xx}(0) + \frac{\Delta x^3}{6} u_{xxx}(0) + \frac{\Delta x^4}{24} u_{xxxx}(0) + \dots \right) + \tau \end{aligned}$$

$$\rightarrow u_x(0) = u_x(0) + \frac{\Delta x^2}{6} u_{xxx}(0) + \frac{\Delta x^4}{5!} u_{xxxx}(0) + \dots + \tau$$

$$\rightarrow \boxed{\tau = -\frac{\Delta x^2}{3!} u_{xxx} - \frac{\Delta x^4}{5!} u_{xxxx} + \dots}$$

Notice, $\tau \rightarrow 0$ as $\Delta x \rightarrow 0$. The order of the approximation is 2.

As another example, consider the second derivative

$$u_{xx} = a(u - \Delta x u_x + \frac{\Delta x^2}{2} u_{xx} + \dots) + bu + c(u + \Delta x u_x + \frac{\Delta x^2}{2} u_{xx} + \dots) + \tau$$

match the corresponding coefficients

$$a + b + c = 0$$

$$-\Delta x a + \Delta x c = 0$$

$$\frac{\Delta x^2}{2} a + \frac{\Delta x^2}{2} c = 1$$

 \Rightarrow

$$a = c = \frac{1}{\Delta x^2}, \quad b = -\frac{2}{\Delta x^2}$$

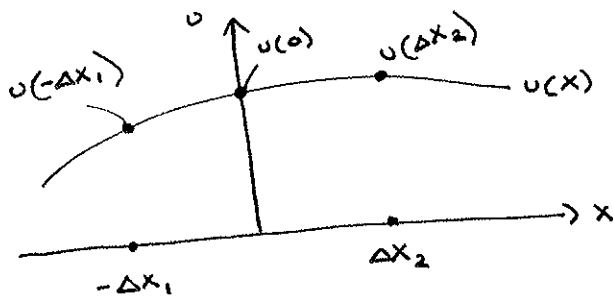
You find that $\tau = -2 \left[\frac{\Delta x^2}{4!} u_{xxx} + \frac{\Delta x^4}{6!} u_{xxxxx} + \dots \right]$

This idea can be used to find one-sided approximations, higher order appxs (requiring more points)

9/11/03

Generating Discrete Approximations

- 1) Taylor series approach
- 2) Interpolation approach:



consider some smooth function $u(x)$, assume $u(0)$ is known. To approximate derivative at $x=0$, need local information. Note, Δx_1 and Δx_2 need not be equal. In addition, notice that the curve is relatively smooth.

Interpolate $u(x)$ using data ~~($-\Delta x_1, u(-\Delta x_1)$)~~
 $(-\Delta x_1, u(-\Delta x_1))$, $(0, u(0))$, $(\Delta x_2, u(\Delta x_2))$

call the interpolant $\tilde{u}(x)$. Suppose we want to approximate $u_x(0)$, then use $u_x(0) \approx \tilde{u}_x(0)$

Standard approach is polynomial interpolation.

Set

$$u_x(0) = \underbrace{A u(-\Delta x_1) + B u(0) + C u(\Delta x_2)}_{\tilde{u}_x(0)} + E \quad \rightarrow \text{truncation error}$$

this is equivalent to setting truncation error $E=0$

For $u(x) = 1, x, x^2$

$$u=1: \quad 0 = A + B + C$$

$$u=x: \quad 1 = -A\Delta x_1 + C\Delta x_2$$

$$u=x^2: \quad 0 = A\Delta x_1^2 + C\Delta x_2^2$$

Solve the linear equations to find

$$A = \frac{-1}{\Delta x_1 \left(\frac{\Delta x_1}{\Delta x_2} + 1 \right)}, \quad C = \frac{1}{\Delta x_2 \left(\frac{\Delta x_2}{\Delta x_1} + 1 \right)}, \quad B = -A - C$$

Notice, if $\Delta x_1 = \Delta x_2 = \Delta x$ then $A = \frac{-1}{2\Delta x}$, $C = \frac{1}{2\Delta x}$, $B = 0$

and so $u_x(0) \approx \frac{1}{2\Delta x} (u(\Delta x) - u(-\Delta x))$, which is the same formula derived by Taylor series approach.

What about the error? Need to consider higher ~~order~~ degree polynomials, $u = x^3, x^4, \dots$. You let $E = K u^{(p)}(\xi)$ for $-\Delta x_1 \leq \xi \leq \Delta x_2$

Finite Volumes

Finite volume discretization is an integral form of an equation.

$$\frac{d}{dt} \int_a^b A e(x,t) dx = \text{rate of change of 'heat' energy between } x=a, x=b$$

where $A = \text{constant cross-sectional area}$
 $e(x,t) = \text{'heat energy' per unit volume}$

Define heat Flux $\Phi(x, t)$, $\left(\frac{\text{heat}}{\text{area} \cdot \text{time}}\right)$

define heat generation $f(x, t) = \left(\frac{\text{heat}}{\text{volume} \cdot \text{time}}\right)$

$$\frac{d}{dt} \int_a^b A e(x, t) dx = A [\Phi(a, t) - \Phi(b, t)] + \int_a^b A f(x, t) dx$$

$A = \text{constant}$, cancel out, Formula is per unit length now

$$\frac{d}{dt} \int_a^b e(x, t) dx = \Phi(a, t) - \Phi(b, t) + \int_a^b f(x, t) dx$$

Discretize this formula:

~~$a = x_{j-1}$~~ choose $a = x_{j-1}$
 $b = x_j$

define $E_j(t) \equiv \frac{1}{\Delta x} \int_{x_{j-1}}^{x_j} e(x, t) dx$, cell average energy

$F_j(t) \equiv \frac{1}{\Delta x} \int_{x_{j-1}}^{x_j} f(x, t) dx$, cell average of heat generation term

$\rightarrow \frac{d}{dt} E_j(t) = - \left(\frac{\Phi(x_j, t) - \Phi(x_{j-1}, t)}{\Delta x} \right) + F_j(t)$ this is still an exact Formula

Define $u(x, t)$ to be a temperature

set $e(x, t) = \rho(x) c(x) u(x, t)$

density $\rho(x)$ specific heat $c(x)$

Set $\Phi(x, t) = -k(x) u_x$

heat conductivity $k(x)$

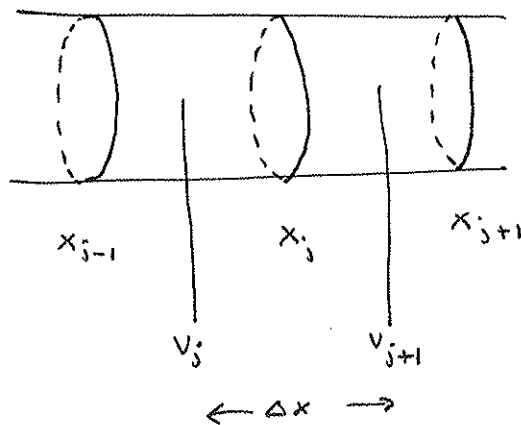
→ $E_j(t)$ is the cell average of $e(x,t)$

$$E_j(t) = \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_j} \rho(x) C(x) u(x,t) dx, \text{ use midpoint rule}$$

$$E_j(t) = \rho(x_{j-1/2}) C(x_{j-1/2}) u(x_{j-1/2}, t) + O(\Delta x^2)$$

Define $v_j(t) \approx u(x_{j-1/2}, t)$, approximation to temperature

$$\Phi(x_j, t) \approx -K(x_j) \left(\frac{v_{j+1}(t) - v_j(t)}{\Delta x} \right)$$



The discrete equation becomes, for some range of j

Conservation Form $\left\{ \rho(x_{j-1/2}) C(x_{j-1/2}) \frac{d}{dt} v_j(t) = -\frac{1}{\Delta x} \left[K(x_j) \frac{v_{j+1} - v_j}{\Delta x} - K(x_{j-1}) \frac{v_j - v_{j-1}}{\Delta x} \right] + F_j(t) \right.$

These are now a bunch of ODEs.

Convergence, Consistency, stability

Differential equation, $Lu = F$

~~where~~ Approximate the differential equation using $L_h u_h = F_h$

where L_h is the difference operator

u_h is the grid function

F_h is the Forcing Function

h denotes the grid spacing

For example,

$$\frac{1}{\Delta t} \delta_{+t} v_j^n - \frac{\nu}{\Delta x^2} \delta_x^2 v_j^n = F(x_j, t_n)$$

The basic idea of convergence is that $v_h \rightarrow u_h$ as $h \rightarrow 0$.

Suppose we want an initial value problem. Let us not concern ourselves with boundary conditions.

$$Lu = F, \quad |x| < \infty, \quad t > 0, \quad u(x, 0) = F(x)$$

The approximation is

$$L_h u_h = F_h, \quad |j| < \infty, \quad n > 0, \quad u_h = F(x_j) \text{ when } n = 0$$

$$\text{with } x_j = j \Delta x, \quad t_n = n \Delta t$$

Def: The difference approximation is said to be convergent

$$\text{if } \|v_h - u_h\|_h \rightarrow 0 \text{ as } h \rightarrow 0$$

For any $x, t \in [0, t_{\text{final}}]$

$\|\cdot\|_h$ is some norm that is defined on the grid,

$$\text{For example the infinity norm } \|e_h\|_\infty = \max_j |e(x_j, t)|$$

The discrete approximation is convergent of order (P, q)

$$\text{if } \|v_h - u_h\| = O(\Delta x^P, \Delta t^q)$$

To show that the method is convergent, often consider whether the method is consistent and stable. However, you can in some instances show convergence directly.

Example : discretization of heat eqn.

$$\frac{1}{\Delta t} \delta_{+t} v_j^n = \frac{\nu}{\Delta x^2} \delta_x^2 v_j^n$$

show that the approximation is a convergent approximation of

$$u_t = \nu u_{xx}$$

provided that $r = \frac{\nu \Delta t}{\Delta x^2} \leq \frac{1}{2}$

$$\text{Set } e_j^n = v_j^n - u(x_j, t_n), \quad x_j = j \Delta x, \quad t_n = n \Delta t$$

$$\text{denote } u(x_j, t_n) = u_j^n \quad \text{then } e_j^n = v_j^n - u_j^n$$

$\rightarrow v_j^n = u_j^n + e_j^n$, substitute into Finite difference approximation

$$\frac{1}{\Delta t} \delta_{+t} (u_j^n + e_j^n) = \frac{\nu}{\Delta x^2} \delta_x^2 (u_j^n + e_j^n)$$

$$\Rightarrow \frac{1}{\Delta t} \delta_{+t} e_j^n - \frac{\nu}{\Delta x^2} \delta_x^2 e_j^n = - \left[\frac{1}{\Delta t} \delta_{+t} u_j^n - \frac{\nu}{\Delta x^2} \delta_x^2 u_j^n \right]$$

We have shown (9/4/03) that

$$\frac{1}{\Delta t} \delta_{+t} u_j^n - \frac{\nu}{\Delta x^2} \delta_x^2 u_j^n = O(\Delta t, \Delta x^2)$$

$$\Rightarrow \delta_{+t} e_j^n - r \delta_x^2 e_j^n = O(\Delta t^2, \overset{\Delta t}{\Delta x^2})$$

$$\rightarrow e_j^{n+1} = e_j^n + r \delta_x^2 e_j^n + O(\Delta t^2, \overset{\Delta t}{\Delta x^2})$$

Take the absolute value, use Δ inequality

$$e_j^{n+1} = e_j^n + r(e_{j-1}^n - 2re_j^n + e_{j+1}^n) + O(\Delta t^2, \Delta x^2)$$

$$\rightarrow |e_j^{n+1}| \leq |1-2r||e_j^n| + r|e_{j-1}^n| + r|e_{j+1}^n| + A(\Delta t^2, \Delta x^2)$$

if $r \leq \frac{1}{2}$, $|1-2r| = 1-2r$

define $E^n = \max_j |e_j^n|$

$$|e_j^{n+1}| \leq (1-2r)E^n + rE^n + rE^n + A(\Delta t^2, \Delta x^2)$$

$\rightarrow |e_j^{n+1}| \leq E^n + A(\Delta t^2, \Delta t \Delta x^2)$ For all j and in particular
For that value of j where e is max

$$\rightarrow E^{n+1} \leq E^n + A(\Delta t^2, \Delta t \Delta x^2)$$

$$\leq E^{n-1} + 2A(\Delta t^2, \Delta t \Delta x^2)$$

$$\leq E^{n-2} + 3A(\Delta t^2, \Delta t \Delta x^2)$$

\vdots

$$\leq E^0 + (n+1)A(\Delta t^2, \Delta t \Delta x^2)$$

but $E^0 = 0$, therefore

$$E^n \leq n A(\Delta t^2, \Delta t \Delta x^2)$$

but $n \Delta t \leq t_{\text{final}} \rightarrow$

$$E^n \leq K(\Delta t, \Delta x^2) \rightarrow \text{convergent}$$

9/15/03

We're discussing convergence, consistency, stability

PDE, initial value problem, $L u = F$, $|x| < \infty$, $t > 0$

initial condition, $u(x, 0) = f(x)$

we have a discrete approximation, $L_h v_h = F_h$

For some grid (x_j, t_n) , $x_j = j \Delta x$, $t_n = n \Delta t$, $|j| < \infty$, $n > 0$

$v_h = f(x_j)$ at $n = 0$

The idea of convergence is that the numerical approximation approaches the exact solution restricted to the grid as $h \rightarrow 0$

$$v_h \rightarrow u_h \text{ as } h \rightarrow 0$$

that was the case for a pure IVP, but if you have a pure BVP,

the idea is essentially the same with a few ' wrinkles'.

- Consider the BVP,
$$\begin{aligned} L u &= F, & 0 < x < 1 \\ u(x, 0) &= f(x) + \text{BCs at } x=0 \text{ and } x=1 \end{aligned}$$

Absorb the BCs into the problem.

→
$$\begin{aligned} \hat{L} u &= \hat{F}, & 0 \leq x \leq 1, t > 0 \\ u(x, 0) &= f(x) \end{aligned}$$
 (notice problem now includes boundary)

Example: $u_t = \nu u_{xx} + F(x, t)$

$u(0, t) = a(t)$

$u_x(1, t) - u(1, t) = b(t)$

Now define \hat{L} and \hat{F}

$$\hat{L} u = \begin{cases} u_t - \nu u_{xx} & , 0 < x < 1 \\ 0 & , x = 0 \\ u_x - u & , x = 1 \end{cases}, \quad \hat{F} = \begin{cases} F & , 0 < x < 1 \\ a & , x = 0 \\ b & , x = 1 \end{cases}$$

discrete approximation: $\hat{L}_h v_h = \hat{F}_h$, $0 \leq j \leq N$, $h > 0$

$$v_h = F(x_j), \quad h = 0$$

to define convergence, introduce a sequence $\{\Delta x_k\}_{k=1}^{\infty}$ such that $\Delta x_k \rightarrow 0$ as $k \rightarrow \infty$. This is required due to finite interval.

Therefore we need $V_{h_k} \rightarrow U_{h_k}$ as $k \rightarrow \infty$

So the idea between pure IVPs and BVPs (concerning convergence) is the same!

Consistency

The idea for consistency is that the discrete approximation should approach the PDE (+BCs) as $h \rightarrow 0$.

Def A discrete approximation $L_h v_h = F_h$ of a differential equation $Lu = F$ is consistent if

$$\|(L\phi - F)_h - (L_h \phi_h - F_h)\|_h \rightarrow 0$$

as $h \rightarrow 0$ For any smooth function ϕ .

Note, ϕ is not necessarily a solution of PDE.

However, often take $\phi = u$, the solution of the PDE.

then $L\phi - F = 0$ and then

$$L_h \phi_h - F_h = L_h u_h - F_h = -\tau_h, \text{ the local truncation error}$$

Therefore this definition of consistency really amounts to whether local truncation error vanishes as grid is refined

$$\|\tau_h\| \rightarrow 0 \text{ as } h \rightarrow 0$$

→ The amount by which the exact solution fails to satisfy the difference approximation must tend to zero as $h \rightarrow 0$.

The order of accuracy of an approximation depends on the rate at which $\tau_h \rightarrow 0$.

eg: if $\|\tau_h\|_h = O(\Delta x^p, \Delta t^q)$

then the order of accuracy is p^{th} order in space and q^{th} order in time.

Example consider the discretization

$$\frac{1}{\Delta t} \delta_{+t} v_j^n = \frac{\nu}{2\Delta x^2} \delta_x^2 (v_j^n + v_j^{n+1}) \quad \text{Crank-Nicholson method}$$

show that the discrete approximation is consistent with the equation $u_t = \nu u_{xx}$ and determine the order of accuracy.

Define truncation error, τ_j^n

$$\tau_j^n = \frac{1}{\Delta t} \delta_{+t} u(x_j, t_n) - \frac{\nu}{2\Delta x^2} \delta_x^2 (u(x_j, t_n) + u(x_j, t_{n+1}))$$

First term:

$$\delta_{+t} u_j^n = u_j^{n+1} - u_j^n = \Delta t u_t + \frac{\Delta t^2}{2} u_{tt} + \frac{\Delta t^3}{6} u_{ttt} + \dots$$

$$\rightarrow \frac{1}{\Delta t} \delta_{+t} u_j^n = u_t + \frac{\Delta t}{2} u_{tt} + \frac{\Delta t^2}{6} u_{ttt} + \dots$$

Second term:

$$\frac{1}{\Delta x^2} \delta_x^2 u_j^n = \frac{1}{\Delta x^2} (u_{j+1}^n - 2u_j^n + u_{j-1}^n)$$

$$= \frac{1}{\Delta x^2} \left(\Delta x^2 u_{xx} + \frac{\Delta x^4}{12} u_{xxxx} + O(\Delta x^6) \right)$$

$$= u_{xx} + \frac{\Delta x^2}{12} u_{xxxx} + O(\Delta x^4)$$

Third term:

$$\frac{1}{\Delta x^2} \delta_x^2 u_j^{n+1} = \left[u_{xx} + \frac{\Delta x^2}{12} u_{xxxx} + O(\Delta x^4) \right]_{(x_j, t_{n+1})}$$

$$= u_{xx} + \Delta t u_{xxt} + \frac{\Delta t^2}{2} u_{xxtt} + O(\Delta t^3)$$

$$+ \frac{\Delta x^2}{12} u_{xxxx} + \frac{\Delta x^2 \Delta t}{12} u_{xxxxt} + O(\Delta x^2 \Delta t^2) + O(\Delta x^2, \Delta t^2)$$

$$+ O(\Delta x^4)$$

$$\rightarrow \tau_j^n = \cancel{u_t} + \frac{\Delta t}{2} u_{tt} + \frac{\Delta t^2}{6} u_{ttt} - \frac{\nu}{2} \left[\cancel{2u_{xx}} + \cancel{\Delta t u_{xxt}} + \frac{2\Delta x^2}{12} u_{xxxx} + O(\Delta t^2, \Delta x^2 \Delta t, \Delta x^4) \right]$$

$\cancel{u_{tt}} = \nu u_{xxt}$

$\rightarrow \tau_j^n = O(\Delta t^2, \Delta x^2) \rightarrow$ therefore consistent because as $\Delta x \rightarrow 0, \Delta t \rightarrow 0$, $\tau_j^n \rightarrow 0$ and 2nd order in both time and space

Stability

Stability is a form of well-posedness for difference equations.

Initial-value problems: $Lu = F, |x| < \infty, t > 0, u(x, 0) = f(x)$

discrete approximation: $L_h v_h = F_h, |j| < \infty, n > 0, v_h = f(x_j)$ at $n=0$

Suppose the problems are linear. Define an error e_h ,

$$e_h = v_h - u_h$$

Error solves $L_h e_h = F_h - L_h u_h = \tau_h$, with $e_h = 0$ at $n=0$

If the discrete approximation is consistent then $\tau_h \rightarrow 0$ as $h \rightarrow 0$.

We would like that $e_h \rightarrow 0$ as $h \rightarrow 0$, which would imply convergence.

For stability, we need only consider the homogeneous problem, $L_h v_h = 0$, v_h given at $n=0$.

Include a superscript n to denote time:

$$L_h v_h^n = 0, \text{ with } v_h^0 \text{ given}$$

We require this equation to be well-behaved for stability.

Def A difference scheme $L_h v_h^n = 0$ is stable if $k \geq 0$, $\beta \geq 0$ exist such that

$$\|v_h^n\|_h \leq \|v_h^0\|_h K e^{\beta t}, \quad t = n \Delta t$$

Suppose you're only interested in integrating to some finite time, i.e. you restrict $t \in (0, t_{\text{final}})$ then

$$\|v_h\|_h \leq \hat{K} \|v_h^0\|_h$$

or if the difference approximation is meant to be an approximation of a differential equation whose solutions are "bounded", "stable", then you might apply a stability criterion such as

$$\|v_h^n\|_h \leq K \|v_h^0\| \text{ for all } t$$

The bottom line is that there are various definitions of stability that make sense depending on the situation.

Often, time-stepping schemes are two level schemes of the form $v_h^{n+1} = Q_h v_h^n$, where Q_h is some difference operator. The general form for stability for the two level schemes are

$$\|Q_h^n\| \leq K e^{\beta t}, \quad t = n \Delta t \quad \text{or} \quad \|Q_h\| \leq 1 + \alpha \Delta t, \quad \alpha = \text{const.}$$

where

$$\|Q_h\| = \max \frac{\|Q_h v_h\|}{\|v_h\|} \text{ for all } v_h$$

Reduce the question of stability to a consideration of the behavior of the stepping function.

Example: consider the scheme $v_j^{n+1} = v_j^n - \sigma \delta_x v_j^n$

Show that this scheme is stable if $0 \leq \sigma \leq 1$

Take absolute value of both sides:

$$|v_j^{n+1}| \leq |v_j^n| + \sigma |\delta_x v_j^n|$$

$$\leq |v_j^n| + \sigma (|v_j^n| + |v_{j-1}^n|)$$

$$v_j^{n+1} = v_j^n - \sigma (v_j^n - v_{j-1}^n)$$

$$v_j^{n+1} = (1-\sigma) v_j^n + \sigma v_{j-1}^n$$

$$\rightarrow |v_j^{n+1}| \leq (1-\sigma) |v_j^n| + \sigma |v_{j-1}^n|$$

$$\text{define } Z^n = \max_j |v_j^n|$$

$$|v_j^{n+1}| \leq (1-\sigma) Z^n + \sigma Z^n = Z^n, \forall j$$

$$\text{For } j \text{ such that } |v_j^{n+1}| = Z^{n+1}$$

$$\rightarrow Z^{n+1} \leq Z^n \leq \dots \leq Z^0 \rightarrow \text{stability}$$

Lax Theorem (Two versions)

9/18/03

Lax Equivalence Theorem

A consistent, two-level difference scheme for a well posed linear IVP is convergent if and only if it is stable.

Lax Theorem

If a two-level difference scheme having the form: $v_h^{n+1} = Q_h v_h^n + \Delta t G_h^n$, $n \Delta t = t \leq t_{\text{final}}$, is accurate of order (p, q) in the norm $\|\cdot\|_h$ to a well-posed IVP and is stable with respect to $\|\cdot\|_h$, then it is convergent of order (p, q) with respect to $\|\cdot\|_h$. (Q is a linear operator)

consistency + stability = convergence

why does this work?

$$\text{let } w_h^n = v_h^n - u_h^n = \text{error}$$

$$\rightarrow w_h^{n+1} = Q_h w_h^n + \Delta t G_h^n - u_h^{n+1}$$

$$w_h^{n+1} = Q_h (w_h^n + u_h^n) + \Delta t G_h^n - u_h^{n+1}$$

$$w_h^{n+1} = Q_h w_h^n + \underbrace{[\Delta t G_h^n + Q_h u_h^n - u_h^{n+1}]}_{\Delta t \tau_h^n, \text{ truncation error}}$$

$$w_h^{n+1} = Q_h w_h^n + \Delta t \tau_h^n$$

by recursion,

$$w_h^n = Q_h w_h^{n-1} + \Delta t \tau_h^{n-1} = Q_h (Q_h w_h^{n-2} + \Delta t \tau_h^{n-2}) + \Delta t \tau_h^{n-1}$$

$$= \dots = |Q_h|^n w_h^0 + \Delta t \sum_{k=1}^n Q_h^{n-k} \tau_h^{n-k}$$

$$\rightarrow \boxed{w_h^n = \Delta t \sum_{k=1}^n Q_h^{n-k} \tau_h^{n-k}}$$

Take $\|\cdot\|_h$ and use

$$\|z_n^n\|_h \leq A(\Delta t^p + \Delta x^p) \quad \forall n \quad (\text{From consistency})$$

~~From consistency~~

From stability, we have $\|Q_n^n\|_h \leq K e^{pt}$

$$\rightarrow \|w_n^n\| \leq \Delta t n K e^{pt} A(\Delta x^p + \Delta t^p)$$

$$\rightarrow \boxed{\|w_n^n\| \leq t_{\text{Final}} K e^{pt} A(\Delta x^p + \Delta t^p)}$$

Fourier Stability Analysis

let us consider an IVP For the heat equation

$$u_t = \nu u_{xx}, \quad |x| < \infty, \quad t > 0, \quad u(x, 0) = f(x)$$

$$\mathcal{F}[u(x, t)] = \hat{u}(\omega, t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-i\omega x} u(x, t) dx$$

$$\mathcal{F}[u_t(x, t)] = \hat{u}_t(\omega, t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-i\omega x} u_t(x, t) dx$$

$$\mathcal{F}[u_{xx}(x, t)] = (i\omega)^2 \hat{u}(\omega, t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-i\omega x} u_{xx}(x, t) dx$$

the transformed PDE becomes $\hat{u}_t = -\nu^2 \omega^2 \hat{u}$

with initial condition $\hat{f}(\omega) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} f(x) e^{-i\omega x} dx$

$$\text{inverse transform, } \mathcal{F}^{-1}[\hat{u}(\omega, t)] = u(x, t) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{i\omega x} \hat{u}(\omega, t) d\omega$$

solving for $u(x,t)$, we obtain

$$u(x,t) = \frac{1}{\sqrt{4\pi ut}} \int_{-\infty}^{\infty} f(x_0) e^{-\frac{(x-x_0)^2}{4ut}} dx_0$$

For constant coefficient linear difference equations a similar analysis involving discrete Fourier transforms is used to analyze stability.

Define the discrete Fourier transform:

$$F v_n = \hat{v}(\xi) = \frac{1}{\sqrt{2\pi}} \sum_{j=-\infty}^{\infty} e^{-ij\xi} v_j, \quad \xi \in [-\pi, \pi]$$

The inverse Fourier transform

$$F^{-1} \hat{v}(\xi) = v_j = \frac{1}{\sqrt{2\pi}} \int_{-\pi}^{\pi} e^{ij\xi} \hat{v}(\xi) d\xi$$

The norm is preserved under the transformations

$$\|v_n\|^2 = \sum_{j=-\infty}^{\infty} |v_j|^2 = \int_{-\pi}^{\pi} |\hat{v}(\xi)|^2 d\xi = \|\hat{v}(\xi)\|^2$$

$$\rightarrow \boxed{\|v_n\| = \|\hat{v}(\xi)\|}$$

$$\|v_n^k\| \leq k e^{pt} \|v_n^0\|$$

$$\|\hat{v}^k\| \leq k e^{pt} \|\hat{v}^0\|$$

Example:

Consider the stability of the difference scheme

$$v_i^{n+1} = (1-2r) v_i^n + r(v_{i-1}^n + v_{i+1}^n) \quad , \quad r = \frac{v \Delta t}{\Delta x^2}$$

calculate some DFTs:

$$F v_i^{n+1} = \frac{1}{\sqrt{2\pi}} \sum_{j=-\infty}^{\infty} e^{-ij\xi} v_j^{n+1} = \hat{v}^{n+1}(\xi)$$

$$F v_i^n = \frac{1}{\sqrt{2\pi}} \sum_{j=-\infty}^{\infty} e^{-ij\xi} v_j^n = \hat{v}^n(\xi)$$

$$\begin{aligned} F v_{i-1}^n &= \frac{1}{\sqrt{2\pi}} \sum_{j=-\infty}^{\infty} e^{-ij\xi} v_{j-1}^n = \frac{1}{\sqrt{2\pi}} \sum_{m=-\infty}^{\infty} e^{-i(m+1)\xi} v_m^n \\ &= \frac{e^{-i\xi}}{\sqrt{2\pi}} \sum_{j=-\infty}^{\infty} e^{-ij\xi} v_j^n = e^{-i\xi} \hat{v}^n(\xi) \end{aligned}$$

and similarly, $F v_{i+1}^n = e^{i\xi} \hat{v}^n(\xi)$

substitute into difference scheme

$$\hat{v}^{n+1}(\xi) = \left[(1-2r) + r(e^{-i\xi} + e^{i\xi}) \right] \hat{v}^n(\xi)$$

$$\hat{v}^{n+1}(\xi) = [1-2r + 2r \cos \xi] \hat{v}^n(\xi) = (1-2r(1-\cos \xi)) \hat{v}^n(\xi)$$

$\Rightarrow \hat{v}^{n+1}(\xi) = p(\xi) \hat{v}^n(\xi)$ where $p(\xi)$ is the "symbol" of the difference equation

$$\rightarrow p(\xi) = 1-2r(1-\cos \xi) = 1-4r \sin^2\left(\frac{\xi}{2}\right)$$

If we require $|p(\xi)| \leq 1$ for $\xi \in [-\pi, \pi]$ then we have stability (with $K=1$, $\beta=0$)

Note, for $\xi \in [-\pi, \pi] \rightarrow \sin^2\left(\frac{\xi}{2}\right) \in [0, 1] \rightarrow |p(\xi)| \leq 1 \Rightarrow 1-4r \geq -1$

$$\rightarrow \boxed{r \leq \frac{1}{2} \text{ , as before}}$$

we had $\hat{v}^{n+1}(\xi) = p(\xi) \hat{v}^n(\xi)$, $\hat{v}^n(\xi) = [p(\xi)]^n v^0(\xi)$

$$\rightarrow \|\hat{v}^n(\xi)\| = \|p(\xi)^n v^0(\xi)\| \leq \|p(\xi)^n\| \|v^0(\xi)\|$$

but $\|p(\xi)^n\| \leq K e^{\beta t}$, $t = n \Delta t$

For stability, clearly if $|p(\xi)| \leq 1$, the inequality for stability is satisfied.

Example : convection-diffusion

Analyze the stability of $v_j^{n+1} = v_j^n - \frac{\sigma}{2} \delta_{0n} v_j^n + r \delta_x^2 v_j^n$

which approximates $u_t + cu_x = \nu u_{xx}$, $\sigma = \frac{c \Delta t}{\Delta x}$, $r = \frac{\nu \Delta t}{\Delta x^2}$

Take DFT :

$$\hat{v}^{n+1} = \hat{v}^n - \frac{\sigma}{2} (e^{i\xi} - e^{-i\xi}) \hat{v}^n + r (e^{i\xi} - 2 + e^{-i\xi}) \hat{v}^n$$

$$\hat{v}^{n+1} = \underbrace{\left[1 + i\sigma \sin(\xi) + 2r(\cos \xi - 1) \right]}_{p(\xi)} \hat{v}^n$$

Notice $p(\xi)$ is complex

$$|p(\xi)|^2 = (1 - 2r(1 - \cos \xi))^2 + \sigma^2 \sin^2 \xi$$

Let $z = \cos \xi \rightarrow |p(\xi)|^2 = (1 - 2r(1 - z))^2 + \sigma^2(1 - z^2)$

$$|p(\xi)| \leq 1 \rightarrow \cancel{1 - 4r(1 - z) + 4r^2(1 - z)^2 + \sigma^2(1 - z^2)} \leq \cancel{1}$$

$$\rightarrow \underbrace{-4r + 4r^2(1 - z) + \sigma^2(1 + z)}_{g(z)} \leq 0$$

since $g(z)$ is linear, need only check endpoints

$\therefore g(1) = -4r + 2\sigma^2 \leq 0 \rightarrow \sigma^2 \leq 2r$

$g(-1) = -4r + 8r^2 \leq 0 \rightarrow 2r \leq 1$

$\rightarrow \sigma^2 \leq 2r \leq 1$ for stability

Lax Equivalence Theorem

- A consistent, two-level difference scheme for a well-posed linear IVP is convergent if and only if it is stable.

Lax Theorem (special case of above theorem)

If a two-level difference scheme having the form $v_h^{n+1} = Q_h v_h^n + \Delta t G_h^n$, $n\Delta t \leq t_{\text{final}}$ is accurate of order (p, q) in the grid norm $\|\cdot\|_h$ to a well posed linear IVP and is stable wrt $\|\cdot\|_h$, then it is convergent of order (p, q) .

$$\Rightarrow \boxed{\text{consistency} + \text{stability} \Rightarrow \text{convergent}}$$

Fourier Stability Analysis

discrete Fourier transform (DFT):

$$\hat{v}(\xi) = \frac{1}{\sqrt{2\pi}} \sum_{j=-\infty}^{\infty} e^{-ij\xi} v_j, \quad \xi \in [-\pi, \pi]$$

$$\boxed{\hat{v}(\xi) = F v_j} \quad - \text{discrete FT of grid function } v_j$$

$$v_j = \frac{1}{\sqrt{2\pi}} \int_{-\pi}^{\pi} \hat{v}(\xi) e^{ij\xi} d\xi, \quad \text{for all } j$$

$$\boxed{v_j = F^{-1} \hat{v}}$$

where we assume v_j decays sufficiently quickly in order for the summation to exist

Parseval's Equality:

$$\|v_h\|^2 = \sum_{j=-\infty}^{\infty} |v_j|^2 = \int_{-\pi}^{\pi} |\hat{v}(\xi)|^2 d\xi = \|\hat{v}(\xi)\|^2$$

→ if we can show regularity of \hat{v} with respect to $\|\cdot\|^2 \Rightarrow$ regularity of v_h in the analogous norm.

Example An implicit scheme:

$$v_j^{n+1} = v_j^n + r \delta_x^2 v_j^{n+1}, \quad r > 0$$

this is a discretization of heat equation with $r = \frac{\Delta t}{\Delta x^2}$.

Take a DFT of the equation:

$$F v_j^{n+1} = F v_j^n + r F \delta_x^2 v_j^{n+1}$$

$$\hat{v}^{n+1} = \hat{v}^n + r F \left(v_{j-1}^{n+1} - 2v_j^{n+1} + v_{j+1}^{n+1} \right)$$

$$\hat{v}^{n+1} = \hat{v}^n + r F v_{j-1}^{n+1} - 2r F v_j^{n+1} + r F v_{j+1}^{n+1}$$

$$\hat{v}^{n+1} = \hat{v}^n + r e^{-i\xi} \hat{v}^{n+1} - 2r \hat{v}^{n+1} + r e^{i\xi} \hat{v}^{n+1}$$

$$\hat{v}^{n+1} = \hat{v}^n + r \left(e^{i\xi} + e^{-i\xi} - 2 \right) \hat{v}^{n+1}$$

$$\hat{v}^{n+1} = \hat{v}^n + r \left(2 \cos \xi - 2 \right) \hat{v}^{n+1}$$

$$\left[1 + 2r(1 - \cos \xi) \right] \hat{v}^{n+1} = \hat{v}^n$$

$$\rightarrow \hat{v}^{n+1} = \frac{\hat{v}^n}{1 + 2r(1 - \cos \xi)}$$

$\rho(\xi) = \frac{1}{1 + 2r(1 - \cos \xi)}$, the "symbol" of the difference scheme, describes growth or decay of discrete grid function in transformed space

→ notice, $1 \leq 1 + 2r(1 - \cos \xi) \leq 1 + 4r$

→ $|p(\xi)| \leq 1$

→ $\|\hat{v}^n\| \leq \|\hat{v}^0\| \rightarrow$ stability, independent of r

unconditionally stable

Implicit methods often have favorable stability properties as compared to explicit schemes.

Example:

Exponential growth or decay. Consider the heat equation with a linear source term.

$$u_t = \gamma u_{xx} + bu$$

(growth or decay depending on sign of b)

Approximate using:

$$v_j^{n+1} = v_j^n + r \delta_x^2 v_j^{n+1} + b \Delta t v_j^n$$

Analyze stability: take DFT

$$\hat{v}^{n+1} = \hat{v}^n + 2r(1 - \cos \xi) \hat{v}^{n+1} + b \Delta t \hat{v}^n$$

$$\rightarrow \hat{v}^{n+1} = \frac{1 + b \Delta t}{\underbrace{1 + 2r(1 - \cos \xi)}_{p(\xi)}} \hat{v}^n$$

$$|p(\xi)| \leq |1 + b \Delta t|$$

there are two cases:

if $b > 0$, (growth) $\rightarrow |p| \leq 1 + b \Delta t \rightarrow$ stable
(because stability allows some growth)

if $b < 0$, (decay) $\rightarrow |p| \leq |1 - (-b) \Delta t|$

$$\rightarrow -b \Delta t \leq 2 \rightarrow \Delta t \leq \frac{2}{-b} \quad (b < 0)$$

$$\begin{aligned} -2x &< 1 \\ x &< -\frac{1}{2} \end{aligned}$$

von Neumann condition

Using a Fourier analysis, you obtain

$$\hat{v}^{n+1} = \rho(\xi) \hat{v}^n$$

The von Neumann ^{stability} condition is that

$$|\rho(\xi)| \leq 1 + C \Delta t, \text{ for } \xi \in [-\pi, \pi]$$

where C is a constant independent of Δx and Δt and for $\Delta x, \Delta t$ sufficiently small.

Tighter condition: $|\rho(\xi)| \leq 1$ (for no growth)

Fourier Mode Analysis

$$\text{inverse DFT: } v_j^n = \frac{1}{\sqrt{2\pi}} \int_{-\pi}^{\pi} \hat{v}^n(\xi) e^{i\xi j} d\xi \quad \forall j$$

The integral represents a sum of Fourier modes of the form

$$\underbrace{\left(\frac{1}{\sqrt{2\pi}} \hat{v}^n(\xi) \right)}_{\text{temporal component}} \underbrace{\left(e^{i\xi j} \right)}_{\text{spatial component}}$$

$$\text{spatial component: } e^{i\xi j} = e^{ikx_j}, \quad x_j = j \Delta x, \quad k = \xi / \Delta x$$

$$= \cos(kx_j) + i \sin(kx_j)$$

These represent spatial oscillations where k is a spatial frequency = wave number.

Recall, $\xi \in [-\pi, \pi] \rightarrow k \in \left[-\frac{\pi}{\Delta x}, \frac{\pi}{\Delta x} \right] \rightarrow$ Finite range of wave numbers,
(higher values of k are aliased to lower ones)

the temporal component: $\frac{1}{\sqrt{2\pi}} \hat{v}^n(\xi) = \frac{1}{\sqrt{2\pi}} \left[a(k\Delta x) \right]^n \hat{v}^0(\xi)$ power, not index 25

"growth" Factor

or
"amplitude" Factor

tells you how the mode
with wave number k grows
or decays $\frac{1}{2}$

The idea: in a Fourier mode analysis, consider solutions of
the difference equation of the form

$$v_j^n = a^n e^{ikx_j}$$

separable solution
of difference eqn.

$x_j = j\Delta x$, k = wave number
 a = amplitude factor (complex)

von Neumann condition in this analysis becomes

$$|a| \leq 1 + c\Delta t, \text{ for } |k\Delta x| \leq \pi$$

or more stringently, $|a| \leq 1$, for $|k\Delta x| \leq \pi$

(the dimensions of k is $1/\text{length}$, such that $k\Delta x$ is non-dimensional)

$k\Delta x$ - "gridnumber"

Example: Convection-Diffusion problem

$$v_j^{n+1} = v_j^n - \frac{\tau}{2} \delta_{0x} v_j^n + r \delta_x^2 v_j^n$$

explicit scheme for $u_t + cu_x = \nu u_{xx}$.

$$\tau = \frac{\Delta t c}{\Delta x} \text{ (positive or negative)}, \quad r = \frac{\Delta t \nu}{\Delta x^2} > 0$$

Notice: if diffusive term is not present, this scheme
is unstable

$$v_j^{n+1} = v_j^n - \frac{\sigma}{2} \delta_{0x} v_j^n + r \delta_x^2 v_j^n$$

Mode analysis , $v_j^n = a^n e^{ikx_j}$

Example,
continued

$$\rightarrow v_j^{n+1} = a^{n+1} e^{ikx_j} = a v_j^n$$

$$\begin{aligned} \rightarrow \delta_{0x} v_j^n &= v_{j+1}^n - v_{j-1}^n = a^n e^{ikx_{j+1}} - a^n e^{ikx_{j-1}} \\ &= a^n e^{ik(x_j + \Delta x)} - a^n e^{ik(x_j - \Delta x)} \\ &= \cancel{e^{ikx_j}} e^{ik\Delta x} v_j^n - e^{-ik\Delta x} v_j^n \\ &= 2i \sin(k\Delta x) v_j^n \end{aligned}$$

$$\begin{aligned} \delta_x^2 v_j^n &= v_{j+1}^n - 2v_j^n + v_{j-1}^n = e^{ik\Delta x} v_j^n - 2v_j^n + e^{-ik\Delta x} v_j^n \\ &= -2(1 - \cos(k\Delta x)) v_j^n \end{aligned}$$

$$\rightarrow a v_j^n = v_j^n - \frac{\sigma}{2} (2i \sin(k\Delta x) v_j^n) - 2r(1 - \cos(k\Delta x)) v_j^n$$

$$\boxed{a = 1 - i\sigma \sin(k\Delta x) - 2r(1 - \cos(k\Delta x))}$$

For no growth we would like $|a| \leq 1 \quad \forall |k\Delta x| \leq \pi$

$$|a|^2 = (\operatorname{Re}(a))^2 + (\operatorname{Im}(a))^2 \leq 1$$

trick is to let $\nu = \cos(k\Delta x)$ and $\sin(k\Delta x) = \sqrt{1 - \nu^2}$ and so on...

Stability Analysis For an Initial-BVP problem

Example: $u_t = \nu u_{xx}$, $0 \leq x \leq 1$, $t \geq 0$

$$u(x, 0) = f(x)$$

$$u(0, t) = a(t) \quad \} \text{ Dirichlet BCs}$$

$$u(1, t) = b(t)$$

discrete approximation:

$$\frac{1}{\Delta t} \delta_{+t} v_j^n = \frac{\nu}{\Delta x^2} \delta_x^2 v_j^n \quad - \quad 1 \leq j \leq N-1, \quad n \geq 0$$

$$x_j = j \Delta x, \quad \Delta x = \frac{1}{N}, \quad t_n = n \Delta t$$

$$v_j^0 = f(x_j), \quad 0 \leq j \leq N$$

$$v_0^n = a(t_n), \quad v_N^n = b(t_n), \quad n > 0$$

Midterm Exam: Tentative Date, Tues 10/14/03

What is the stability of the approximation?

Consider a perturbation to v_j^n , it would solve the discrete problem with homogeneous boundary conditions.

Call the perturbed solution \tilde{v}_j^n . It solves

$$\tilde{v}_j^{n+1} = \tilde{v}_j^n + r \delta_x^2 \tilde{v}_j^n, \quad r = \frac{\nu \Delta t}{\Delta x^2}$$

with \tilde{v}_j^0 given and $\tilde{v}_0^n = \tilde{v}_N^n = 0$.

Consider the behavior of this perturbation:

$$\tilde{V}_1^{n+1} = \tilde{V}_1^n + r (\cancel{\tilde{V}_0^n} - 2\tilde{V}_1^n + \tilde{V}_2^n) \quad \nearrow 0, \text{ due to BC}$$

$$\tilde{V}_1^{n+1} = (1-2r)\tilde{V}_1^n + r\tilde{V}_2^n$$

$$\tilde{V}_2^{n+1} = (1-2r)\tilde{V}_2^n + r\tilde{V}_1^n + r\tilde{V}_3^n$$

$$\vdots$$

$$\tilde{V}_{N-1}^{n+1} = (1-2r)\tilde{V}_{N-1}^n + r\tilde{V}_{N-2}^n + r\cancel{\tilde{V}_N^n}$$

this forms a linear system for $\underline{\tilde{V}}^{n+1} =$

$$\begin{bmatrix} \tilde{V}_1^{n+1} \\ \tilde{V}_2^{n+1} \\ \vdots \\ \tilde{V}_{N-1}^{n+1} \end{bmatrix}$$

$$\Rightarrow \underline{\tilde{V}}^{n+1} = \underline{A} \underline{\tilde{V}}^n$$

where

$$A = \begin{bmatrix} 1-2r & r & & & \\ r & 1-2r & r & & \\ & r & 1-2r & \ddots & \\ & & \ddots & \ddots & \\ & & & r & 1-2r & r \\ & & & & r & 1-2r \end{bmatrix}$$

TRIDIAGONAL
MATRIX

\Rightarrow repeated application shows that

$$\underset{\text{index}}{\underline{\tilde{V}}}^n = \underset{\text{power}}{\underline{A}}^n \underset{\text{index}}{\underline{\tilde{V}}}^0$$

Therefore For the perturbation to be "well behaved,"
we want

$$\|\underline{A}^n\| \leq 1$$

The task is to bound powers of \underline{A} , the growth matrix, independent of N (spatial grid) and Δt . May be difficult to do.

Notice that \underline{A} is symmetric, for our example.

\Rightarrow a orthogonal matrix \underline{R} exists such that $\underline{R}^T \underline{A} \underline{R} = \underline{\Lambda}$ where $\underline{\Lambda}$ is a matrix whose diagonal entries are real and are the eigenvalues of \underline{A}

$$\underline{\Lambda} = \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_n \end{bmatrix}, \quad \underline{R} = \begin{bmatrix} | & | & & | \\ \underline{r}_1 & \underline{r}_2 & \dots & \underline{r}_n \\ | & | & & | \end{bmatrix}$$

λ_j - eigenvalues of \underline{A} , \underline{r}_j - right eigenvectors of \underline{A}

For this case, the 2-norm of \underline{A} is equal to:

$$\|\underline{A}\| = \|\underline{R} \underline{\Lambda} \underline{R}^T\| = \|\underline{\Lambda}\| = \max_j |\lambda_j| \leftarrow \text{spectral radius}$$

The problem of stability becomes analyzing the eigenvalues of the growth matrix. This is true for all symmetric matrices \underline{A} .

Let us calculate eigenvalues for \underline{A} : $\underline{A} \underline{z} = \lambda \underline{z}$

set ~~$\underline{A} = \underline{I} + r \underline{T}$~~ $\underline{A} = \underline{I} + r \underline{T}$, where \underline{T}

$$\underline{T} = \begin{bmatrix} -2 & 1 & & \\ 1 & -2 & 1 & \\ & 1 & -2 & 1 \\ & & 1 & -2 \end{bmatrix} \rightarrow (\underline{I} + r \underline{T}) \underline{z} = \lambda \underline{z}$$

$$\rightarrow r \underline{T} \underline{z} = (\lambda - 1) \underline{z} \rightarrow \underline{T} \underline{z} = \frac{\lambda - 1}{r} \underline{z}$$

$$\rightarrow \underline{T} \underline{z} = \left(\frac{\lambda - 1}{r} \right) \underline{z} \rightarrow \underline{T} \underline{z} = \mu \underline{z}, \quad \mu = \frac{\lambda - 1}{r}$$

eigenvalues of \underline{T} :

$$\begin{bmatrix} -2 & 1 & & \\ 1 & -2 & 1 & \\ & 1 & -2 & 1 \\ & & \ddots & \ddots \end{bmatrix} \begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ \vdots \end{pmatrix} = \lambda \begin{pmatrix} z_1 \\ z_2 \\ \vdots \end{pmatrix}$$

the k^{th} row of this system is:

$$z_{k-1} - 2z_k + z_{k+1} = \lambda z_k \quad \text{for } k=1 \text{ to } N-1$$

assign $z_0 = z_N = 0$

$$\rightarrow z_{k-1} - (2+\lambda)z_k + z_{k+1} = 0$$

Assign $2+\lambda = 2\cos\theta$

$$\rightarrow z_{k-1} - 2\cos\theta z_k + z_{k+1} = 0, \quad z_0 = z_N = 0 \rightarrow$$

Difference
Equation
with
constant
coeffs

Set $z_k = \xi^k$, $\xi = \text{constant}$, possibly complex
and ξ^k is the k^{th} power of ξ , not subscript

$$\rightarrow \xi^{k-1} - 2\cos\theta \xi^k + \xi^{k+1} = 0 \rightarrow$$

$$\xi^{k-1} (1 - 2\cos\theta \xi + \xi^2) = 0$$

$$\xi = \frac{2\cos\theta \pm \sqrt{4\cos^2\theta - 4}}{2} = \frac{\cos\theta \pm \sqrt{\cos^2\theta - 1}}{1} = \cos\theta \pm i\sin\theta$$

$$\rightarrow \xi = e^{\pm i\theta}$$

$$\rightarrow \text{general solution: } z_k = c_1 e^{i\theta k} + c_2 e^{-i\theta k}$$

$$\text{set } z_0 = 0 \rightarrow c_1 = -c_2 \rightarrow z_k = c_1 (e^{ik\theta} - e^{-ik\theta})$$

$$\rightarrow z_k = 2ic_1 \sin(k\theta)$$

$$\text{set } z_N = 0 \rightarrow 2ic_1 \sin(N\theta) = 0 \rightarrow N\theta = p\pi$$

$$p = 1, 2, \dots, N-1$$

$$\rightarrow \theta = \frac{p\pi}{N}$$

work backwards:

$$2 + \nu = 2 \cos \Theta \rightarrow \nu = -2(1 - \cos \Theta)$$

$$\rightarrow \nu = -2(1 - \cos \frac{p\pi}{N})$$

$$\text{and now } \lambda = 1 + r\nu = 1 - 2r(1 - \cos \frac{p\pi}{N}), p = 1, \dots, N-1$$

$$\rightarrow \boxed{\lambda = 1 - 2r(1 - \cos \frac{p\pi}{N}), p = 1, \dots, N-1}$$

What is restriction on r so that $|\lambda| \leq 1$???

$$0 < 1 - \cos(\frac{p\pi}{N}) < 2, \text{ notice strict inequalities because } p = 1, \dots, N-1$$

$$\rightarrow 1 - 4r \geq -1 \rightarrow \boxed{r \leq \frac{1}{2}}$$

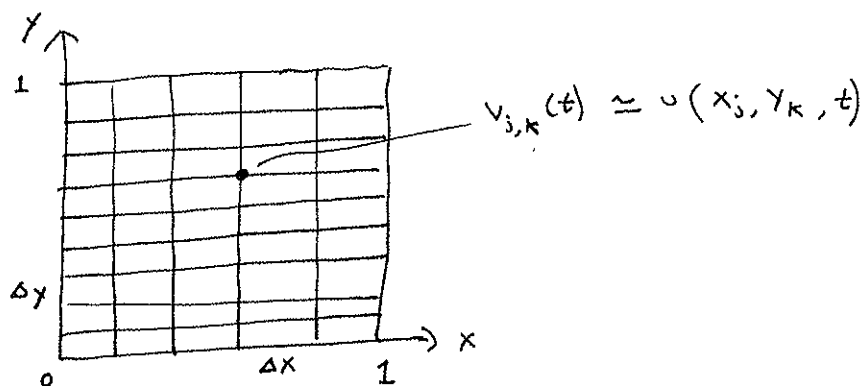
Finite Difference Methods for Parabolic Methods

$$u_t = v(u_{xx} + u_{yy}), \quad 0 < x < 1, \quad 0 < y < 1$$

$$\text{ICs: } u(x, y, 0) = f(x, y), \quad t > 0$$

$$\text{BCs: } u(x, y, t) = g(x, y, t) \text{ for } (x, y) \text{ on the boundary}$$

Approximate solution, grid (x_j, y_k) , $j = 0, \dots, N$, $k = 0, \dots, M$
 $x_j = j\Delta x$, $\Delta x = \frac{1}{N}$, $y_k = k\Delta y$, $\Delta y = \frac{1}{M}$



Discrete equations: (method of lines)

$$v'_{j,k}(t) = \nu \left[\frac{1}{\Delta x^2} \delta_x^2 + \frac{1}{\Delta y^2} \delta_y^2 \right] v_{j,k}(t), \quad \begin{matrix} j=1, \dots, N-1 \\ k=1, \dots, M-1 \end{matrix} \quad \left. \vphantom{\begin{matrix} j=1, \dots, N-1 \\ k=1, \dots, M-1 \end{matrix}} \right\} \begin{matrix} \text{SYSTEM} \\ \text{OF} \\ \text{ODES} \end{matrix}$$

$$v_{j,k}(t) = g(x_j, y_k, t) \quad \text{For } x_j, y_k \text{ on boundary}$$

$$\text{Initial conditions, } v_{j,k}(0) = f(x_j, y_k)$$

Now integrate in time. The simplest method of numerical integration is Forward Euler:

$$\frac{v_{j,k}(t+\Delta t) - v_{j,k}(t)}{\Delta t} = \nu \left[\frac{1}{\Delta x^2} \delta_x^2 + \frac{1}{\Delta y^2} \delta_y^2 \right] v_{j,k}(t)$$

$$\text{Set } t_n = n\Delta t, \quad v_{j,k}(t_n) = v_{j,k}^n$$

$$\rightarrow v_{j,k}^{n+1} = v_{j,k}^n + \left[r_x \delta_x^2 v_{j,k}^n + r_y \delta_y^2 v_{j,k}^n \right]$$
$$r_x = \frac{\nu \Delta t}{\Delta x^2}, \quad r_y = \frac{\nu \Delta t}{\Delta y^2}$$

explicit method
advantage is that
the implementation is
straight forward

order of accuracy: $O(\Delta t, \Delta x^2, \Delta y^2)$

For stability, conduct a mode analysis:

$$v_{j,k}^n = a^n e^{i\alpha x_j} e^{i\beta y_k}$$

where α, β are wave numbers and a - growth factor

substitute into numerical scheme:

$$a = 1 + 2r_x(\cos \alpha \Delta x - 1) + 2r_y(\cos \beta \Delta y - 1)$$

$$a = 1 - 2 \left[r_x(1 - \cos \alpha \Delta x) + r_y(1 - \cos \beta \Delta y) \right]$$

we want to restrict r_x and r_y so that
 $|a| \leq 1$ For $|\alpha \Delta x| \leq \pi$ and $|\beta \Delta y| \leq \pi$

when $\alpha \Delta x = \alpha \Delta y = 0$, $a = 1$,

we need $1 - 4(r_x + r_y) \geq -1 \rightarrow \frac{1}{2} \geq r_x + r_y$

$$\rightarrow \frac{\Delta t \nu}{\Delta x^2} + \frac{\Delta t \nu}{\Delta y^2} \leq \frac{1}{2} \rightarrow \Delta t \nu \left(\frac{\Delta x^2 + \Delta y^2}{\Delta x^2 \Delta y^2} \right) \leq \frac{1}{2}$$

$$\rightarrow \boxed{\Delta t \nu \leq \frac{1}{2} \frac{\Delta x^2 \Delta y^2}{\Delta x^2 + \Delta y^2}} \text{ stability constraint}$$

suppose $\Delta x = \Delta y = h$, then $\Delta t \nu \leq \frac{h^2}{4}$

this is a severe restriction on Δt

suggests the consideration of an implicit method

eg Backward Euler, Trapezoidal Rule, etc.....

9/29/03

Crank - Nicholson Method

$$u_t = \nu(u_{xx} + u_{yy}) \quad 0 \leq x \leq 1, \quad 0 \leq y \leq 1, \quad t \geq 0$$

$$u(x, y, 0) = f(x, y), \quad u(x, y, t) = g(x, y, t) \text{ on boundary}$$

$$v'_{j,k}(t) = \nu \left[\frac{1}{\Delta x^2} \delta_x^2 + \frac{1}{\Delta y^2} \delta_y^2 \right] v_{j,k}(t) \rightarrow O(\Delta t, \Delta x^2, \Delta y^2) \left. \vphantom{\begin{matrix} v'_{j,k}(t) \\ \delta_x^2 \\ \delta_y^2 \end{matrix}} \right\} \begin{array}{l} \text{For} \\ \text{Forward} \\ \text{Euler} \end{array}$$

$$\text{For stability } \Delta t \leq \frac{1}{2\nu} \frac{\Delta x^2 \Delta y^2}{(\Delta x^2 + \Delta y^2)}$$

If we integrate using trapezoidal rule, obtain
Crank - Nicholson method

$$\frac{v_{j,k}^{n+1} - v_{j,k}^n}{\Delta t} = \nu \left[\frac{1}{\Delta x^2} \delta_x^2 + \frac{1}{\Delta y^2} \delta_y^2 \right] \left(\frac{v_{j,k}^{n+1} + v_{j,k}^n}{2} \right)$$

$$\frac{V_{i,k}^{n+1} - V_{i,k}^n}{\Delta t} = \nu \left[\frac{1}{\Delta x^2} \delta_x^2 + \frac{1}{\Delta y^2} \delta_y^2 \right] \left(\frac{V_{i,k}^{n+1} + V_{i,k}^n}{2} \right)$$

$$\rightarrow V_{i,k}^{n+1} \left(1 - \frac{r_x}{2} \delta_x^2 - \frac{r_y}{2} \delta_y^2 \right) = \left(1 + \frac{r_x}{2} \delta_x^2 + \frac{r_y}{2} \delta_y^2 \right) V_{i,k}^n$$

$$\text{where } r_x = \frac{\nu \Delta t}{\Delta x^2}, \quad r_y = \frac{\nu \Delta t}{\Delta y^2}$$

this is a linear system, $\underline{A} \underline{V}^{n+1} = \underline{C}^n$

$$\text{where } \underline{V}^n = \left[V_{1,1}^n \cdots V_{N-1,1}^n \mid V_{1,2}^n \cdots V_{N-1,2}^n \mid \cdots \mid V_{1,M-1}^n \cdots V_{N-1,M-1}^n \right]^T$$

\underline{V}^n is an $(M-1)(N-1)$ vector

the matrix \underline{A} is a block tri-diagonal matrix
where each "sub" matrix is $(N-1) \times (M-1)$

$$\underline{A} = \begin{bmatrix} \underline{A}_{11} & \underline{A}_{12} & & \\ \underline{A}_{21} & \underline{A}_{22} & \underline{A}_{23} & \\ & \ddots & \ddots & \ddots \\ & & \underline{A}_{N-2,M-2} & \underline{A}_{N-1,M-1} \end{bmatrix}$$

where \underline{A}_{kk} =

$$\begin{bmatrix} 1+r_x+r_y & -\frac{r_x}{2} & & \\ -\frac{r_x}{2} & 1+r_x+r_y & -\frac{r_x}{2} & \\ & -\frac{r_x}{2} & 1+r_x+r_y & -\frac{r_x}{2} \\ & & \ddots & \ddots \\ & & & -\frac{r_x}{2} & 1+r_x+r_y \end{bmatrix}$$

matrix on the diagonal

the matrices on sub-diagonal and super-diagonal
are $\underline{A}_{k-1,k} = \underline{A}_{k,k-1} = -\frac{r_y}{2} \underline{I}$

$$\underline{C}^n = \begin{bmatrix} C_1^n \\ C_2^n \\ \vdots \\ C_{N-1}^n \end{bmatrix}, \quad C_k^n = \left(1 + \frac{r_x}{2} \delta_x^2 + \frac{r_y}{2} \delta_y^2 \right) V_{i,k}^n$$

the system is $(M-1)(M-1) \times (N-1)(N-1)$

boundary terms

when $k=1$:

$$\frac{r_y}{2} \begin{bmatrix} g(x_1, 0, (n+1)\Delta t) \\ \vdots \\ g(x_{N-1}, 0, (n+1)\Delta t) \end{bmatrix} + \frac{r_x}{2} \begin{bmatrix} g(0, y_N, (n+1)\Delta t) \\ \vdots \\ g(0, y_2, (n+1)\Delta t) \end{bmatrix}$$

when $k=2$:

$$\frac{r_y}{2} \begin{bmatrix} g(x_1, 0, (n+1)\Delta t) \\ \vdots \\ g(x_{N-1}, 0, (n+1)\Delta t) \end{bmatrix} + \frac{r_x}{2} \begin{bmatrix} g(0, y_N, (n+1)\Delta t) \\ \vdots \\ g(0, y_2, (n+1)\Delta t) \end{bmatrix}$$

and so on for $k=3, 4, \dots$

Algorithm

- 1) specification $\Delta t, N, M, r_x, r_y$
- 2) set $v_{i,k}^0 = F(x_i, y_k)$
- 3) build \underline{A} (Factor it)
- 4) for $n=0, 1, \dots, n_{\text{final}}$
- 5) build \underline{c}^n using $v_{i,k}^n + \text{BCs}$
- 6) solve $\underline{A} \underline{v}^{n+1} = \underline{c}^n$
- 7) construct $v_{i,k}^{n+1}$ from $\underline{v}^{n+1} + \text{BCs}$
- 8) output solution

Order of Accuracy

$$\tau_{i,k}^n = v \left[\frac{1}{\Delta x^2} \delta_x^2 + \frac{1}{\Delta y^2} \delta_y^2 \right] \left(\frac{v_{i,k}^n + v_{i,k}^{n+1}}{2} \right) - \left(\frac{v_{i,k}^{n+1} - v_{i,k}^n}{\Delta t} \right)$$

the order is $O(\Delta t^2, \Delta x^2, \Delta y^2)$

Stability

$$\text{set } v_{i,k}^n = a^n e^{i(\alpha x_i + \beta y_k)} \quad \text{wave numbers}$$

$$\begin{aligned} & \left[1 + r_x(1 - \cos \alpha \Delta x) + r_y(1 - \cos \beta \Delta y) \right] a \\ &= \left[1 - r_x(1 - \cos \alpha \Delta x) - r_y(1 - \cos \beta \Delta y) \right] \end{aligned}$$

$$a = \frac{1 - r_x(1 - z_1) - r_y(1 - z_2)}{1 + r_x(1 - z_1) + r_y(1 - z_2)} = \frac{N}{D}, \quad \begin{aligned} z_1 &= \cos \alpha \Delta x \\ z_2 &= \cos \beta \Delta y \end{aligned}$$

$$|a| < 1 \Rightarrow |N| \leq |D| \Rightarrow |N^2| \leq |D^2|$$

$$\begin{aligned} & \cancel{1 + r_x^2(1 - z_1)^2 + r_y^2(1 - z_2)^2} - 2r_x(1 - z_1) - 2r_y(1 - z_2) + \cancel{2r_x r_y(1 - z_1)(1 - z_2)} \leq \\ & \cancel{1 + r_x^2(1 - z_1)^2 + r_y^2(1 - z_2)^2} + 2r_x(1 - z_1) + 2r_y(1 - z_2) + \cancel{2r_x r_y(1 - z_1)(1 - z_2)} \end{aligned}$$

$$\rightarrow -r_x(1 - z_1) - r_y(1 - z_2) \leq r_x(1 - z_1) + r_y(1 - z_2)$$

$$\rightarrow 0 \leq r_x(1 - z_1) + r_y(1 - z_2)$$

$$\left. \begin{aligned} 1 - z_1 &\geq 0 \\ 1 - z_2 &\geq 0 \\ r_x &\geq 0 \\ r_y &\geq 0 \end{aligned} \right\} \text{unconditional stability}$$

↓

accuracy can now dictate
the time step

Computational Cost

the main cost is in solving the linear system,
but it is sparse and banded

If we want a solution at an $O(1)$ time,
 $N \sim M \rightarrow O(N)$ time steps are required

total cost is $N \cdot (\text{cost per step})$

cost per step is based on the matrix
solve of $\underline{A} \underline{v}^{n+1} = \underline{c}^n$

assume, $N \sim M$:

$$N^2 \left\{ \begin{array}{c} \left[\begin{array}{c} \swarrow \quad \nwarrow \\ \leftarrow N \rightarrow \\ \searrow \quad \swarrow \end{array} \right] \\ \underbrace{\hspace{1.5cm}}_{N^2} \end{array} \right.$$

direct solve, such as
block tri-diagonal
would require
 $O(\text{dim of system} \times \text{bandwidth}^2)$
 $= O(N^4)$
per step

\rightarrow cost = $O(N^5)$ total

Iterative solvers: conjugate gradients, SOR

$O(\underbrace{\text{ops per sweep}}_{N^2} \times \underbrace{\# \text{ of sweeps}}_N) = O(N^3) \Rightarrow O(N^4) \text{ total}$

some Fancy methods may achieve, $O(N^2 \log N)$
 $\Rightarrow O(N^3 \log N)$ total

Alternating Direction Implicit (ADI)

commonly used for parabolic equations in 2, 3 dimensions

consider the heat equation

$$u_t = \nu(u_{xx} + u_{yy}), \quad 0 \leq x \leq 1, \quad 0 \leq y \leq 1, \quad t \geq 0$$

$$u(x, y, 0) = f(x, y)$$

$$u(x, y, t) = g(x, y, t) \text{ on boundary}$$

the ADI method splits the implicit part into
2-half steps

The ADI scheme of Peaseman & Rachford

$$\frac{v_{j,k}^{n+1/2} - v_{j,k}^n}{\Delta t/2} = \nu \left[\frac{1}{\Delta x^2} \delta_x^2 v_{j,k}^{n+1/2} + \frac{1}{\Delta y^2} \delta_y^2 v_{j,k}^n \right]$$

$$\frac{v_{j,k}^{n+1} - v_{j,k}^{n+1/2}}{\Delta t/2} = \nu \left[\frac{1}{\Delta x^2} \delta_x^2 v_{j,k}^{n+1/2} + \frac{1}{\Delta y^2} \delta_y^2 v_{j,k}^{n+1} \right]$$

boundary condition, $v_{j,k}^n = g(x_j, y_k, n\Delta t)$, (x_j, y_k) on boundary

initial condition $v_{j,k}^0 = f(x_j, y_k)$

The advantage is that we only perform an implicit calculation in one dimension at a time, and each half step requires a tri-diagonal solve.

First half-step:

$$v_{j,k}^{n+1/2} - v_{j,k}^n = r_x \delta_x^2 v_{j,k}^{n+1/2} + r_y \delta_y^2 v_{j,k}^n, \quad r_x = \frac{\nu \Delta t}{2\Delta x^2}, \quad r_y = \frac{\nu \Delta t}{2\Delta y^2}$$

$$\rightarrow (1 - r_x \delta_x^2) v_{j,k}^{n+1/2} = (1 + r_y \delta_y^2) v_{j,k}^n$$

$$\underline{v}_k^n = \begin{bmatrix} v_{1,k}^{n+1/2} \\ \vdots \\ v_{N-1,k}^{n+1/2} \end{bmatrix}, \quad \underline{A}_n = \begin{bmatrix} 1+2r_x & r_x & & \\ -r_x & 1+2r_x & -r_x & \\ & \ddots & \ddots & \\ & & -r_x & 1+2r_x \end{bmatrix}$$

$$\underline{c}_k^n = \begin{bmatrix} (1 + r_y \delta_y^2) v_{1,k}^n + r_x g(x_0, y_k, t_{n+1/2}) \\ (1 + r_y \delta_y^2) v_{2,k}^n \\ \vdots \\ (1 + r_y \delta_y^2) v_{N-1,k}^n + r_x g(x_N, y_k, t_{n+1/2}) \end{bmatrix}$$

the First half-step becomes:

$$\underline{A}_k \underline{v}_k^n = \underline{c}_k^n, \quad k = 1, \dots, M-1$$

\hookrightarrow tridiagonal $\rightarrow O(NM)$ ops

Similarly, the second half step

$$\underline{A}_k \underline{v}_j^n = \underline{c}_j^n, \quad j=1, \dots,$$

$\uparrow O(NM)$, so each step requires $O(NM)$ operations

Consistency and Order of Accuracy & Stability

$$1^{st} \text{ step: } (1 - r_x \delta_x^2) v_{j,k}^{n+1/2} = (1 + r_y \delta_y^2) v_{j,k}^n$$

$$2^{nd} \text{ step: } (1 - r_y \delta_y^2) v_{j,k}^{n+1} = (1 + r_x \delta_x^2) v_{j,k}^{n+1/2}$$

$$\begin{aligned} \underline{\text{Note:}} \quad (1 - r_x \delta_x^2)(1 - r_y \delta_y^2) v_{j,k}^{n+1} &= (1 - r_x \delta_x^2)(1 + r_x \delta_x^2) v_{j,k}^{n+1/2} \\ &= (1 + r_x \delta_x^2)(1 + r_y \delta_y^2) v_{j,k}^n \end{aligned}$$

$$\rightarrow \Delta t \tau_{j,k}^n = (1 - r_x \delta_x^2)(1 - r_y \delta_y^2) u_{j,k}^{n+1} - (1 + r_x \delta_x^2)(1 + r_y \delta_y^2) u_{j,k}^n$$

$$\begin{aligned} \text{note, } r_x \delta_x^2 u_{j,k}^n &= r_x (\Delta x^2 u_{xx} + O(\Delta x^4)) \\ &= \frac{v \Delta t}{2} u_{xx} + O(\Delta x^2) \end{aligned}$$

$$r_y \delta_y^2 u_{j,k}^n = \frac{v \Delta t}{2} u_{yy} + O(\Delta y^2)$$

$$\begin{aligned} \rightarrow \Delta t \tau_{j,k}^n &= \left[u - \frac{v \Delta t}{2} (u_{xx} + u_{yy}) + \frac{v^2 \Delta t^2}{4} u_{xx} u_{yy} + O(\Delta x^2 \Delta t, \Delta y^2 \Delta t) \right]^{n+1} \\ &\quad - \left[u + \frac{v \Delta t}{2} (u_{xx} + u_{yy}) + \frac{v^2 \Delta t^2}{4} u_{xx} u_{yy} + O(\Delta x^2 \Delta t, \Delta y^2 \Delta t) \right]^n \end{aligned}$$

$$\begin{aligned} \rightarrow \Delta t \tau_{j,k}^n &= \Delta t u_t + \frac{\Delta t^2}{2} u_{tt} + O(\Delta t^3) - v \Delta t (u_{xx} + u_{yy}) \\ &\quad - \frac{v \Delta t^2}{2} (u_{xx} + u_{yy})_t + O(\Delta t^3) + \frac{v^2 \Delta t^3}{4} u_{xx} u_{yy} + O(\Delta x^2 \Delta t, \Delta y^2 \Delta t) \end{aligned}$$

$$\tau_{j,k}^n = O(\Delta t^2, \Delta x^2, \Delta y^2)$$

Methods For Heat Equation in Multi-space dimensions

for example, $u_t = \nu(u_{xx} + u_{yy})$, $0 < x < 1$, $0 < y < 1$

$$u(x, y, 0) = F(x, y)$$

u is given on boundary

the Method of Lines leads to a whole family of methods

need a grid: $x_j = j \Delta x$, $\Delta x = \frac{1}{N}$

$$y_k = k \Delta y, \quad \Delta y = \frac{1}{M}$$

$$v_{j,k}(t) = u(x_j, y_k, t)$$

let $v_{j,k}$ solve the ODEs $v'_{j,k}(t) = \nu \left(\frac{1}{\Delta x^2} \delta_x^2 + \frac{1}{\Delta y^2} \delta_y^2 \right) v_{j,k}$

for $1 \leq j \leq N-1$, $1 \leq k \leq M-1$, with initial condition $v_{j,k}(0) = F(x_j, y_k)$

and boundary conditions given at $j=0, N$, $k=0, M$

Time integration

a) Forward Euler: $\tau_{j,k}^n = O(\Delta t, \Delta x^2, \Delta y^2)$ accuracy

• stability restriction: $\nu \frac{\Delta t}{\Delta x^2} + \nu \frac{\Delta t}{\Delta y^2} \leq 1$ ← restrictive because explicit method

• easy to code, but time step would necessarily be too small.

b) Trapezoidal Rule \Rightarrow Crank-Nicholson Method

• $\tau_{j,k}^n = O(\Delta t^2, \Delta x^2, \Delta y^2)$ accuracy

• unconditionally stable ("can take any Δt you want, and scheme will be stable, but not necessarily accurate, so take Δt on order of $\Delta x, \Delta y$ ")

• implicit method so a little harder to code, (must solve linear systems)

c) ADI, Peaseman - Rachford

- 2-level alternating implicit scheme
- $\tau_{j,k}^n = O(\Delta t^2, \Delta x^2, \Delta y^2)$
- unconditionally stable
- moderately hard to code
- choose $\Delta t \sim \Delta x \sim \Delta y$ to balance

$$\frac{v_{j,k}^* - v_{j,k}^n}{\Delta t/2} = \nu \left[\frac{1}{\Delta x^2} \delta_x^2 v_{j,k}^* + \frac{1}{\Delta y^2} \delta_y^2 v_{j,k}^n \right]$$

implicit
explicit

First order operator δ_x

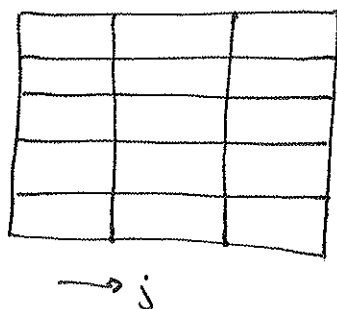
$$\frac{v_{j,k}^{n+1} - v_{j,k}^*}{\Delta t/2} = \nu \left[\frac{1}{\Delta x^2} \delta_x^2 v_{j,k}^* + \frac{1}{\Delta y^2} \delta_y^2 v_{j,k}^{n+1} \right]$$

explicit
implicit

You can show that this is 2nd order accurate in both time and space.

ADI Algorithm Structure

- 1) Initialization - grid, time-step, and other parameters...
set initial grid function
- 2) Main time-stepping loop
- 3) For $n=0, \dots, n_{\text{Final}}$



- 4) For $k=1, \dots, M-1$
 - 5) build tri-diagonal matrix system, solve it,
- $O(N)$ { save solution on grid line
 $O(MN)$ {
 • end
 we have $v_{j,k}^*$, don't need $v_{j,k}^n$

$O(MN)$ {

 6) For $j = 1, \dots, N-1$

 $O(M)$ {

 7) build tridiagonal matrix, solve it,

 save solution on grid line

 • end

 you have $v_{j,k}^{n+1}$ and don't need $v_{j,k}^*$

7) output results

 8) end time loop

\Rightarrow each time step costs $O(MN)$

 \Rightarrow total cost is $O(MN)$ per time step

 \Rightarrow this is optimal

Polar Coordinates

the domain is now a disk, $0 < r < 1$, $0 \leq \theta \leq 2\pi$, $t > 0$

$$v_t = \frac{1}{r}(rv_r)_r + \frac{1}{r^2}v_{\theta\theta}$$

where the diffusivity, $v = 1$

initial condition: $v(r, \theta, 0) = f(r, \theta)$

boundary condition: $v(1, \theta, t) = g(\theta, t)$

there are additional boundary conditions that are implied by this problem

$$v(r, \theta, t) = v(r, \theta + 2\pi, t) \quad (2\pi \text{ periodic})$$

the singularity at $r = 0$ is "ok"

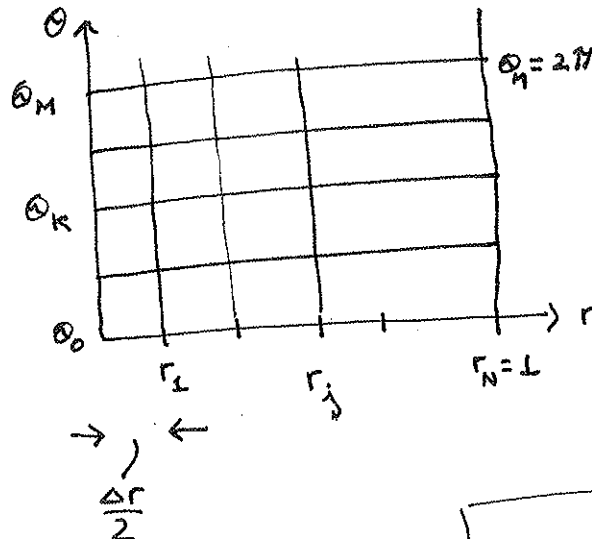
then enforce $\lim_{r \rightarrow 0} rv_r = 0$

grid: $r_j = (j - \frac{1}{2}) \Delta r$, want $r_N = 1$

→ solve for Δr , $1 = (N - \frac{1}{2}) \Delta r \rightarrow \boxed{\Delta r = \frac{1}{N - \frac{1}{2}}}$

$\theta_k = k \Delta \theta$, want $\theta_M = 2\pi \rightarrow \boxed{\Delta \theta = \frac{2\pi}{M}}$

polar grid:



let $v_{j,k}(t) \approx u(r_j, \theta_k, t)$

→
$$v'_{j,k}(t) = \frac{1}{r_j \Delta r^2} \left[r_{j+1/2} \delta_{+r} - r_{j-1/2} \delta_{-r} \right] v_{j,k} + \frac{1}{r_j^2} \frac{1}{\Delta \theta^2} \delta_{\theta}^2 v_{j,k}$$

to approximate $\frac{1}{r} (r u_r)_r$,

use

$$\frac{1}{r} \left[\frac{(r u_r)_{j+1/2} - (r u_r)_{j-1/2}}{\Delta r} \right]$$

where $v_{j,k}(0) = f(r_j, \theta_k)$

and $v_{N,k}(t) = g(\theta_k, t)$

and implied BCs: • periodicity, since $v_{j,0}(t) = v_{j,M}(t)$

then solve ODEs on grid for $k=1, \dots, M$

likewise, $\boxed{v_{j,M+1}(t) = v_{j,1}(t)}$

• at $j=1$, forward difference term is ok, but backward difference term will have $r_{j-1/2}$

$$\boxed{r_{j-1/2} \delta_{-r} v_{j,k} = 0 \text{ for } j=1}$$

Time-integration

want to handle @-differences implicitly because of $\frac{1}{r_j^2}$ coefficient

consider the following scheme:

$$\frac{V_{j,k}^{n+1} - V_{j,k}^n}{\Delta t} = \frac{1}{\Delta r^2} D V_{j,k}^n + \frac{1}{r_j^2} \cdot \frac{1}{\Delta \theta^2} \delta_\theta^2 V_{j,k}^{n+1}$$

where $D = \frac{1}{r_j} \left[r_{j+1/2} \delta_{+r} - r_{j-1/2} \delta_{-r} \right]$

combining terms, obtain

$$\left[1 - \frac{\Delta t}{r_j^2 \Delta \theta^2} \delta_\theta^2 \right] V_{j,k}^{n+1} = \left(1 + \frac{\Delta t}{\Delta r^2} D \right) V_{j,k}^n$$

consider a Fixed j grid line

due to periodicity

$$\begin{bmatrix} \times & \times & & \times \\ \times & \times & \times & \\ & \times & \times & \times \\ \times & & & \times & \times \end{bmatrix} \begin{pmatrix} V_{j,1}^{n+1} \\ V_{j,2}^{n+1} \\ \vdots \\ V_{j,M}^{n+1} \end{pmatrix} =$$

to solve this system, partition the matrix into actual tridiagonal matrix and border parts

$$\begin{bmatrix} \times & \times & & & & \times \\ & \times & \times & \times & & \\ & & \times & \times & \times & \\ & & & \times & \times & \times \\ & & & & \times & \times & \times \\ \times & & & & & \times & \times & \times \end{bmatrix}$$

Bordering Algorithms

denote this matrix as

$$\begin{bmatrix} \underline{A} & \vdots & \underline{b} \\ \vdots & \ddots & \vdots \\ \underline{c}^T & \vdots & \underline{d} \end{bmatrix} \begin{pmatrix} \underline{x} \\ \vdots \\ \underline{y} \end{pmatrix} = \begin{pmatrix} \underline{r} \\ \vdots \\ \underline{s} \end{pmatrix}$$

$$\left[\begin{array}{c|c} \underline{\underline{A}} & \underline{\underline{b}} \\ \hline \underline{\underline{c}}^T & d \end{array} \right] \begin{pmatrix} \underline{\underline{x}} \\ y \end{pmatrix} = \begin{pmatrix} \underline{\underline{r}} \\ s \end{pmatrix}, \text{ assume that } \underline{\underline{A}} \text{ is nonsingular}$$

$$\text{solve } \underline{\underline{A}} \underline{\underline{p}} = \underline{\underline{r}}, \quad \underline{\underline{A}} \underline{\underline{q}} = \underline{\underline{b}} \Rightarrow \underline{\underline{p}}, \underline{\underline{q}} \text{ are known}$$

$$\left[\begin{array}{c|c} \underline{\underline{A}}^{-1} & \underline{\underline{0}} \\ \hline \underline{\underline{0}} & 1 \end{array} \right] \left[\begin{array}{c|c} \underline{\underline{A}} & \underline{\underline{b}} \\ \hline \underline{\underline{c}}^T & d \end{array} \right] \begin{pmatrix} \underline{\underline{x}} \\ y \end{pmatrix} = \left[\begin{array}{c|c} \underline{\underline{A}}^{-1} & \underline{\underline{0}} \\ \hline \underline{\underline{0}} & 1 \end{array} \right] \begin{pmatrix} \underline{\underline{r}} \\ s \end{pmatrix}$$

$$\Rightarrow \left[\begin{array}{c|c} \underline{\underline{I}} & \underline{\underline{q}} \\ \hline \underline{\underline{c}}^T & d \end{array} \right] \begin{pmatrix} \underline{\underline{x}} \\ y \end{pmatrix} = \begin{pmatrix} \underline{\underline{p}} \\ s \end{pmatrix}$$

$$\text{now multiply all of this times } \left[\begin{array}{c|c} \underline{\underline{I}} & \underline{\underline{0}} \\ \hline \underline{\underline{c}}^T & 1 \end{array} \right] \Rightarrow \left[\begin{array}{c|c} \underline{\underline{I}} & \underline{\underline{0}} \\ \hline \underline{\underline{c}}^T & 1 \end{array} \right] \left[\begin{array}{c|c} \underline{\underline{I}} & \underline{\underline{q}} \\ \hline \underline{\underline{c}}^T & d \end{array} \right]$$

$$\Rightarrow \left[\begin{array}{c|c} \underline{\underline{I}} & \underline{\underline{q}} \\ \hline \underline{\underline{0}} & \tilde{d} \end{array} \right] \begin{pmatrix} \underline{\underline{x}} \\ y \end{pmatrix} = \begin{pmatrix} \underline{\underline{p}} \\ \tilde{s} \end{pmatrix} \quad \text{where } \tilde{d} = d - \underline{\underline{c}}^T \underline{\underline{q}}, \quad \tilde{s} = s - \underline{\underline{c}}^T \underline{\underline{p}}$$

$$\Rightarrow \text{results : } y = \frac{\tilde{s}}{\tilde{d}}, \quad \underline{\underline{x}} = \underline{\underline{p}} - y \underline{\underline{q}}$$

A nonlinear Problem For the Heat Eqn

$$u_t = u_{xx} + b(u), \quad 0 < x < 1, \quad t > 0$$

where $b(u)$ is some nonlinear source

initial condition, $u(x, 0) = f(x)$

boundary conditions: $u(0, t) = u(1, t) = 0$

Finite Difference approximation.

grid: $x_j = j \Delta x$, $\Delta x = \frac{1}{N}$, $v_j(t) \approx u(x_j, t)$ - method of lines

let $v_j(t)$ solve

$$v_j'(t) = \frac{1}{\Delta x^2} \delta_x^2 v_j^n + b(v_j^n), \quad 1 \leq j \leq N-1 \quad (\text{interior grid lines})$$

$$v_j(0) = f(x_j), \quad v_0(t) = v_N(t) = 0$$

We want to consider a 2nd order implicit integration

$$\frac{v_j^{n+1} - v_j^n}{\Delta t} = \frac{1}{\Delta x^2} \delta_x^2 \left(\frac{v_j^n + v_j^{n+1}}{2} \right) + b\left(\frac{v_j^n + v_j^{n+1}}{2} \right), \quad 1 \leq j \leq N, \quad n \geq 0$$

$$v_j^0 = f(x_j), \quad v_0^n = v_N^n = 0 \quad \text{"Fully discrete difference scheme"}$$

The question is how do we solve for v_j^{n+1}

$$F_j \equiv v_j^{n+1} - v_j^n - r \delta_x^2 \left(\frac{v_j^n + v_j^{n+1}}{2} \right) - \Delta t b\left(\frac{v_j^n + v_j^{n+1}}{2} \right) = 0, \quad r = \frac{\Delta t}{\Delta x^2}$$

this is a system of nonlinear equations:

$$F(v) = 0, \quad v = [v_1^{n+1} \quad v_2^{n+1} \quad \dots \quad v_{N-1}^{n+1}]^T$$

($N-1$ nonlinear algebraic equations)

use Newton's method to solve

Initial guess: $v^{(0)} = v$ at the n^{th} time level

For $p = 0, 1, \dots$

solve $\underline{F}_v(v^{(p)}) \underline{dv}^{(p)} = \underline{F}(v^{(p)})$ where \underline{F}_v is the Jacobian

update $v^{(p+1)} = v^{(p)} - \underline{dv}^{(p)}$

stop when $\|\underline{dv}^{(p)}\| \leq \text{tol} \cdot \|v^{(p)}\|$ where $\text{tol} \sim 10^{-8}$

if \underline{F}_v is nonsingular at root and initial guess sufficiently close, then convergence is quadratic

Form of \underline{F}_v : known

$$F_1 = v_1^{n+1} - v_1^n - \frac{r}{2} (v_2^{n+1} - 2v_1^{n+1} + v_0^{n+1}) - \frac{r}{2} \delta_x^2 v_1^n - \Delta t b b' \left(\frac{v_1^{n+1} + v_1^n}{2} \right)$$

$$\frac{\partial F_1}{\partial v_1} = 1 + r - \Delta t b' \left(\frac{v_1^{n+1} + v_1^n}{2} \right) \cdot \frac{1}{2}$$

$$\frac{\partial F_2}{\partial v_2} = -\frac{r}{2}, \quad \frac{\partial F_1}{\partial v_j} = 0, \quad j \geq 3$$

similarly

$$\frac{\partial F_2}{\partial v_1} = -\frac{r}{2}, \quad \frac{\partial F_2}{\partial v_2} = 1 + r - \Delta t b' \left(\frac{v_1^{n+1} + v_1^n}{2} \right) \cdot \frac{1}{2}, \quad \frac{\partial F_2}{\partial v_3} = -\frac{r}{2}, \quad \frac{\partial F_2}{\partial v_j} = 0, \quad j \geq 4$$

therefore \underline{F}_v is tridiagonal

$$\frac{\partial F}{\partial v} = \begin{bmatrix} \beta_1 & \gamma_1 & & \\ \alpha_2 & \beta_2 & \gamma_2 & \\ & \alpha_3 & \beta_3 & \gamma_3 \\ & & \ddots & \ddots \end{bmatrix}$$

where $\alpha_j = \gamma_j = -\frac{r}{2}$

$$\beta_j = 1 + r - \frac{\Delta t}{2} b' \left(\frac{v_j^{n+1} + v_j^n}{2} \right)$$

Convergence, consistency, stability

Error: $w_j^n = v_j^n - u_j^n$ where $u_j^n = u(x_j, t_n)$

substitute $v_j^n = w_j^n + u_j^n$ into difference scheme

$$u_j^{n+1} + w_j^{n+1} - (u_j^n + w_j^n) = r \delta_x^2 \left(\frac{u_j^n + u_j^{n+1}}{2} \right) + r \delta_x^2 \left(\frac{w_j^n + w_j^{n+1}}{2} \right) + \Delta t b \left(\frac{u_j^n + u_j^{n+1}}{2} + \frac{w_j^n + w_j^{n+1}}{2} \right)$$

expand about $b \left(\frac{u_j^n + u_j^{n+1}}{2} \right)$

$$\rightarrow \Delta t b \left(\frac{u_j^n + u_j^{n+1}}{2} + \frac{w_j^n + w_j^{n+1}}{2} \right) = b \left(\frac{u_j^n + u_j^{n+1}}{2} \right) + b'(\theta_j^n) \frac{w_j^n + w_j^{n+1}}{2}$$

\rightarrow Taylor series, one term, with remainder

where θ_j^n is between $\frac{u_j^n + u_j^{n+1}}{2}$ and $\frac{u_j^n + u_j^{n+1}}{2} + \frac{w_j^n + w_j^{n+1}}{2}$

Truncation error:

$$\Delta t \tau_j^n = u_j^{n+1} - u_j^n - r \delta_x^2 \left(\frac{u_j^n + u_j^{n+1}}{2} \right) - \Delta t b \left(\frac{u_j^n + u_j^{n+1}}{2} \right), \quad r = \frac{\Delta t}{\Delta x^2}$$

$$= \Delta t u_t + \frac{\Delta t^2}{2} u_{tt} + O(\Delta t^3) - \frac{\Delta t}{2} \left[u_{xx} + \frac{\Delta x^2}{12} u_{xxxx} + O(\Delta x^4) + u_{xx} + \Delta t u_{xxt} + O(\Delta t^2) + \frac{\Delta x^2}{12} u_{xxxx} + O(\Delta x^2 \Delta t) \right] - \Delta t b \left[u + \frac{\Delta t}{2} u_t + O(\Delta t^2) \right]$$

Note, $b(u + \frac{\Delta t}{2} u_t + \dots) = b(u) + b'(u) \frac{\Delta t}{2} u_t + \dots$

$$\Rightarrow \cancel{\tau_j^n = \cancel{u_t} + \frac{\Delta t}{2} \cancel{u_{tt}} + O(\Delta t^2) - \cancel{u_{xx}} - \frac{\Delta t}{2} \cancel{u_{xxt}} + O(\Delta x^2)}$$

$$\quad \quad \quad \underbrace{-b(u) - \frac{\Delta t}{2} b'(u) u_t + O(\Delta t^2)}_{\substack{\text{time derivative} \\ \text{of} \\ \text{PDE}}}$$

from PDE

$$\Rightarrow \boxed{\tau_j^n = O(\Delta t^2) + O(\Delta x^2)}$$

So back to the error equation

$$w_j^{n+1} = w_j^n + \delta x^2 \left(\frac{w_j^n + w_j^{n+1}}{2} \right) + \Delta t b'(\theta_j^n) \left(\frac{w_j^n + w_j^{n+1}}{2} \right) - \Delta t \tau_j^n$$

to analyze this error equation, need to freeze $b'(\theta_j^n)$ to be some constant and consider the stability of the resulting linear constant coeff equation for all possible values for $\tau = b'(\theta_j^n)$

\Rightarrow necessary condition for stability (in order to obtain sufficient condition, need to do nonlinear analysis)

Midterm

closed book, closed notes, 1 crib sheet

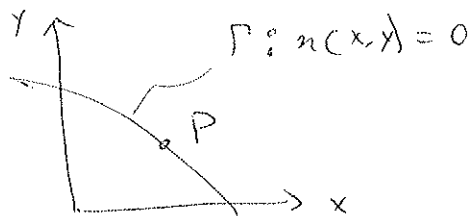
1) classification of PDEs

$$Au_{xx} + 2Bu_{xy} + Cu_{yy} = D$$

nonlinear PDE, but linear in highest-derivatives

so we call it quasi-linear

issue of classification boils down to being able to construct equation locally



can you construct solution u locally about point P , given u and u_n on Γ at P

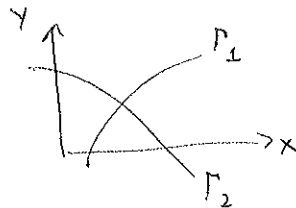
Cases

1) $B^2 - AC > 0$ - hyperbolic, may or may not be able to construct (from Taylor series) solution u depending on whether Γ is a characteristic

differential form of characteristic equation

$$A dy^2 - 2B dx dy + C dx^2 = 0$$

has real solutions



issues include domain of dependence region of influence

2) $B^2 - AC < 0$ - elliptic - no real characteristics, solutions of equations are analytic

3) $B^2 - AC = 0$, parabolic case - may or may not be able to build solution locally

◆ Canonical Forms

$$u_{xx} + u_{yy} = \delta, \text{ elliptic}$$

$$u_{xy} = \delta, \text{ hyperbolic}$$

$$u_{xx} = \delta, \text{ parabolic}$$

where δ only depends on lower derivatives

◆ Boundary conditions for well-posed problems

2) Finite Difference Methods

- heat eqn w/ Dirichlet BCs

- heat eqn w/ Neumann BCs

- ghost lines, staggered grids

- generating discrete approximation

- Taylor series approach

- interpolation approach

- finite-volumes, discrete conservation

- Theoretical issues

- consistency, order of accuracy, truncation error

- convergence, global error

- stability - do errors in difference equation grow?
well-posedness of the difference equation

- Lax theorems

- if consistent

- stability analysis

- Fourier stability analysis, Discrete Fourier Transform

- Fourier mode analysis

- Matrix analysis, if boundaries included

- Parabolic Eqns

- Crank-Nicholson method

- ADI

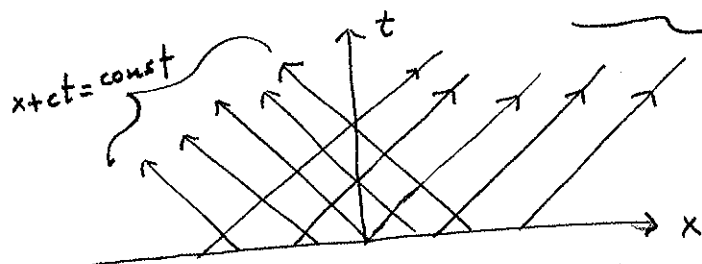
- Polar coordinates

- Non linear Eqn

Hyperbolic PDEs

A hyperbolic PDE is one that possesses a full set of real characteristics.

Example, wave equation, $u_{tt} = c^2 u_{xx}$. This PDE has two families of real characteristics, $x - ct = \text{constant}$, $x + ct = \text{constant}$.



Can also consider the wave equation as a system of two first-order equations.

$$u_t + \underline{A} u_x = 0$$

where $u = \begin{bmatrix} u^{(1)} \\ u^{(2)} \end{bmatrix}$, $\underline{A} = \begin{bmatrix} 0 & c \\ c & 0 \end{bmatrix} \Rightarrow \begin{aligned} u_t^{(1)} + c u_x^{(2)} &= 0 \quad (1) \\ u_t^{(2)} + c u_x^{(1)} &= 0 \quad (2) \end{aligned}$

Take $\frac{\partial}{\partial t}$ of eqn (1) and subtract $c \frac{\partial}{\partial x}$ of eqn (2) $\rightarrow u_{tt}^{(1)} - c^2 u_{xx}^{(1)} = 0$

Therefore the first component of u satisfies the wave equation. Similarly, you can show that the second component also satisfies the wave equation.

Consider the linear system of first order PDEs

$$u_t + \underline{A} u_x = 0, \quad \underline{A} \in \mathbb{R}^{m \times m}, \quad u(x, t) = [u^{(1)} \ u^{(2)} \ \dots \ u^{(m)}]^T$$

This system is called hyperbolic if A is diagonalizable with real eigenvalues

Assume that \underline{A} is diagonalizable with real eigenvalues. Then a nonsingular matrix \underline{R} exists such that $\underline{R}^{-1} \underline{A} \underline{R} = \underline{\Lambda}$ where $\underline{\Lambda} = \text{diag} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_m \end{bmatrix}$, a ^{diagonal} matrix with ^{real} eigenvalues of \underline{A} on diagonals.

Set $\underline{w}(x, t) = \underline{R}^{-1} \underline{u}(x, t)$ = "characteristic variables", "Riemann variables"

Note: rows of \underline{R}^{-1} are left eigenvectors and columns of \underline{R} are the right eigenvectors

$$\rightarrow \underline{R}^{-1} \left(\underline{u}_t + \underbrace{\underline{A}}_{\underline{R} \underline{R}^{-1}} \underline{u}_x \right) = 0 \rightarrow \underline{R}^{-1} \underline{u}_t + \underline{R}^{-1} \underline{A} \underline{R} \underline{R}^{-1} \underline{u}_x = 0$$

$$\rightarrow \left(\underline{R}^{-1} \underline{u} \right)_t + \underline{\Lambda} \left(\underline{R}^{-1} \underline{u} \right)_x = 0 \rightarrow \underline{w}_t + \underline{\Lambda} \underline{w}_x = 0$$

$$\rightarrow \frac{\partial}{\partial t} \begin{bmatrix} w^{(1)} \\ w^{(2)} \\ \vdots \\ w^{(m)} \end{bmatrix} + \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_m \end{bmatrix} \frac{\partial}{\partial x} \begin{bmatrix} w^{(1)} \\ w^{(2)} \\ \vdots \\ w^{(m)} \end{bmatrix}$$

the equations have decoupled, the p^{th} equation is

$$w_t^{(p)} + \lambda_p w_x^{(p)} = 0, p=1, \dots, m$$

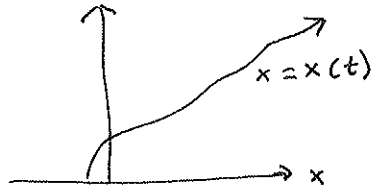
\rightarrow each characteristic variable satisfies a scalar linear advection equation

Linear Advection Equation

$$u_t + c u_x = 0, \quad c \in \mathbb{R}, \quad |x| < \infty, \quad t > 0$$

Solve using the method of characteristics. Consider the rate of change of u with respect to t along a path

$$x = x(t)$$



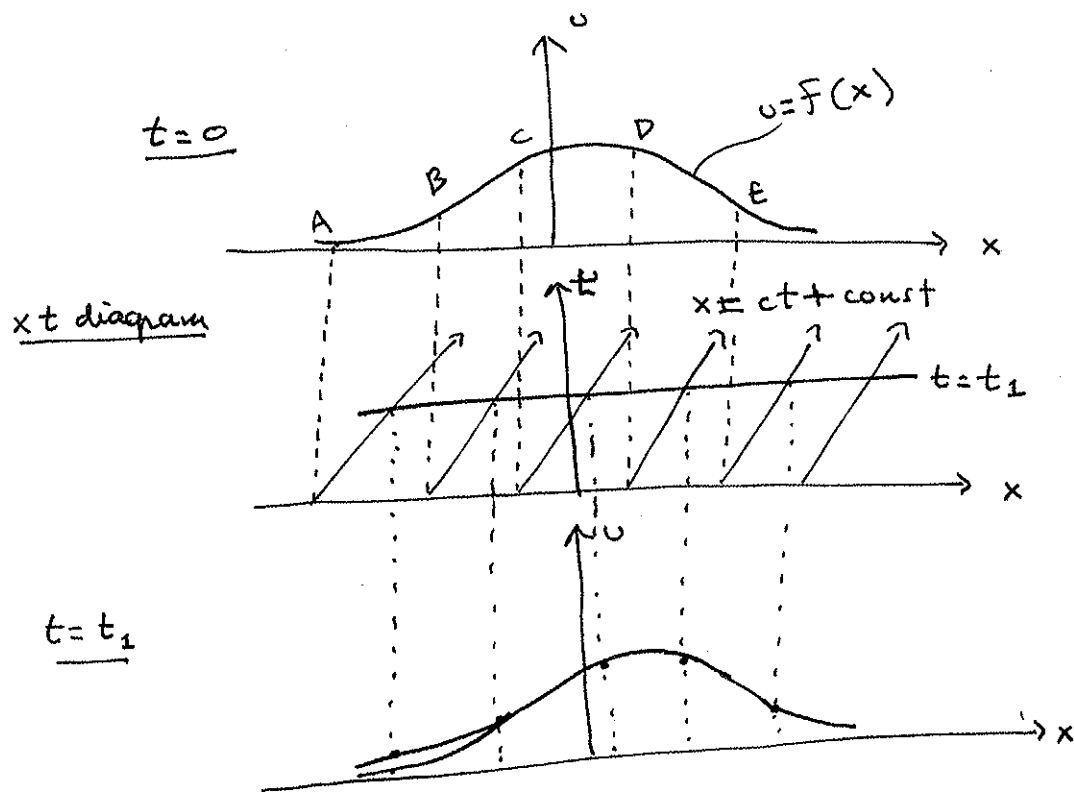
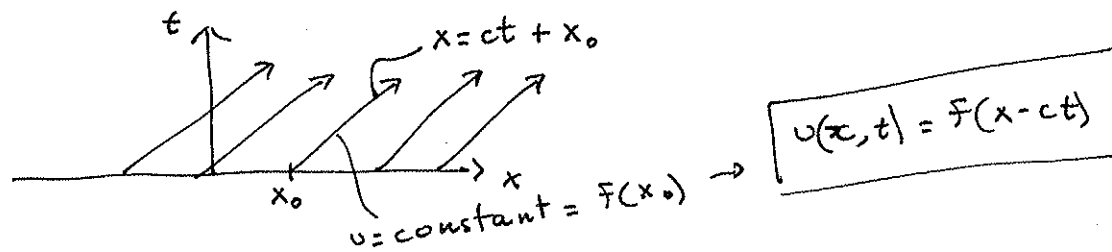
$$\frac{d}{dt} u(x(t), t) = \frac{\partial u}{\partial x} \frac{dx}{dt} + \frac{\partial u}{\partial t}$$

select $\frac{dx}{dt} = c$ then $\frac{d}{dt} u(x(t), t) = 0$

→ characteristic equations are $\frac{du}{dt} = 0$ along $\frac{dx}{dt} = c$

solve the characteristic equations

→ $u = \text{constant}$ along $x = ct + \text{constant}$



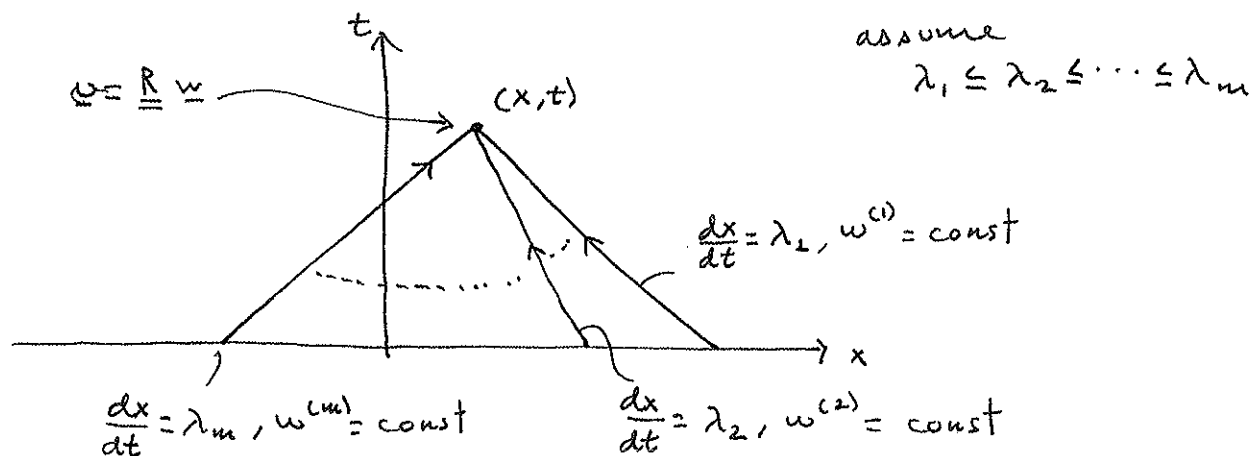
the solution translates to the right at speed $c > 0$ and unchanged shape

Return to the linear system

$$\frac{\partial}{\partial t} w^{(p)} + \lambda_p \frac{\partial}{\partial x} w^{(p)} = 0, \quad p = 1, 2, \dots, m$$

solution

$$w^{(p)} = \text{constant along } x = \lambda_p t + \text{constant}$$



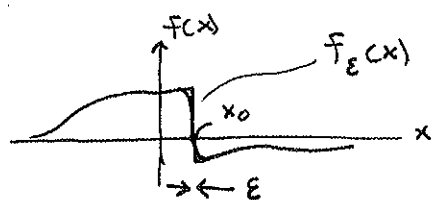
if $u(x, 0)$ is given, then $w(x, 0) = R^{-1} u(x, 0)$,
the characteristic variables initial values are known

Behavior of Discontinuities For Linear Equations

$$u_t + cu_x = 0, \quad |x| < \infty, \quad t > 0, \quad u(x, 0) = F(x)$$

solution: $u(x, t) = F(x - ct)$

Suppose F has the form:

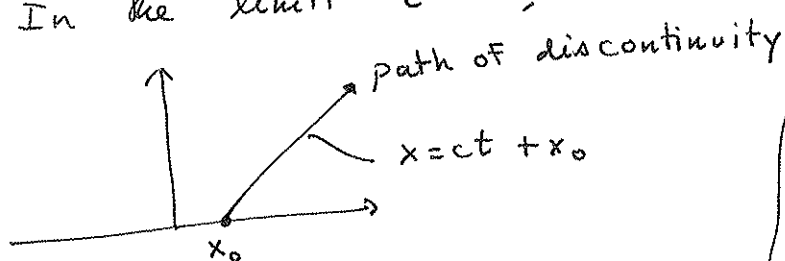


In the limit $\epsilon \rightarrow 0$, $F(x)$ becomes discontinuous at x_0

so $F_\epsilon(x)$ is smooth if $\epsilon > 0$

$$\rightarrow u(x, t) = F_\epsilon(x - ct), \quad \epsilon > 0$$

In the limit $\epsilon \rightarrow 0$, $u(x, t) = F(x - ct)$



For a linear problem, the discontinuities are a result of a discontinuity in the initial condition

Numerical Methods

scalar advection equation, $u_t + cu_x = 0$, $|x| < \infty$, $t > 0$, $f(x) = u(x, 0)$

grid: $x_j = j \Delta x$, $\Delta x = \text{given}, > 0$

$t_n = n \Delta t$, $\Delta t = \text{given}, > 0$

let $v_j^n \approx u(x_j, t_n)$

standard methods:

◆ "centered" methods:

Lax-Friedrichs method:

$$v_j^{n+1} = \frac{1}{2} (v_{j-1}^n + v_{j+1}^n) - \frac{c \Delta t}{2 \Delta x} (v_{j+1}^n - v_{j-1}^n)$$

Lax-Wendroff method:

$$v_j^{n+1} = v_j^n - \frac{c \Delta t}{2 \Delta x} (v_{j+1}^n - v_{j-1}^n) + \frac{1}{2} \left(\frac{c \Delta t}{\Delta x} \right)^2 (v_{j+1}^n - 2v_j^n + v_{j-1}^n)$$

◆ "one-sided" methods (upwind methods)

$$v_j^{n+1} = v_j^n - \frac{c \Delta t}{\Delta x} (v_j^n - v_{j-1}^n)$$

$$v_j^{n+1} = v_j^n - \frac{c \Delta t}{\Delta x} (v_{j+1}^n - v_j^n)$$

◆ Remarks

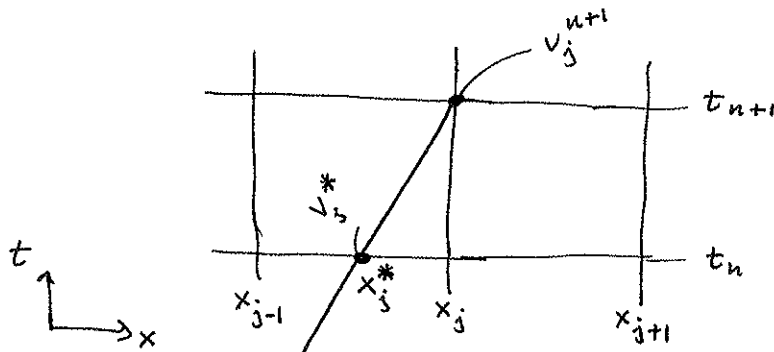
1) Most methods for hyperbolic equations are explicit because the stability constraints for these methods are $\Delta t \leq \text{const} \cdot \Delta x$ which is probably what you would pick for accuracy

2) First choice might be a method of lines construction
 set $v_j(t) \approx u(x_j, t) \rightarrow v_j'(t) + c \underbrace{\frac{v_{j+1} - v_{j-1}}{2 \Delta x}}_{\text{2nd order centered approximation for } u_x} = 0$
 then use Forward Euler to integrate:

$$\frac{v_j^{n+1} - v_j^n}{\Delta t} + \frac{c}{2 \Delta x} (v_{j+1}^n - v_{j-1}^n) = 0$$

$$\rightarrow v_j^{n+1} = v_j^n - \frac{c \Delta t}{2 \Delta x} (v_{j+1}^n - v_{j-1}^n) \rightarrow \text{UNCONDITIONALLY UNSTABLE}$$

3) Previous methods can be derived by consideration of characteristics



characteristic, $\frac{dx}{dt} = c$

From characteristic construction,

$$v_j^{n+1} = v_j^* \quad (\text{exact})$$

Approximate value of v_j^* using grid data at t_n via some form of interpolation:

- linear interpolation yields the upwind method
consider data points (x_{j-1}, v_{j-1}^n) and (x_j, v_j^n)
or perhaps data points (x_j, v_j^n) and (x_{j+1}, v_{j+1}^n) (1st order)
- Interpolate v_j^* using a linear fit to
 $(x_{j-1}^n, v_{j-1}^n), (x_{j+1}^n, v_{j+1}^n) \Rightarrow$ Lax-Friedrichs (1st order)
- Interpolate v_j^* using a quadratic fit to
 ~~(x_{j-k}^n, v_{j-k}^n)~~ For $k = -1, 0, 1 \Rightarrow$ Lax-Wendroff
2nd order accurate

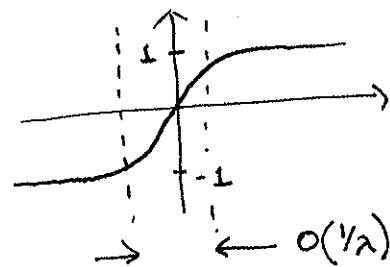
Example: linear advection equation

$$u_t + cu_x = 0, \quad t > 0,$$

$$u(x, 0) = \tanh \lambda x, \quad \lambda = \text{const.}$$

solution: $u(x, t) = \tanh \lambda(x - ct)$

the parameter λ determines the transition width, $O(1/\lambda)$



Compare with various numerical approximations for different values of Δx and λ .

1st order scheme: Lax-Friedrichs

$$v_j^{n+1} = \frac{1}{2}(v_{j-1}^n + v_{j+1}^n) - \frac{c\Delta t}{2\Delta x}(v_{j+1}^n - v_{j-1}^n)$$

2nd order scheme: Lax-Wendroff

$$v_j^{n+1} = v_j^n - \frac{c\Delta t}{2\Delta x}(v_{j+1}^n - v_{j-1}^n) + \frac{1}{2}\left(\frac{c\Delta t}{\Delta x}\right)^2(v_{j+1}^n - 2v_j^n + v_{j-1}^n)$$

Compare and contrast these two schemes.

Integrate from $t=0$ to $t=5$

with $\Delta t = \tau \Delta x$ where $\tau \approx 1$

(Note, $\tau=1$ is the stability limit, where both methods become exact.)

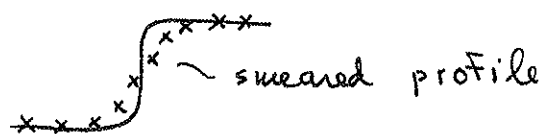
Observations:

- 1) when λ is small and the solution is well represented on the grid, then

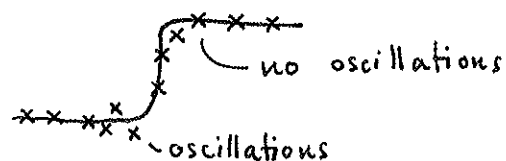
$$\left. \begin{array}{l} \text{error (LF) is } O(\Delta t, \Delta x) \\ \text{error (LW) is } O(\Delta t^2, \Delta x^2) \end{array} \right\} \text{expected behavior}$$

- 2) when λ is large and the solution is not well represented on the grid, then we observe significant error near the transition

LF, the error is smooth



LW, the error shows oscillations



- 3) IF $|v| = 1 \Rightarrow$ numerical solution becomes exact for this simple equation (because the characteristic is exactly aligned with the grid)

Analyze The Behavior of the Error For Both Methods

Lax-Friedrichs:
$$v_j^{n+1} = \frac{1}{2} (v_{j-1}^n + v_{j+1}^n) - \frac{c \Delta t}{2 \Delta x} (v_{j+1}^n - v_{j-1}^n)$$

solve for truncation error, τ_j^n :

$$v_j^{n+1} = \frac{1}{2} (v_{j-1}^n + v_{j+1}^n) - \frac{c \Delta t}{2 \Delta x} (v_{j+1}^n - v_{j-1}^n) + \Delta t \tau_j^n$$

$$\begin{aligned} \rightarrow v + \Delta t v_t + \frac{\Delta t^2}{2} v_{tt} + \dots &= \frac{1}{2} (2v + \Delta x^2 v_{xx} + \dots) \\ &- \frac{c \Delta t}{2 \Delta x} (2 \Delta x v_x + O(\Delta x^3)) + \Delta t \tau_j^n \end{aligned}$$

Divide through by Δt

$$u_t + \frac{\Delta t}{2} u_{tt} + O(\Delta t^2) = \frac{\Delta x^2}{2\Delta t} u_{xx} + O(\Delta x^4/\Delta t) - cu_x + O(\Delta x^2) + \tau$$

- If we assign $u(x,t)$ as a solution of $u_t + cu_x = 0$, then the truncation error, $\tau = (\Delta t, \frac{\Delta x^2}{\Delta t})$ = First order accurate provided $\frac{\Delta x}{\Delta t} = \text{constant}$, which it is.

- If we assign $u(x,t)$ as a solution of the modified equation

$$u_t + cu_x = -\frac{\Delta t}{2} u_{tt} + \frac{\Delta x^2}{2\Delta t} u_{xx}$$

then the truncation error is smaller. In other words, the numerical solution approximates the solution of the modified equation better than the original.

$$u_t = -cu_x + \dots \leftarrow \text{modified equations}$$

$$\rightarrow u_{tt} = -cu_{xt} + \dots = -c(u_t)_x + \dots = -c(-cu_x)_x + \dots$$

$$\rightarrow u_{tt} = c^2 u_{xx} + \dots \quad \text{higher order in } \Delta t, \Delta x \text{ so } \underline{\text{ignored}}$$

$$\Rightarrow u_t + cu_x = -\frac{\Delta t}{2} (c^2 u_{xx} + \dots) + \frac{\Delta x^2}{2\Delta t} u_{xx}$$

Modified equation (to leading order)

$$u_t + cu_x = \frac{\Delta t}{2} \left(\frac{\Delta x^2}{\Delta t^2} - c^2 \right) u_{xx}$$

$$\rightarrow u_t + cu_x = \frac{c^2 \Delta t}{2} \left(\frac{\Delta x^2}{c^2 \Delta t^2} - 1 \right) u_{xx}$$

$$\rightarrow u_t + cu_x = \frac{c^2 \Delta t}{2} \left(\frac{1}{r^2} - 1 \right) u_{xx}, \quad r = \frac{c \Delta t}{\Delta x}$$

The general form of the modified equation is thus $u_t + cu_x = \nu u_{xx}$, where $\nu = \frac{c\Delta t}{2} \left(\frac{1}{r^2} - 1 \right)$

The term νu_{xx} provides the leading behavior of the error in Lax-Friedrichs. It is a diffusive term,

and therefore the leading behavior of error is diffusive (or dissipative) for Lax-Friedrichs

If we consider ν to be a diffusivity, then we'll need $\nu \geq 0$, otherwise end up with ill-posed backwards heat equation, so a stability constraint falls out of the analysis.

Fundamental solutions

let $u(x,t) = e^{\lambda t + ikx}$. Substitute into modified equation:

$$\lambda e^{\lambda t + ikx} + ikc e^{\lambda t + ikx} = -\nu k^2 e^{\lambda t + ikx}$$

$$\rightarrow \lambda = -ikc - \nu k^2 \Rightarrow u(x,t) = e^{ik(x-ct) - \nu k^2 t}$$

the $e^{ik(x-ct)}$ represents advection whereas as the $e^{-\nu k^2 t}$ represents a small amplitude decay

$$\text{general solution: } u(x,t) = \int_{-\infty}^{\infty} c(k) e^{ik(x-ct) - \nu k^2 t} dk$$

Behavior of the Error For Numerical Methods of the Linear Advection Equation

$$u_t + u_x = 0, \quad (c=1), \quad u(x,0) = F(x) = \tanh \lambda x$$

$$\text{exact solution } u(x,t) = F(x-ct)$$

Lax-Friedrichs (typical first order method)

$$v_i^{n+1} = \frac{1}{2} (v_{i-1}^n + v_{i+1}^n) - \frac{\Delta t}{2\Delta x} (v_{i+1}^n - v_{i-1}^n)$$

one method of studying error is to derive the modified equation. The idea is to replace the leading order term of the truncation error as a forcing term.

Modified Equation:

$$v_i^{n+1} = \frac{1}{2} (v_{i-1}^n + v_{i+1}^n) - \frac{\Delta t}{2\Delta x} (v_{i+1}^n - v_{i-1}^n) + \Delta t \tau_i^n$$

Expand in Taylor series and retain the leading order terms in Δt :

$$u_t + u_x = \nu u_{xx}, \quad \nu = \frac{\Delta t}{2} \left(\frac{1}{\sigma^2} - 1 \right), \quad \sigma = \frac{\Delta t}{\Delta x}$$

advection-diffusion equation

the leading behavior of the error is diffusive

Fundamental solution of modified eqn, $u = e^{\lambda t + i k x}$

$$\rightarrow u_k = e^{i k (x-t) - \nu k^2 t}$$

by superposition,

$$u(x,t) = \int_{-\infty}^{\infty} c(k) e^{i k (x-t) - \nu k^2 t} dk$$

where $c(k)$ determined by initial solution

$$\text{At } t=0, u(x, 0) = \int_{-\infty}^{\infty} c(k) e^{ikx} dk = F(x)$$

The question is, when does the term $e^{-\nu k^2 t}$ contribute?

The error is related to the size of $\nu k^2 t$.

If our initial state $F(x)$ is well represented on the grid, then the coefficient $|c(k)|$ would be small when $|k\Delta x|$ is large \Rightarrow error is small.

$$e^{ik(x-t) - \nu k^2 t} = e^{k(i(x-t) - \nu kt)} \quad , \quad \nu = \frac{\Delta t}{2} \left(\frac{1}{\sigma^2} - 1 \right) \quad , \quad \sigma = \frac{\Delta t}{\Delta x}$$

All of this is true for t sufficiently small.

If $F(x)$ is not well represented on the grid then $|c(k)|$ is large when $|k\Delta x|$ is large \Rightarrow significant error.

Lax-Wendroff Methods

$$v_j^{n+1} = v_j^n - \frac{\Delta t}{2\Delta x} \left(v_{j+1}^n - v_{j-1}^n \right) + \frac{1}{2} \left(\frac{\Delta t}{\Delta x} \right)^2 \left(v_{j+1}^n - 2v_j^n + v_{j-1}^n \right)$$

with initial condition $v_j^0 = F(x_j)$

Modified equation: $u_t + u_x = \nu u_{xxx}$, $\nu = -\frac{\Delta t^2}{6} \left(\frac{1}{\sigma^2} - 1 \right)$, $\sigma = \frac{\Delta t}{\Delta x}$

the dominant behavior of the error is no longer diffusive

$u_t + u_x = \nu u_{xxx} \rightarrow$ linear KDV equation, arises in shallow water waves
can derive solution by separation of variables

— Dispersive effect —

$$u_t + u_x = \nu u_{xxx}$$

substitute $u = e^{\lambda t + i k x}$

$$\rightarrow \lambda = -i k - i \nu k^3 \rightarrow u_k = e^{-i k t - i \nu k^3 t + i k x}$$

$$\rightarrow u_k = e^{i k (x - (1 + \nu k^2) t)} \rightarrow u_k = e^{i k (x - (1 + \nu k^2) t)}$$

→ general solution:
$$u(x, t) = \int_{-\infty}^{\infty} c(k) e^{i k (x - (1 + \nu k^2) t)} dk$$
 phase error, $\nu = -\frac{\Delta t^2}{6} \left(\frac{1}{r^2} - 1 \right)$

IF $f(x)$ is well represented on grid then $|c(k)|$ is small when $|k \Delta x|$ is large \Rightarrow ~~phase error~~ no sign error. IF $f(x)$ is not well represented on the grid then the magnitude of $c, |c|$ is large when $|k \Delta x|$ is large \Rightarrow sign phase error, dispersive.

The next term in the truncation error is a dissipation term, so there is some dissipation, but at high order.

Stability and the CFL condition (Courant - Friedrichs - Levy)

Look at one sided methods for $u_t + c u_x = 0$.

(1) $v_j^{n+1} = v_j^n - \tau (v_j^n - v_{j-1}^n)$, $\tau = \frac{c \Delta t}{\Delta x}$, Backwards Diff

(2) $v_j^{n+1} = v_j^n - \tau (v_{j+1}^n - v_j^n)$, Forward Diff

stability analysis, $v_j^n = a^n e^{i k x_j}$

(1) $\Rightarrow a = 1 - \tau (1 - e^{-i k \Delta x})$

(2) $\Rightarrow a = 1 - \tau (e^{i k \Delta x} - 1)$

$$(1) : a = 1 - \tau(1 - e^{-ik\Delta x})$$

$$(2) : a = 1 - \tau(e^{ik\Delta x} - 1)$$

$$(1) \Rightarrow |a|^2 = \left(1 - \tau(1 - \cos k\Delta x)\right)^2 + \tau^2 \sin^2 k\Delta x \leq 1$$

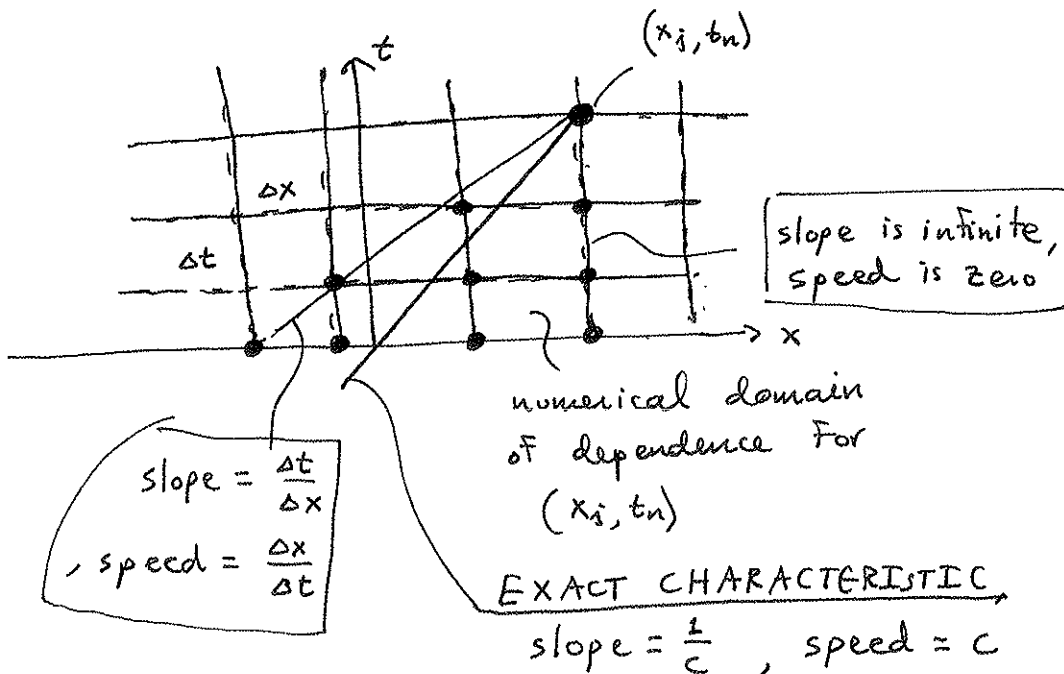
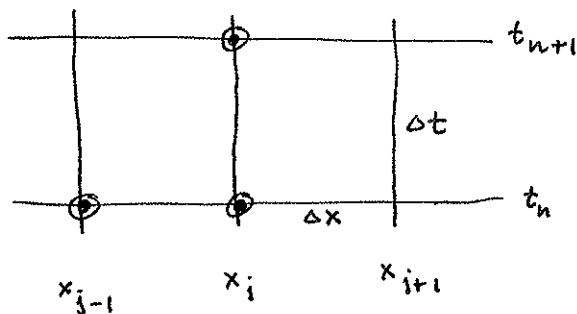
$$\Rightarrow \tau(1 - \tau) \geq 0 \Rightarrow \underline{0 \leq \tau \leq 1 \text{ For stability}}$$

$$(2) \Rightarrow |a|^2 = \left(1 + \tau(1 - \cos k\Delta x)\right)^2 + \tau^2 \sin^2 k\Delta x \leq 1$$

$$\Rightarrow \underline{-1 \leq \tau \leq 0 \text{ For stability}}$$

We can interpret these results in terms of the domain of dependence.

scheme (1) Backward Difference, $v_j^{n+1} = v_j^n - \tau(v_j^n - v_{j-1}^n)$

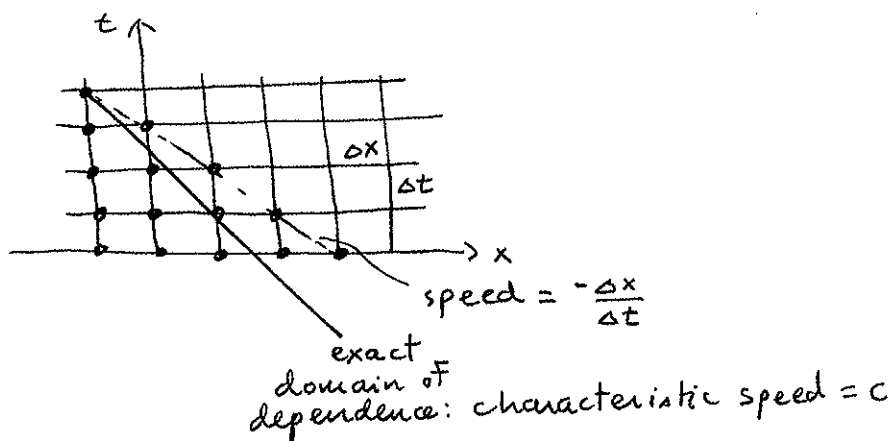


Stability result for (1), backwards diff

$$\Rightarrow 0 \leq r \leq 1 \Rightarrow 0 \leq \frac{c \Delta t}{\Delta x} \leq 1 \Rightarrow 0 \leq c \leq \frac{\Delta x}{\Delta t}$$

\Rightarrow The exact domain of dependence must be contained within the numerical domain of dependence \Rightarrow CFL condition (necessary condition for stability)

Similarly, for (2), Forwards Difference



$$\Rightarrow -1 \leq r \leq 0 \Rightarrow -\frac{\Delta x}{\Delta t} \leq c \leq 0$$

Non-Constant Coeffs

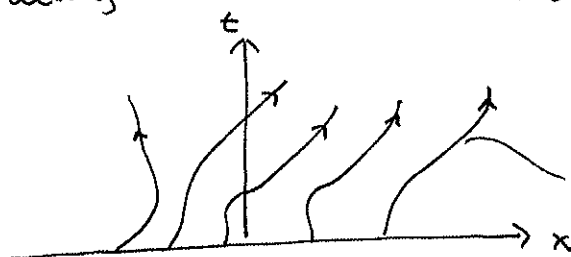
General linear scalar eqn:

$$u_t + c(x,t) u_x = a(x,t) u + b(x,t)$$

characteristic form:

$$\frac{du}{dt} = a(x(t), t) u + b(x(t), t), \quad u(x, 0) = u(x_0, 0)$$

$$\text{along characteristics } \frac{dx}{dt} = c(x, t), \quad x(0) = x_0$$



non-overlapping paths

along each path, u evolves according to the characteristic form

Numerical solution

$$u_t + c(x,t)u_x = a(x,t)u + b(x,t)$$

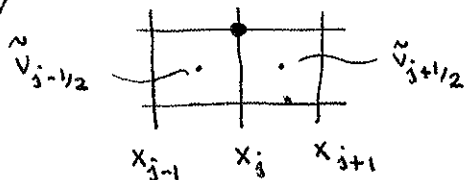
◆ lower order method, such as upwind:

$$v_j^{n+1} = v_j^n - \frac{\Delta t c(x_j, t_n)}{\Delta x} (v_j^n - v_{j-1}^n) + \Delta t (a(x_j, t_n)v_j^n + b(x_j, t_n))$$

stability implies that $0 \leq c(x_j, t_n) \leq \frac{\Delta x}{\Delta t}$

◆ 2nd order scheme: 2-step Lax-Wendroff

For simple advection equation, $u_t + cu_x = 0$



$$\tilde{v}_{j+1/2} = \frac{1}{2} (v_j^n + v_{j+1}^n) - \frac{c \Delta t}{2 \Delta x} (v_{j+1}^n - v_j^n)$$

$$v_j^{n+1} = v_j^n - \frac{\Delta t c}{\Delta x} (\tilde{v}_{j+1/2} - \tilde{v}_{j-1/2}) \quad \text{2nd order accurate}$$

For the general, linear, non-constant coeff eqn:

$$\tilde{v}_{j+1/2} = \frac{1}{2} (v_j^n + v_{j+1}^n) + \frac{c(x_{j+1/2}, t_n) \Delta t}{2 \Delta x} (v_{j+1}^n - v_j^n) + \frac{\Delta t}{2} (a(x_{j+1/2}, t_n) \frac{1}{2} (v_j^n + v_{j+1}^n) + b(x_{j+1/2}, t_n))$$

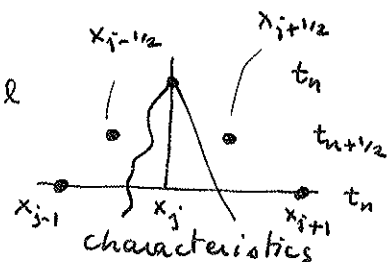
$$v_j^{n+1} = v_j^n - \frac{\Delta t c(x_j, t_{n+1/2})}{\Delta x} (\tilde{v}_{j+1/2} - \tilde{v}_{j-1/2}) + \Delta t (a(x_j, t_{n+1/2}) \frac{1}{2} (\tilde{v}_{j-1/2} + \tilde{v}_{j+1/2}) + b(x_j, t_{n+1/2}))$$

stability implies

must worry
about fastest
wave speed

$$\max_{j,n} |c(x_j, t_n)| \frac{\Delta t}{\Delta x} \leq 1$$

as can be seen by stencil



Systems of Linear Equations

$$\underline{u}_t + \underline{A} \underline{u}_x = 0, \quad \underline{A} \in \mathbb{R}^{m \times m}$$

The system is hyperbolic if \underline{A} is diagonalizable and with real eigenvalues.

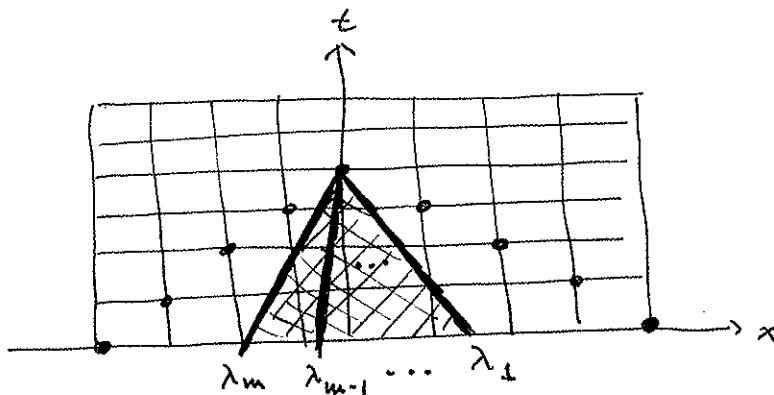
Lax-Friedrichs:

$$\underline{v}_j^{n+1} = \frac{1}{2}(\underline{v}_{j-1}^n + \underline{v}_{j+1}^n) - \underline{A} \frac{\Delta t}{2\Delta x} (\underline{v}_{j+1}^n - \underline{v}_{j-1}^n)$$

Lax-Wendroff:

$$\underline{v}_j^{n+1} = \underline{v}_j^n - \underline{A} \frac{\Delta t}{2\Delta x} (\underline{v}_{j+1}^n - \underline{v}_{j-1}^n) + \underline{A}^2 \left(\frac{\Delta t}{\Delta x}\right)^2 \frac{1}{2} (\underline{v}_{j+1}^n - 2\underline{v}_j^n + \underline{v}_{j-1}^n)$$

CFL stability Analysis



Suppose $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_m$ are the eigenvalues of \underline{A}

$$\max_{1 \leq p \leq m} |\lambda_p| \frac{\Delta t}{\Delta x} \leq 1$$

$$\Rightarrow \Delta t \leq \frac{\Delta x}{\max_{1 \leq p \leq m} |\lambda_p|}$$

Faster the propagation, smaller the time step

Upwind Methods (have less diffusion than Lax-Friedrichs)

$$u_t + \underline{A} u_x = 0,$$

suppose the eigenvalues of \underline{A} are

$$\lambda_1 \leq \dots \leq \lambda_q \leq 0 \leq \lambda_{q+1} \leq \dots \leq \lambda_m$$

the sign of the eigenvalues determines direction of discrete derivatives

$$\underline{R}^{-1} u_t + \underline{R}^{-1} \underline{A} \underline{R} \underline{R}^{-1} u_x = 0, \quad \underline{R}^{-1} u_t + \underline{\Lambda} \underline{R}^{-1} u_x = 0$$

$$\rightarrow \underline{w} = \underline{R}^{-1} u, \quad \rightarrow \underline{w}_t + \underline{\Lambda} \underline{w}_x = 0$$

$$\text{set } \underline{\Lambda} = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_m \end{bmatrix} = \underbrace{\begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_q & 0 & 0 \\ & & 0 & \ddots & 0 \\ & & 0 & & 0 \end{bmatrix}}_{\underline{\Lambda}_-} + \underbrace{\begin{bmatrix} 0 & & \\ & \ddots & \\ & & \lambda_{q+1} & \ddots & \\ & & & \ddots & \lambda_m \end{bmatrix}}_{\underline{\Lambda}_+}$$

$$\rightarrow \underline{w}_t + (\underline{\Lambda}_- + \underline{\Lambda}_+) \underline{w}_x = 0$$

multiply through by \underline{R} to recover u

$$u_t + \underbrace{\underline{R} \underline{\Lambda}_- \underline{R}^{-1}}_{\underline{A}_-} u_x + \underbrace{\underline{R} \underline{\Lambda}_+ \underline{R}^{-1}}_{\underline{A}_+} u_x = 0 \quad \text{with } \underline{A} = \underline{A}_- + \underline{A}_+$$

$\rightarrow u_t + \underline{A}_- u_x + \underline{A}_+ u_x = 0$ decomposed eqn into parts with left running characteristics and right running characteristics

the upwind discretization

$$v_j^{n+1} = v_j^n - \frac{\Delta t}{\Delta x} \underline{A}_- (v_{j+1}^n - v_j^n) - \frac{\Delta t}{\Delta x} \underline{A}_+ (v_j^n - v_{j-1}^n)$$

$$= v_j^n - \frac{\Delta t}{2\Delta x} \underline{A} (v_{j+1}^n - v_{j-1}^n) + \frac{\Delta t}{2\Delta x} |\underline{A}| (v_{j+1}^n - 2v_j^n + v_{j-1}^n)$$

$$\text{where } |\underline{A}| = \underline{A}_+ - \underline{A}_-$$

Boundary Conditions For Linear Hyperbolic Systems

$$\underline{u}_t + \underline{A} \underline{u}_x = 0, \quad 0 < x < 1, \quad t > 0, \quad \underline{u}(x, 0) = \underline{f}(x)$$

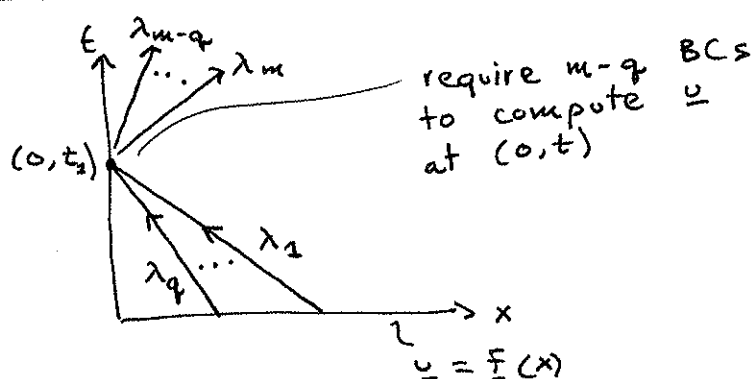
What boundary conditions are allowed?

We must consider the characteristics:

\underline{A} has eigenvalues $\lambda_1, \dots, \lambda_m$

where $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_q < 0 < \lambda_{q+1} \leq \dots \leq \lambda_m$

Near the boundary $x=0$:



the characteristics $\lambda_1, \dots, \lambda_q \rightarrow$ outflow characteristics

the characteristics $\lambda_{q+1}, \dots, \lambda_m$ are inflow characteristics

Rule: Require one boundary condition for each inflow characteristic.

The question is, what do we specify?
The general form of the $m-q$ BCs would be

$$\underline{B} \underline{u}(0, t) = \underline{g}(t)$$

where $\underline{B} = \underbrace{\begin{bmatrix} -\underline{b}_1^T & - \\ \vdots & \\ -\underline{b}_{m-q}^T & - \end{bmatrix}}_m \}^{m-q}, \quad \underline{g}(t) = \left[\begin{array}{c} \\ \\ \end{array} \right] \}^{m-q}$

What are the row vectors $\underline{b}_1, \dots, \underline{b}_{m-q}$?

The characteristic form of the equation, $\underline{w}(x,t) = \underline{R}^{-1} \underline{u}(x,t)$

$$\rightarrow \underline{w}_t + \underline{A} \underline{R} \underline{w}_x = 0, \quad \underline{A} \underline{R} = \underline{R} \underline{A}$$

The boundary conditions for \underline{w} :

$$\underline{B} \underline{R} \underline{w}(0,t) = \underline{g}(t)$$

$$\text{let } \underline{\tilde{B}} = \underline{B} \underline{R} = \begin{bmatrix} \underline{b}_1^T & \dots & \underline{b}_{m-q}^T \\ \vdots & & \vdots \\ \underline{b}_{m-q}^T & \dots & \underline{b}_m^T \end{bmatrix} \begin{bmatrix} \underline{r}_1 & \dots & \underline{r}_m \\ \vdots & & \vdots \\ \underline{r}_1 & \dots & \underline{r}_m \end{bmatrix} \quad \begin{matrix} (m-q, m) \times (m, m) \\ \parallel \\ (m-q, m) \end{matrix}$$

$$\rightarrow \underline{\tilde{B}}_{ij} = \begin{bmatrix} \underline{b}_i^T & \underline{r}_j \end{bmatrix}, \quad i = 1, \dots, m-q, \quad j = 1, \dots, m$$

$$\Rightarrow \underline{\tilde{B}} \underline{w}(0,t) = \underline{g}(t) \Rightarrow \begin{bmatrix} \underline{\tilde{B}} \end{bmatrix} \begin{bmatrix} \underline{w}(0,t) \end{bmatrix} = \begin{bmatrix} \underline{g}(t) \end{bmatrix}$$

We need to partition the matrix:

$$\begin{bmatrix} \underline{\tilde{B}}_1 & \underline{\tilde{B}}_2 \end{bmatrix} \begin{bmatrix} \underline{w}_1 \\ \vdots \\ \underline{w}_2 \end{bmatrix} = \begin{bmatrix} \underline{g}(t) \end{bmatrix}$$

data from ICs

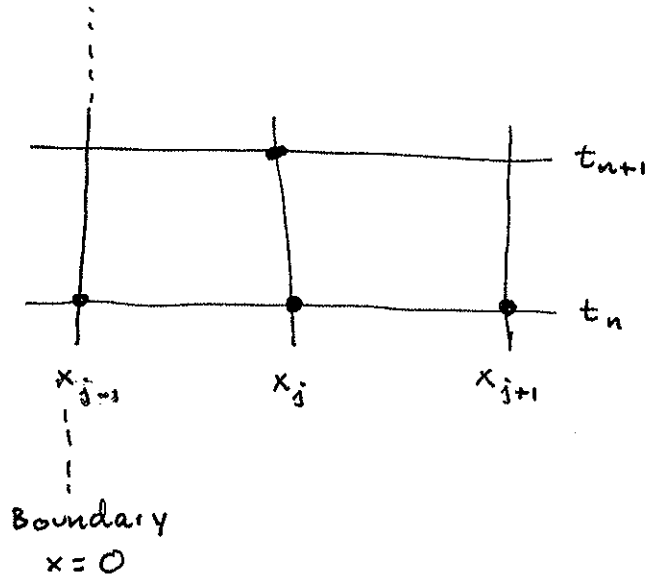
need to be specified by B.C.s

$$\underline{\tilde{B}}_2 \underline{w}_2(0,t) = \underline{g}(t) - \underline{\tilde{B}}_1 \underline{w}_1(0,t)$$

therefore we need that $\underline{\tilde{B}}_2$ be nonsingular.

From a numerical point of view:

Many schemes, such as Lax-Friedrichs and Lax-Wendroff have the following stencils



At $x=0$, require m -component for v_{j-1}^n , but only $m-q$ corresponding to $m-q$ inflow characteristics would be specified by the boundary conditions.

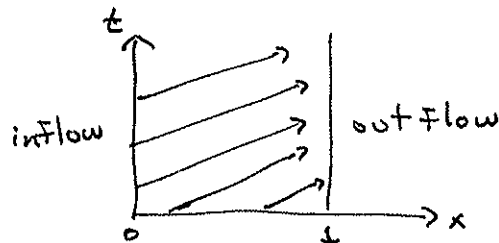
To complete the remaining q components of v_{j-1}^n , use extrapolation from the interior.

Example:

$$v_t + v_x = 0, \quad 0 < x < 1, \quad t > 0$$

$$v(x, 0) = f(x)$$

$$v(0, t) = g(t) \quad (\text{inflow boundary})$$



numerical approximation:

$$x_j = j \Delta x, \quad \Delta x = \frac{1}{N}, \quad t_n = n \Delta t$$

$$v_j^{n+1} = \frac{1}{2}(v_{j-1}^n + v_{j+1}^n) - \frac{\Delta t}{\Delta x}(v_{j+1}^n - v_{j-1}^n) \quad \leftarrow j=1, \dots, N-1 \quad \text{Lax-Friedrichs}$$

$$v_j^0 = f(x_j), \quad 0 \leq j \leq N$$

$$v_0^n = g(t_n), \quad n > 0$$

$$v_N^n = \text{extrapolation, simplest is } v_N^n = v_{N-1}^n$$

Hyperbolic Conservation Laws

Basic equation:
$$\underline{u}_t + \underline{f}(\underline{u})_x = 0$$

conservation form

where $\underline{u}(x,t)$ = m "state" variables

$\underline{f}(\underline{u}(x,t))$ = m "Flux" Functions

This general form is derived from an integral conservation:
integrate differential form from $x=a$ to $x=b$

$$\int_a^b (\underline{u}_t + \underline{f}(\underline{u})_x) dx = 0$$

$$\rightarrow \int_a^b \underline{u}_t dx + \underline{f}(\underline{u}(x,t)) \Big|_{x=a}^{x=b} = 0$$

$$\rightarrow \frac{d}{dt} \int_a^b \underline{u}(x,t) dx + \underline{f}(\underline{u}(x,t)) \Big|_{x=a}^{x=b} = 0$$

integral form

The rate of change of \underline{u} between $x=a$ and $x=b$ is balanced by the net flux of \underline{u} at the boundaries

$$\text{If } \left[\underline{f}(\underline{u}) \right]_a^b = 0 \text{ then } \int_a^b \underline{u}(x,t) dx = \text{constant}$$

and \underline{u} is a vector of conserved quantities.

$$\underline{u}_t + \underline{F}_{\underline{u}}(\underline{u}) \underline{u}_x = 0$$

, quasilinear form

where $\underline{F}_{\underline{u}}$ is a Jacobian matrix. If this Jacobian matrix is diagonalizable with real eigenvalues for any \underline{u} then this equation is hyperbolic.

Examples

• scalar equations

◆ let $F = cv$, $c = \text{const} \rightarrow v_t + (cv)_x = 0$, conservation form
 $v_t + cv_x = 0$, quasilinear form

• let $F = \frac{1}{2}v^2$ (nonlinear) Burgers Inviscid Equation

$$v_t + \left(\frac{1}{2}v^2\right)_x = 0, \text{ conservation form}$$

$$v_t + vv_x = 0, \text{ quasilinear form}$$

• systems

◆ $\underline{F} = \underline{A} \underline{v}$, \underline{A} real eigenvalues

$$\underline{v}_t + (\underline{A} \underline{v}_x) = 0, \text{ conservation form}$$

$$\underline{v}_t + \underline{A} \underline{v}_x = 0, \text{ quasilinear form}$$

◆ $\underline{F} = \begin{bmatrix} v_2 \\ v_2^2/v_1 + a^2 v_1 \end{bmatrix}$, $a > 0$, constant, $\underline{v} = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$

$$\underline{v}_t + \underline{F}(\underline{v})_x = 0$$

$$\underline{v}_t + \underline{F}_{\underline{v}}(\underline{v}) \underline{v}_x = 0$$

v_1 is a density $\Rightarrow v_1 > 0$

v_2 is a momentum

a is a sound speed

calculate Jacobian

$$\underline{F}_{\underline{v}} = \begin{bmatrix} \frac{\partial F_1}{\partial v_1} & \frac{\partial F_1}{\partial v_2} \\ \frac{\partial F_2}{\partial v_1} & \frac{\partial F_2}{\partial v_2} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -\frac{v_2^2}{v_1^2} + a^2 & \frac{2v_2}{v_1} \end{bmatrix}$$

Acoustic waves

$$\det \begin{bmatrix} -\lambda & 1 \\ a^2 - \frac{v_2^2}{v_1^2} & \frac{2v_2}{v_1} - \lambda \end{bmatrix} = \lambda \left(\lambda - \frac{2v_2}{v_1} \right) - \left(a^2 - \frac{v_2^2}{v_1^2} \right)$$

$$= \lambda^2 - \frac{2v_2}{v_1} \lambda - \left(a^2 - \frac{v_2^2}{v_1^2} \right) = 0$$

$$\rightarrow \boxed{\lambda = \frac{v_2}{v_1} \pm a}, \text{ by quadratic equation}$$

$$\Rightarrow \lambda = v \pm a$$

Scalar Conservation Laws

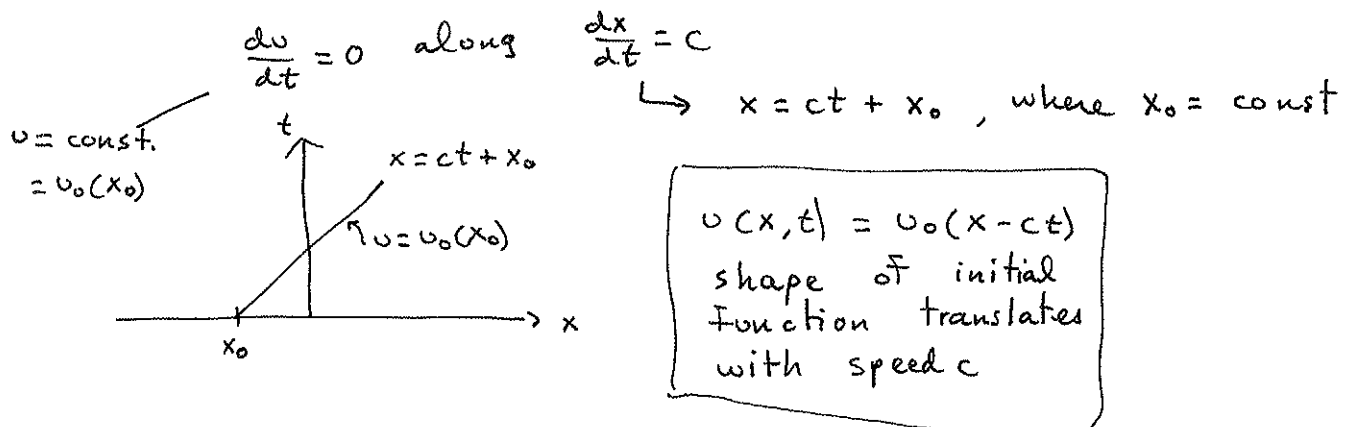
Begin with the scalar case: (pure IVP)

$$u_t + F(u)_x = 0, \quad |x| < \infty, \quad t > 0, \quad u(x, 0) = u_0(x)$$

Consider case $F(u) = cu$, $c = \text{const.}$

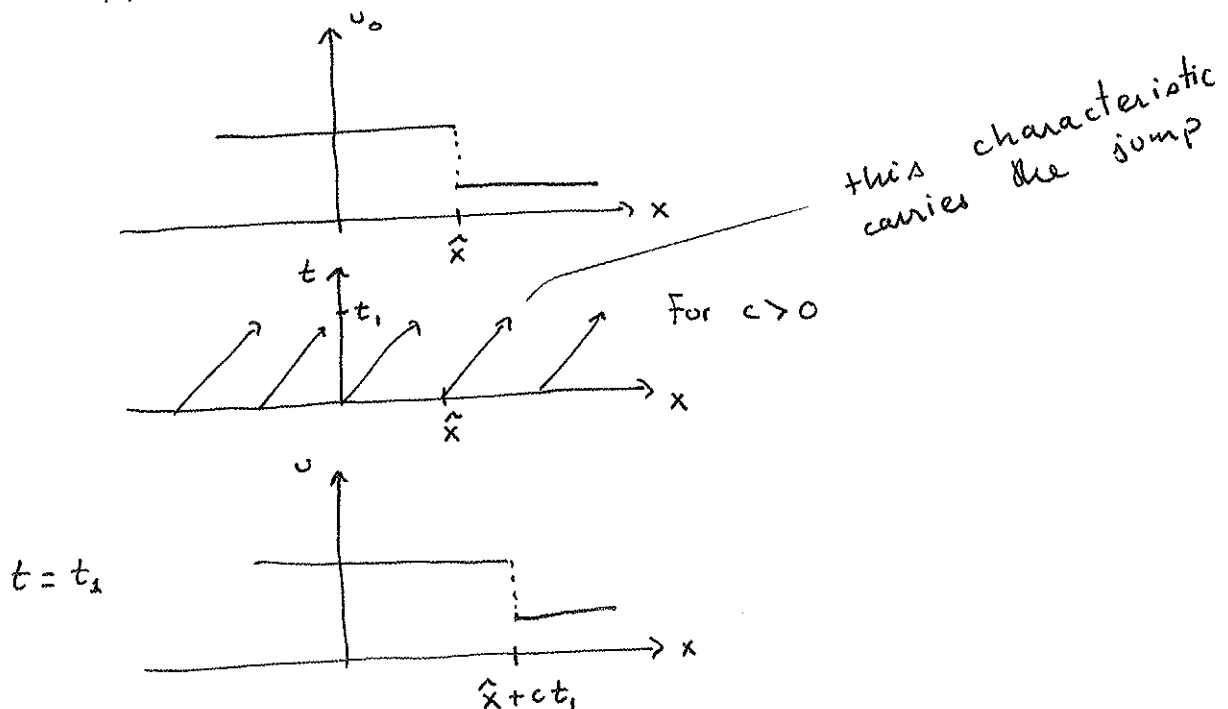
$$\rightarrow u_t + cu_x = 0, \quad u(x, 0) = u_0(x)$$

characteristic description:



of particular interest are solutions with discontinuities

Suppose $u_0(x)$ has a jump discontinuity at \hat{x} .



Need to interpret discontinuities in the PDE.

1) zero limit of a "viscous" PDE

$$\tilde{u}_t + c\tilde{u}_x = \nu \tilde{u}_{xx}, \quad \nu = \text{const} > 0$$

$$\tilde{u}(x, 0) = u_0(x)$$

interpret solution of $u_t + cu_x = 0$, $u(x, 0) = u_0(x)$, the inviscid problem as the limit as $\nu \rightarrow 0$.

change of variables: let $\xi = x - ct$, $\tilde{t} = t$

$$\Rightarrow \tilde{u}_{\tilde{t}} = \nu \tilde{u}_{\xi\xi}, \quad \tilde{u}(\xi, 0) = u_0(\xi) \quad \text{heat equation}$$

solve using Fourier transforms

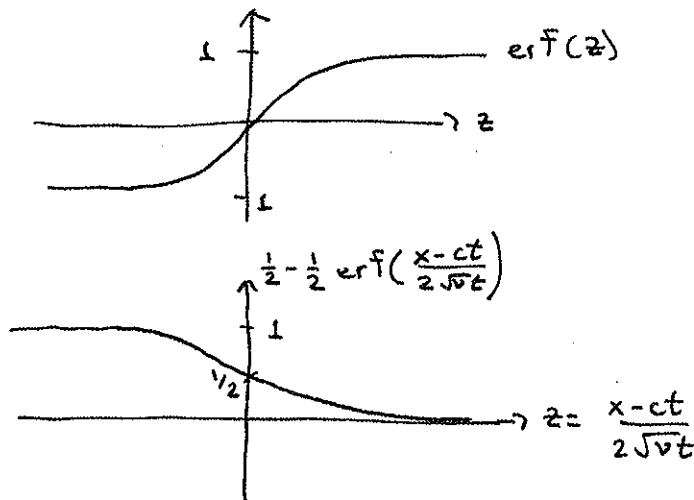
$$\hat{u}(\xi, t) = \frac{1}{\sqrt{4\pi\nu t}} \int_{-\infty}^{\infty} u_0(s) e^{-\frac{(s-\xi)^2}{4\nu t}} ds$$

$$\text{If } u_0(x) = \begin{cases} 1, & x < 0 \\ 0, & x > 0 \end{cases} \rightarrow \hat{u}(\xi, t) = \frac{1}{\sqrt{4\pi\nu t}} \int_{-\infty}^0 e^{-\frac{(s-\xi)^2}{4\nu t}} ds$$

$$\rightarrow \hat{u}(\xi, t) = \frac{1}{2} - \frac{1}{2} \text{erf}\left(\frac{\xi}{2\sqrt{\nu t}}\right)$$

$$\rightarrow \tilde{u}(x, t) = \frac{1}{2} - \frac{1}{2} \text{erf}\left(\frac{x-ct}{2\sqrt{\nu t}}\right) \quad \text{where}$$

$$\boxed{\text{erf } z = \frac{2}{\sqrt{\pi}} \int_0^z e^{-s^2} ds}$$



transition width decreases as ν gets smaller and in the limit as $\nu \rightarrow 0$, becomes a jump.

2) Weak solutions of the integral conservation

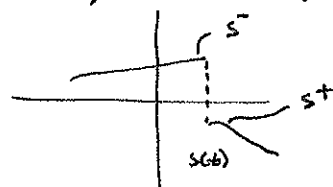
$$\frac{d}{dt} \int_a^b u dx + F(u) \Big|_a^b = 0$$

suppose $u(x,t)$ has a jump discontinuity at $x=s(t)$
 Position a, b on either side of $x=s(t)$.

$$\rightarrow \frac{d}{dt} \int_a^b u dx = \frac{d}{dt} \int_a^{s(t)} u dx + \frac{d}{dt} \int_{s(t)}^b u dx$$

then u is smooth for each interval $a < x < s(t)$, $s(t) < x < b$.
 so move derivative inside

$$\rightarrow \frac{d}{dt} \int_a^b u dx = \int_a^{s(t)} u_t dx + \int_{s(t)}^b u_t dx + u(s^-, t) \frac{ds}{dt} - u(s^+, t) \frac{ds}{dt}$$



substitute $u_t = -F_x(u)$

$$\begin{aligned} \rightarrow \frac{d}{dt} \int_a^b u dx &= - \int_a^{s(t)} F_x(u) dx - \int_{s(t)}^b F_x(u) dx + (u(s^-, t) - u(s^+, t)) \frac{ds}{dt} \\ &= - F(u) \Big|_a^{s(t)} - F(u) \Big|_{s(t)}^b + (u(s^-, t) - u(s^+, t)) \frac{ds}{dt} \end{aligned}$$

$$\begin{aligned} \frac{d}{dt} \int_a^b u dx &= - F(u(s^-, t)) + F(u(a, t)) + F(u(s^+, t)) - F(u(b, t)) \\ &\quad + (u(s^-, t) - u(s^+, t)) \frac{ds}{dt} \end{aligned}$$

Define jump notation:

$$[u] \equiv u(s^+, t) - u(s^-, t)$$

$$[F] \equiv F(u(s^+, t)) - F(u(s^-, t))$$

$$\rightarrow \frac{d}{dt} \int_a^b u \, dx = [f] - f(u) \Big|_a^b - [u] \frac{ds}{dt}$$

$$\text{recall, } \frac{d}{dt} \int_a^b u \, dx + f(u) \Big|_a^b = 0$$

$$\Rightarrow \boxed{[u] \frac{ds}{dt} = [f]} \quad \text{Rankine-Hugoniot jump conditions}$$

Let us apply this to the linear case,

$$u_t + cu_x = 0$$

$$\text{where } f = cu$$

$$\rightarrow [u] \frac{ds}{dt} = [cu] \rightarrow [u] \frac{ds}{dt} = c[u] \rightarrow \frac{ds}{dt} = c$$

weak form: use PDE whenever solution is smooth
and patch u across any jump using
jump conditions

Burgers Equation

(simplest nonlinearity - quadratic)

$$f(u) = \frac{1}{2} u^2$$

$$u_t + \left(\frac{1}{2} u^2\right)_x = 0, \quad u(x, 0) = u_0(x) \text{ where } u_0 \text{ is smooth}$$

$$u_t + f(u)_x = 0$$

$$u_t + f'(u) u_x = 0, \quad c(u) = f'(u)$$

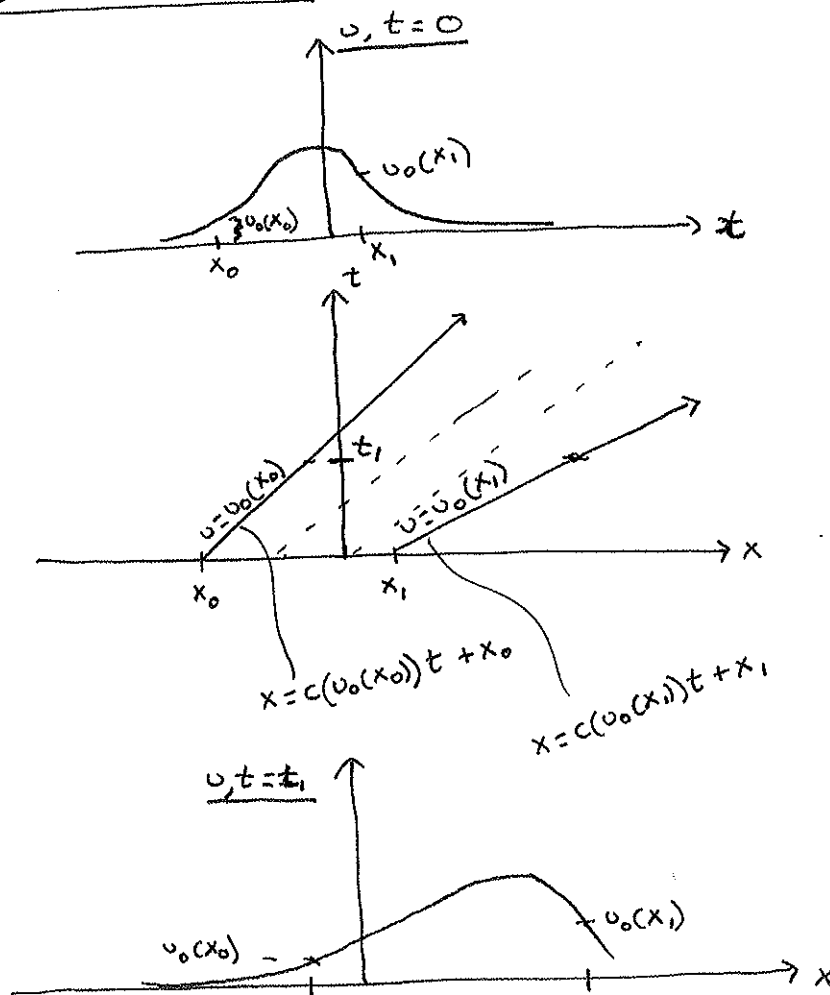
characteristics: $\frac{du}{dt} = 0$ along $\frac{dx}{dt} = c(u)$

→ $u = \text{const.}$ along a characteristic

→ $c(u)$ is a constant along a characteristic

→ $x = c(u)t + \text{const.}$

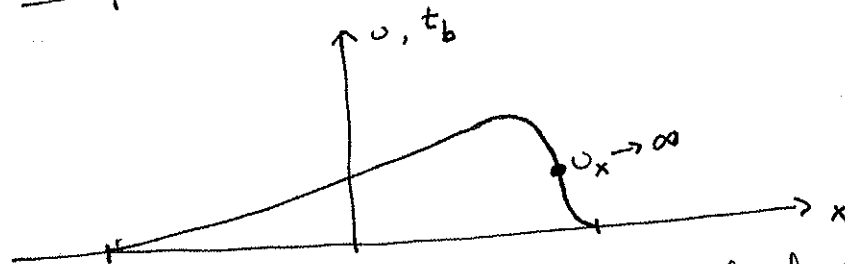
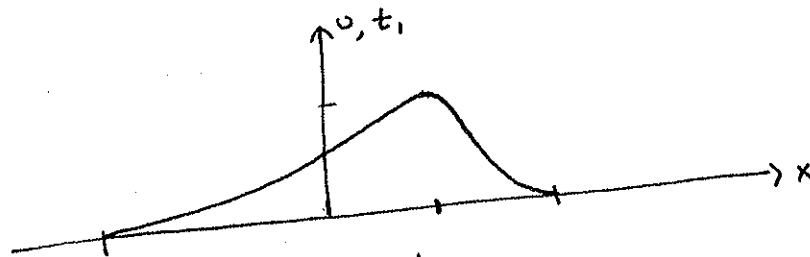
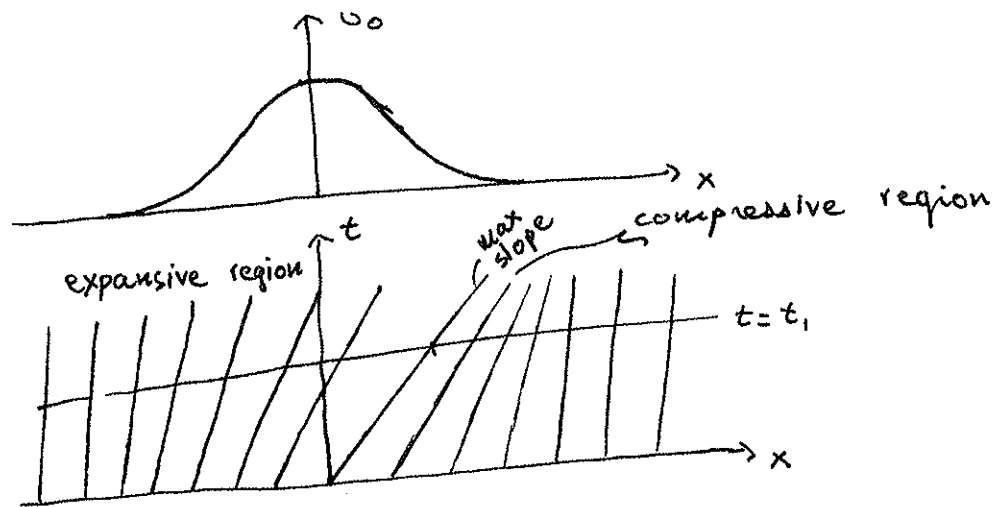
Graphical construction:



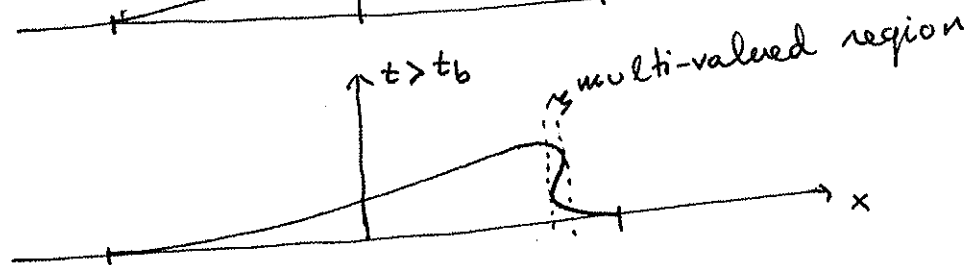
initial shape
given by $u_0(x)$
translates and
distorts

Burger's eqn: $F(u) = \frac{1}{2}u^2$

$$c(u) = F'(u) = u$$



t_b - time of
First breaking,
when characteristics
First cross

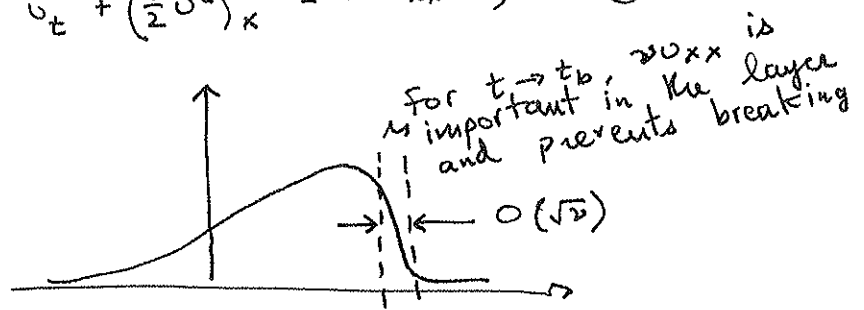


this solution satisfies the PDE for $t < t_b$
At $t = t_b$, a singularity in the PDE occurs,
so this solution is no longer acceptable.

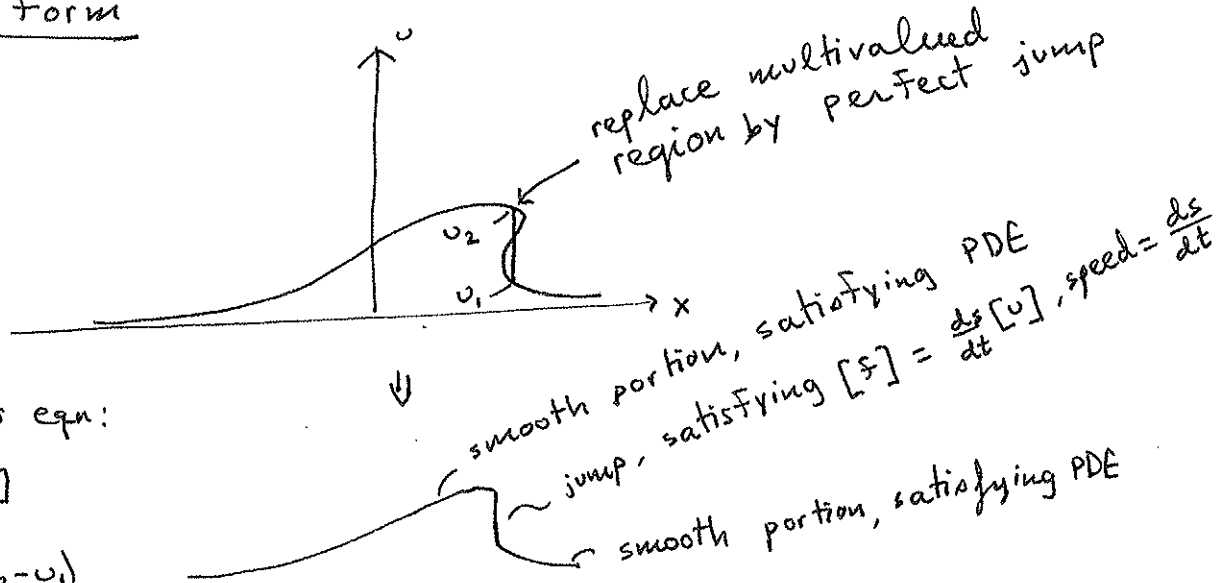
we need to recover the solution

1) viscous limit, $u_t + \left(\frac{1}{2}u^2\right)_x = \nu u_{xx}$, $\nu \rightarrow 0$

near $t = t_b$;



2) weak Form



For Burgers eqn:

$$[f] = \frac{ds}{dt} [u]$$

$$\frac{1}{2}(u_2^2 - u_1^2) = \frac{ds}{dt}(u_2 - u_1)$$

the solution is constructed by patching smooth bits where the PDE holds together with jumps where the jump condition holds.

The issue with the weak solution is that it is NOT UNIQUE.

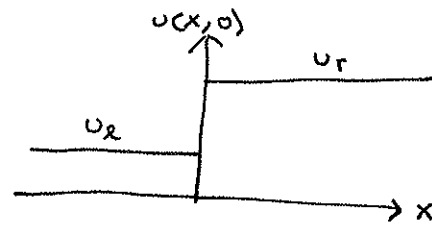
We can force uniqueness by insisting that the weak form is the $\nu \rightarrow 0$ limit of the viscous problem

→ entropy satisfying solutions

Example: "Riemann" problem

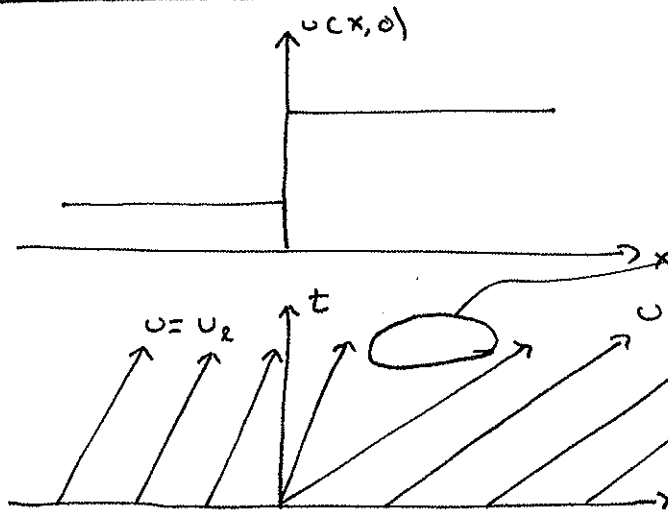
$$u_t + \left(\frac{1}{2}u^2\right)_x = 0, \quad |x| < \infty, \quad t > 0$$

$$u(x,0) = \begin{cases} u_l, & x < 0 \\ u_r, & x > 0 \end{cases}$$



- There are 2 cases $u_l < u_r$ and $u_l > u_r$

case ① $u_l < u_r$



what is solution here?

similarity form: let $u(x,t) = w(\eta)$, $\eta = \frac{x}{t}$

$$\rightarrow u_t = w' \cdot \eta_t = -w' \frac{x}{t^2} = -\eta w'(\eta) \frac{1}{t}$$

$$u_x = w'(\eta) \eta_x = w'(\eta) \cdot \frac{1}{t}$$

$$\rightarrow u_t + u \cdot u_x = 0 \rightarrow -\eta w'(\eta) \cdot \frac{1}{t} + w \cdot w'(\eta) \frac{1}{t} = 0$$

$$\rightarrow (w - \eta) w'(\eta) = 0$$

$$\rightarrow w'(\eta) = 0 \rightarrow w(\eta) = \text{const} \quad \text{or} \quad w = \eta$$

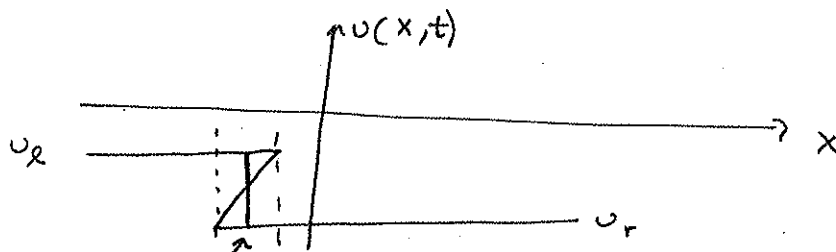
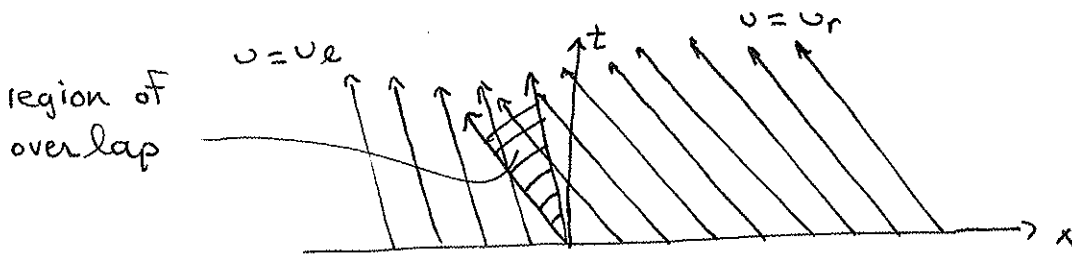
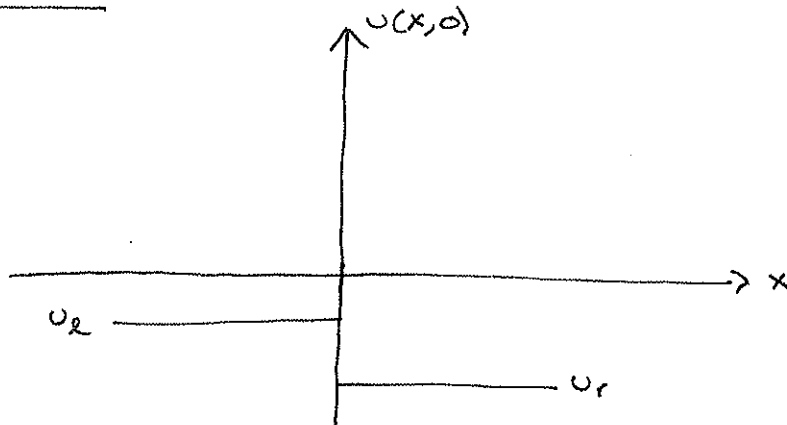
$$\downarrow \quad \quad \quad \downarrow$$

$$u = \text{const} \quad \quad \quad u = \frac{x}{t}$$

in solution, notice that $u = \text{const}$ outside expansive region and inside, we apply $u = \frac{x}{t}$

The solution is
$$u(x,t) = \begin{cases} u_l & , \frac{x}{t} < u_l \\ \frac{x}{t} & , u_l < \frac{x}{t} < u_r \\ u_r & , \frac{x}{t} > u_r \end{cases}$$

Case (2) $u_l > u_r$



replace
multivalued
region by
a jump

that propagates
at speed $\frac{ds}{dt}$,

$[F] = \frac{ds}{dt} [u]$, $s(t)$ - position
of jump

For Burger's equation

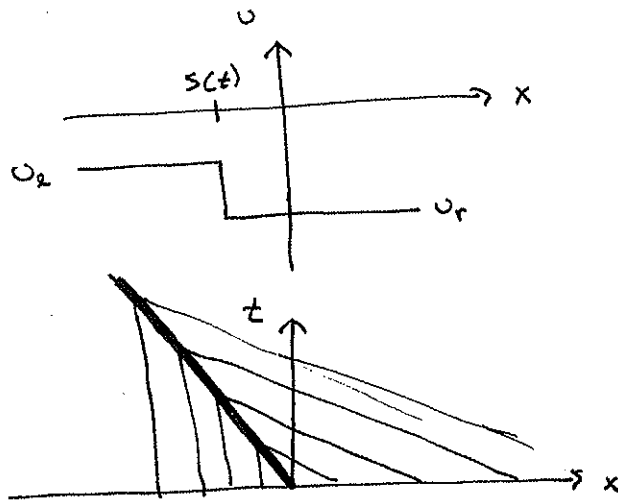
$$[F(u)] = [u] \frac{ds}{dt}$$

$$\frac{1}{2}(u_r^2 - u_l^2) = (u_r - u_l) \frac{ds}{dt}$$

$$\rightarrow \boxed{\frac{ds}{dt} = \frac{u_r + u_l}{2}} \rightarrow s(t) = \left(\frac{u_r + u_l}{2}\right)t + \text{const}$$

$$\downarrow$$

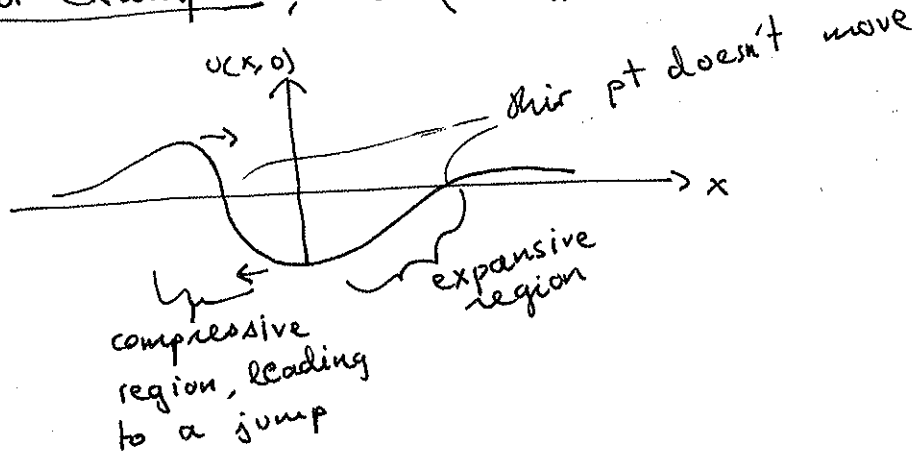
$$\boxed{s(t) = \left(\frac{u_r + u_l}{2}\right)t}$$



$$s(t) = \left(\frac{u_L + u_R}{2} \right) t$$

$$u(x,t) = \begin{cases} u_L, & x < \frac{1}{2}(u_L + u_R)t \\ u_R, & x > \frac{1}{2}(u_L + u_R)t \end{cases}$$

For example, $u_t + \left(\frac{1}{2} u^2 \right)_x = 0$



Numerics (For scalar case)

$$u_t + F(u)_x = 0, \quad |x| < \infty, \quad t > 0$$

$$u(x,0) = u_0(x)$$

Issues

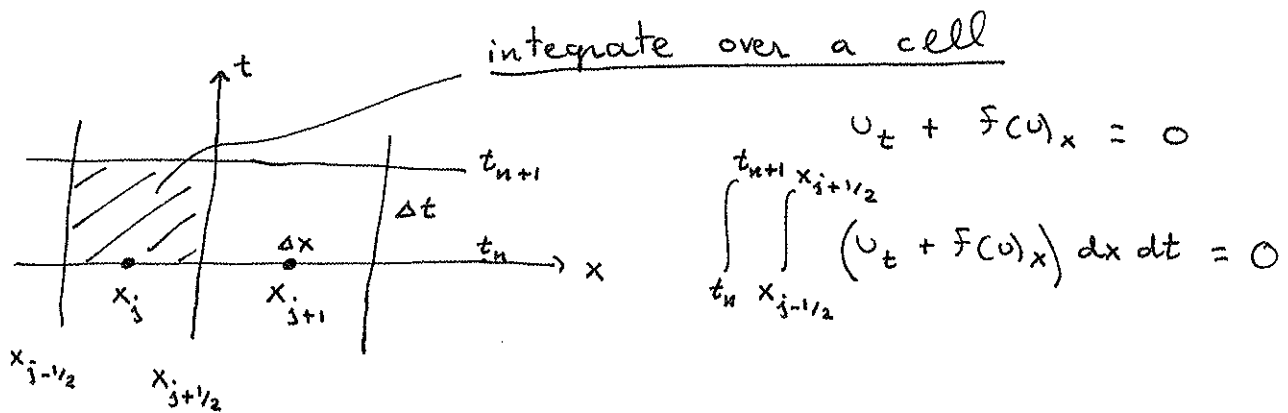
- (1) If solution is smooth, then the nonlinearity does not play a significant role and a numerical approximation should behave as if the equation is linear.
- (2) Nonlinearity can lead to jump discontinuities forming in finite time. We need to worry about how numerical approximation behaves near jumps and what equation is being approximated in the vicinity of a jump.

3) Uniqueness of weak solutions and nonlinear stability.

Finite Volume Formulation

\Rightarrow try to maintain the integral Formulation (integral conservation) in the discrete approximation

\Rightarrow to obtain correct weak Form



$$\rightarrow \left[\int_{x_{j-1/2}}^{x_{j+1/2}} u dx \right]_{t_n}^{t_{n+1}} + \left[\int_{t_n}^{t_{n+1}} F(u) dt \right]_{x_{j-1/2}}^{x_{j+1/2}} = 0 \quad \text{every thing is exact}$$

Define $U_j^n \equiv \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} u(x, t_n) dx$ "cell average"

$$F_{j+1/2}^n \equiv \frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} F(u(x_{j+1/2}, t)) dt$$

$$\rightarrow \Delta x U_j^{n+1} - \Delta x U_j^n + \Delta t F_{j+1/2}^n - \Delta t F_{j-1/2}^n = 0$$

$$\rightarrow U_j^{n+1} = U_j^n - \frac{\Delta t}{\Delta x} (F_{j+1/2}^n - F_{j-1/2}^n) \quad \text{still exact}$$

Numerical approximation comes from the approximations of $F_{j+1/2}^n$ and $F_{j-1/2}^n$ ie numerical Flux Functions

typically, $F_{j+1/2}^n \approx \underset{\substack{\uparrow \\ \text{numerical Flux Function}}}{G}(U_j^n, U_{j+1}^n)$, $F_{j-1/2}^n \approx G(U_{j-1}^n, U_j^n)$

comment, there are other possibilities to define F , eg $F_{j+1/2}^n = \tilde{G}(U_{j-p}^n, \dots, U_{j+q}^n)$

Conservative Finite Volume Scheme

$$V_j^{n+1} = V_j^n - \frac{\Delta t}{\Delta x} \left(G(V_j^n, V_{j+1}^n) - G(V_{j-1}^n, V_j^n) \right)$$

where $V_j^n \approx \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} u(x, t_n) dx$

and $G(V_j^n, V_{j+1}^n) \approx \frac{1}{\Delta t} \int_{t_n}^{t_{n+1}} \tilde{f}(u(x_{j+1/2}, t)) dt$

note, $V_j^0 = \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} u_0(x) dx = u_0(x_j) + O(\Delta x^2)$

Examples

~~$G(u_l, u_r)$~~ $G(u_l, u_r) = \frac{1}{2} \left(f(u_l) + f(u_r) \right) - \frac{\Delta x}{2\Delta t} (u_r - u_l)$ Lax-Friedrichs (1st order)

$G(u_l, u_r) = \frac{1}{2} \left(f(u_l) + f(u_r) \right) - \frac{\Delta t}{2\Delta x} f'(\bar{u}) (f(u_r) - f(u_l))$ Lax-Wendroff (2nd order)

where $\bar{u} = \frac{1}{2} (u_l + u_r)$

Comments

- There are many more choices for G , other than Lax-Friedrichs and Lax-Wendroff.

- Not free to make any choice for $G(u_r, u_l)$.

Consistency requires that

$$G(u, u) = f(u)$$

and smoothness of function G .

Consider the truncation error:

$$\tau_j^n = \frac{u_j^{n+1} - u_j^n}{\Delta t} + \frac{G(u_j^n, u_{j+1}^n) - G(u_{j-1}^n, u_j^n)}{\Delta x}$$

consider function with two arguments $G(u, v)$

$$\tau_j^n = u_t + O(\Delta t) + \frac{G(u, u + \Delta x u_x + \dots) - G(u - \Delta x u_x + \dots, u)}{\Delta x}$$

$$\tau_j^n = u_t + O(\Delta t) + \frac{\cancel{G(u, u)} + G_u(u, u) \Delta x u_x + O(\Delta x^2)}{\Delta x} - \frac{\cancel{G(u, u)} - G_v(u, u) \Delta x u_x + O(\Delta x^2)}{\Delta x}$$

$$\tau_j^n = u_t + O(\Delta t) + (G_u(u, u) + G_v(u, u)) u_x + O(\Delta x)$$

we want $f'(u) = (G_u(u, u) + G_v(u, u))$ for consistency

this is true if $f(u) = G(u, u)$.

- The method is conservative, which means that if jump discontinuities develop, then they will propagate with the correct speed as determined by the integral conservation $\Rightarrow [f(u)] = [u] \frac{ds}{dt}$

Example

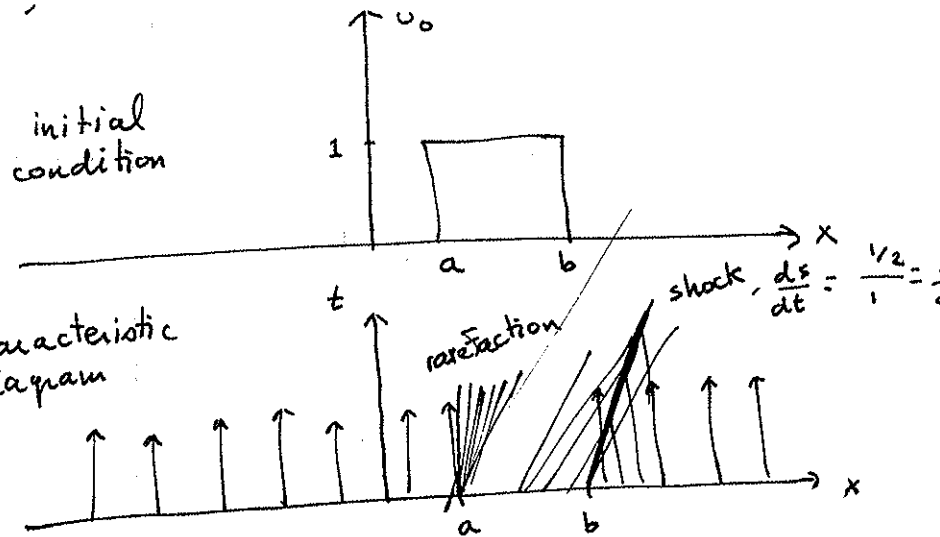
$$u_t + f(u)_x = 0, \quad |x| < \infty, \quad t > 0, \quad u(x, 0) = u_0(x)$$

Solve using

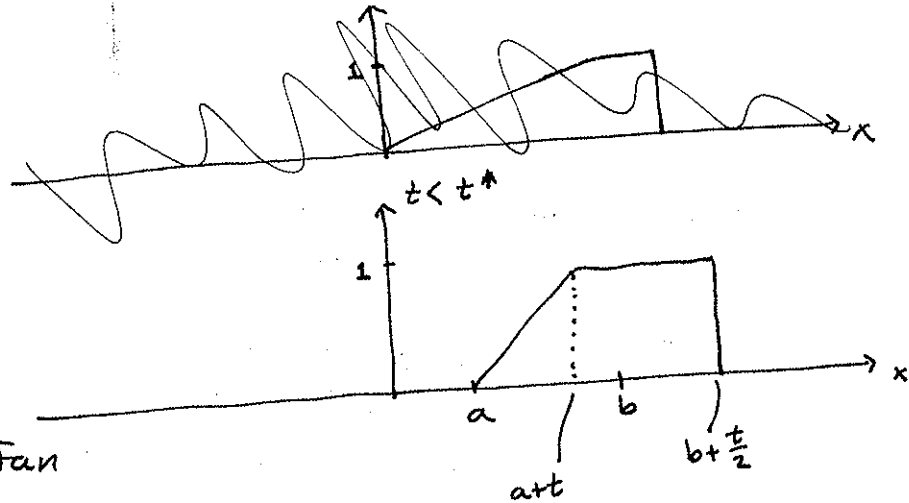
$$v_j^{n+1} = v_j^n - \frac{\Delta t}{\Delta x} (G(v_j^n, v_{j+1}^n) - G(v_{j-1}^n, v_j^n))$$

using both Lax-Friedrichs and Lax-Wendroff for G .

Assume $f(u) = \frac{1}{2} u^2$, inviscid Burgers eqn, and that initial state is

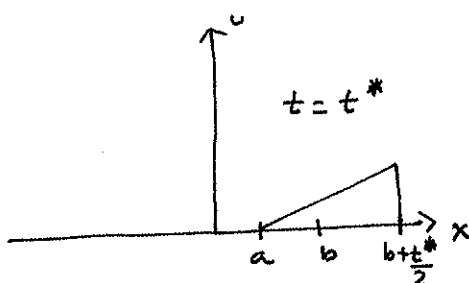


$$\Rightarrow u(x, t) = \begin{cases} 0, & x < a \\ \frac{x-a}{t}, & a < x < a+t \\ 1, & a+t < x < b+\frac{t}{2} \\ 0, & x > b+\frac{t}{2} \end{cases}$$



Notice, the expansion Fan meets the shock when

$$a+t = b+\frac{t}{2} \rightarrow t = 2(b-a) = t^*$$



what about after t^* ? \Rightarrow



After t^* , $\frac{ds}{dt} = \frac{f(\hat{u}) - f(0)}{\hat{u} - 0}$

where \hat{u} is the solution in the expansion fan at $x = s(t)$, ie at the shock

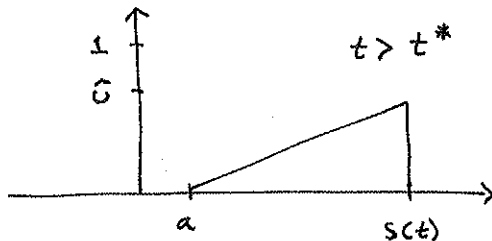
$f(0) = 0, \Rightarrow \frac{ds}{dt} = \frac{\frac{1}{2} \hat{u}^2}{\hat{u}} = \frac{1}{2} \hat{u} = \frac{1}{2} \left(\frac{s-a}{t} \right)$

solve differential equation for s , with initial condition

$s(t^*) = b + \frac{t^*}{2}$

$\Rightarrow s(t) = \sqrt{2t(b-a)} + a, t \geq t^*$

$\hat{u} = \sqrt{\frac{2(b-a)}{t}}$



Test the numerical scheme using following parameters:

$x_j = (j - 1/2) \Delta x, \Delta x = \frac{1}{N}$

$v_j^0 = \frac{1}{2} \tanh(\lambda(x_j - a)) + \frac{1}{2} \tanh(\lambda(b - x_j))$

$a = 0.1, b = 0.4, \lambda = 100, N = 200$

Compare with non-conservative scheme

$v_t + uv_x = 0, \text{ for } u \geq 0$

$v_j^{n+1} = v_j^n - \frac{\Delta t}{\Delta x} v_j^n (v_j^n - v_{j-1}^n)$

results in completely wrong shock location

Lax-Wendroff Theorem (LeVeque)

Consider a sequence of grids indexed by $k=1, 2, \dots$ with Δx_k and Δt_k vanishing as $k \rightarrow \infty$. Let $v_k(x, t)$ be a piecewise constant function taking the value v_j^n when $x \in (x_{j-1/2}, x_{j+1/2}]$, $t \in [t_n, t_{n+1})$ on grid k , where v_j^n is obtained from a consistent, conservative scheme. IF v_k converges to a function $u(x, t)$ as $k \rightarrow \infty$, then u is a weak solution of the conservation law.

Remember, weak solutions aren't necessarily unique!

Remarks

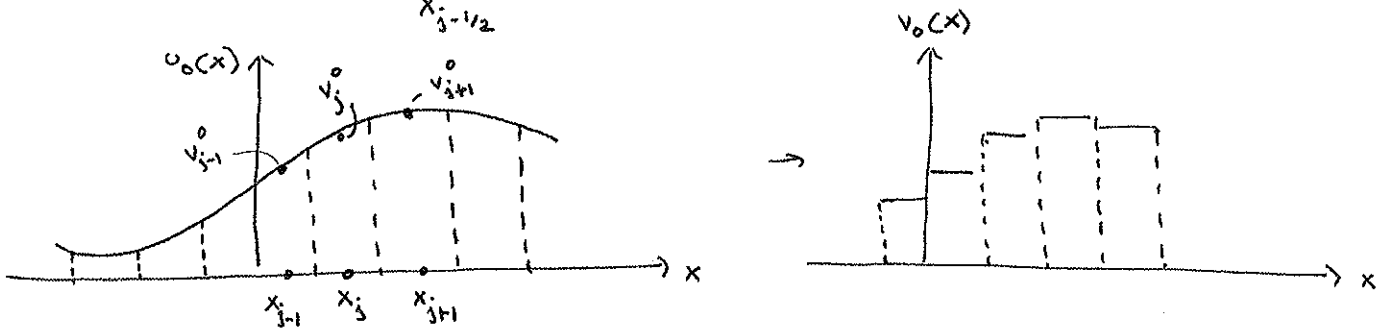
- 1) The Lax-Wendroff theorem does not guarantee convergence. It applies if the scheme converges. (Require nonlinear stability analysis. IF solution is smooth, stability analysis to linearized solution applies, but when jumps occur, this no longer applies.)
- 2) If convergence occurs, then discrete solution converges to weak solution, but the weak solutions are not unique! Require solutions that are entropy satisfying...

Godunov Methods

Essentially a nonlinear version of upwind methods. Numerical approach is still conservative and based on solutions of Riemann problems.

$$u_t + f(u)_x = 0, \quad u(x, 0) = u_0(x)$$

Define $v_j^0 = \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} u_0(x) dx$ = cell average of initial data

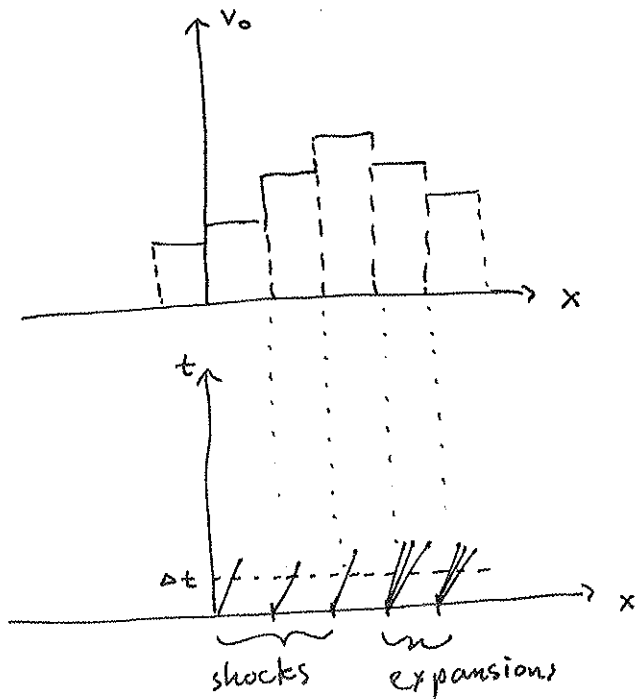


let $v_0(x)$ be piecewise constant, $= v_j^0$ for $x \in [x_{j-1/2}, x_{j+1/2}]$

Consider the exact solution of the problem

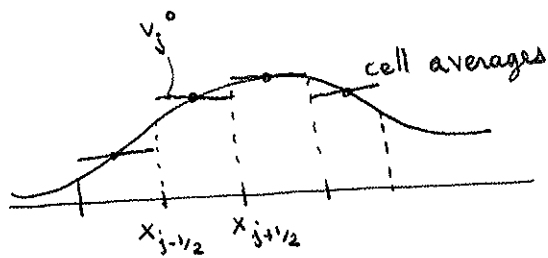
$$u_t + f(u)_x = 0$$

$u(x, 0) = v_0(x)$ - piecewise constant approximation of u_0



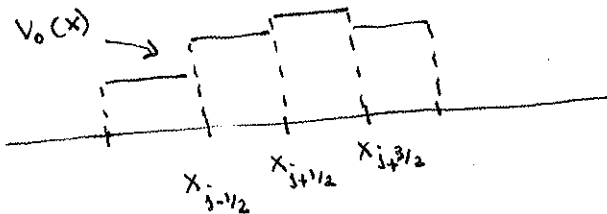
Godunov - reconstruct solution at Δt exactly based on $u(x, 0) = v_0(x)$

cell average:
$$v_j^0 = \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} v_0(x) dx$$



From initial cell averages, construct a piecewise constant function $v_0(x)$ st

$$v_0(x) = v_j^0 \text{ for } x \in (x_{j-1/2}, x_{j+1/2})$$



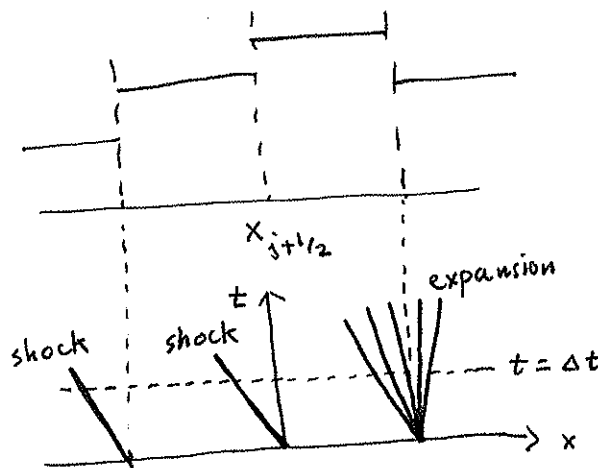
← solve the IVP:

$$u_t + f(u)_x = 0, \quad |x| < \infty, \quad 0 < t < \Delta t$$

$$u(x, 0) = v_0(x)$$

so instead of solving problem where initial data is smooth, solve same problem where initial data is piecewise constant \rightarrow this results in solution of many Riemann problems

Near $x_{j+1/2}$:



For Δt sufficiently small, you can construct the solution exactly. Take

$$v_j^1 = \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} u(x, \Delta t) dx$$

where $u(x, \Delta t)$ is the exact solution.

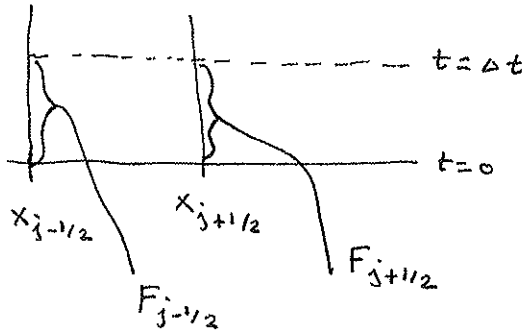
$$v_j^1 = \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} u(x, \Delta t) dx = v_j^0 - \frac{\Delta t}{\Delta x} [F_{j+1/2}^0 - F_{j-1/2}^0]$$

From exact conservation law

$$v_j^1 = v_j^0 - \frac{\Delta t}{\Delta x} [F_{j+1/2}^0 - F_{j-1/2}^0]$$

where $F_{j+1/2}^0 = \frac{1}{\Delta t} \int_0^{\Delta t} f(u(x_{j+1/2}, t)) dt$

$$F_{j-1/2}^0 = \frac{1}{\Delta t} \int_0^{\Delta t} f(u(x_{j-1/2}, t)) dt$$



For Δt sufficiently small u is constant along the interface between each Riemann problem.

$$\Rightarrow F_{j+1/2} = f(u(x_{j+1/2}, t))$$

$$u^*(v_j^0, v_{j+1}^0)$$

Godunov's Method

$$v_j^{n+1} = v_j^n - \frac{\Delta t}{\Delta x} (f(u^*(v_j^n, v_{j+1}^n)) - f(u^*(v_{j-1}^n, v_j^n)))$$

where $u^*(u_l, u_r)$ is the exact solution of the Riemann problem

$$u_t + f(u)_x = 0$$

$$u(x, 0) = \begin{cases} u_l, & x < 0 \\ u_r, & x > 0 \end{cases}$$

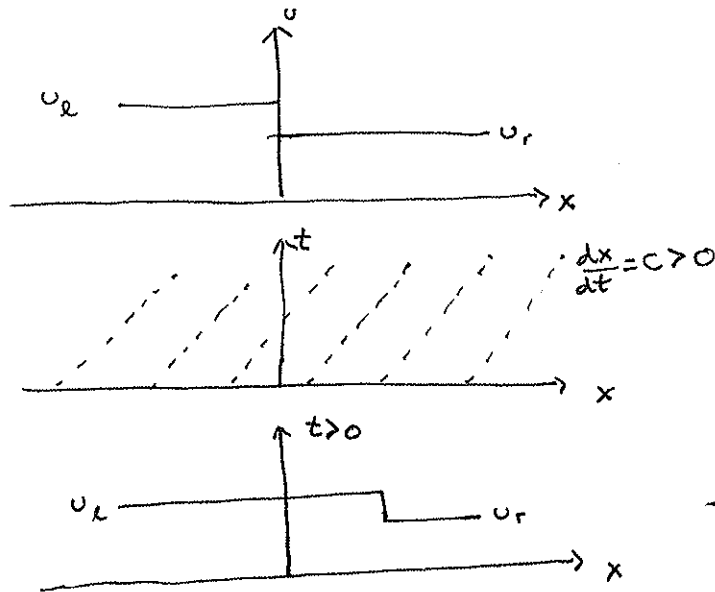
along $x = 0$ for $t > 0$

Apply the method for the linear advection equation:

$$u_t + cu_x = 0, \quad c > 0, \quad u(x, 0) = u_0(x)$$

here, Flux, $f(u) = cu$

The Riemann problem is $u_t + cu_x = 0, \quad u(x, 0) = \begin{cases} u_L, & x < 0 \\ u_R, & x > 0 \end{cases}$



$$\rightarrow u^*(u_L, u_R) = u_L$$

$$\rightarrow f(u^*(u_L, u_R)) = cu_L$$

So Godunov's method becomes

$$v_j^{n+1} = v_j^n - \frac{\Delta t}{\Delta x} (c v_j^n - c v_{j-1}^n)$$

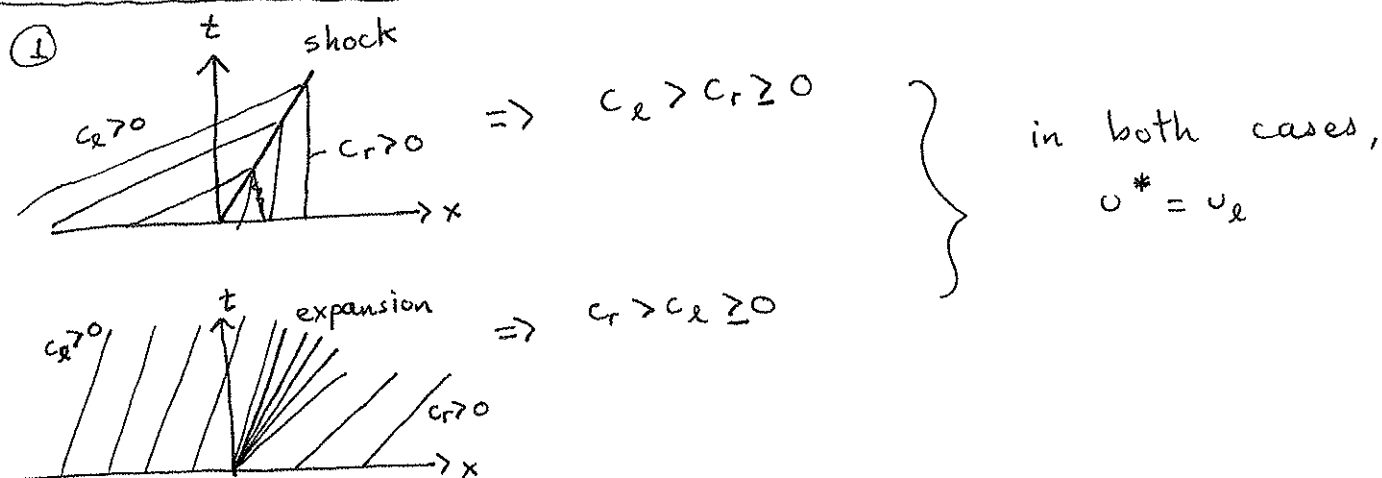
$$= v_j^n - \frac{c \Delta t}{\Delta x} (v_j^n - v_{j-1}^n) \rightarrow \text{First-order upwind method}$$

For a general flux function, Godunov may be regarded as a nonlinear upwind method (first order accurate).

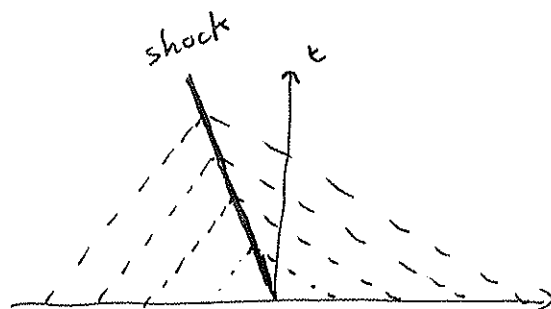
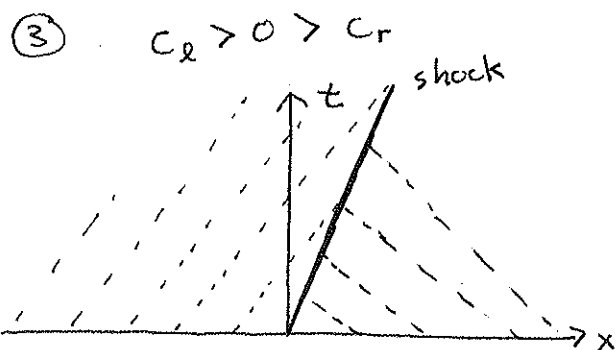
Now consider the Riemann problem for a nonlinear, scalar flux function, with $f''(u) \neq 0$, (a convex flux function). The characteristic speed is monotone if $f''(u) \neq 0$.

There are 4 cases to consider

set $c_l = f'(u_l)$ - characteristic speed associated w/ left state
 $c_r = f'(u_r)$ - characteristic speed associated w/ right state



② $c_l < 0, c_r < 0$
 these are the opposite pictures to case ① - shocks to the left, expansions to the left
 $\Rightarrow u^* = u_r$

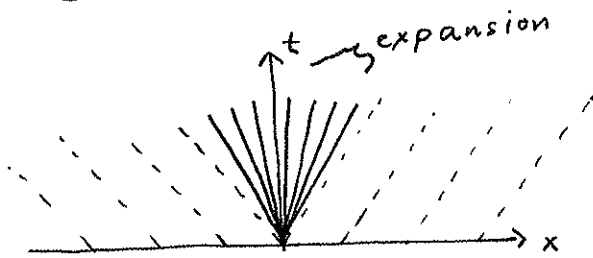


$$\frac{f(u_r) - f(u_l)}{u_r - u_l} > 0$$

$$\frac{f(u_r) - f(u_l)}{u_r - u_l} < 0$$

(4)

$$c_l < 0 < c_r$$



hardest case because
 u^* is neither u_l or u_r

Recall, the expansion solution is

$$u(x, t) = \begin{cases} u_l, & \frac{x}{t} \leq c_l \\ z, & c_l < \frac{x}{t} < c_r \\ u_r, & \frac{x}{t} \geq c_r \end{cases}$$

where z solves $f'(z) = \frac{x}{t}$

→ u^* solves $f'(u^*) = 0$, "sonic case"

You can show that all cases are covered by

$$f(u^*(u_l, u_r)) = \begin{cases} \min_{u_l \leq u \leq u_r} f(u), & u_l \leq u_r \leftarrow \text{inequalities} \\ \max_{u_r \leq u \leq u_l} f(u), & u_r < u_l \leftarrow \text{strict inequalities} \end{cases}$$

Review of Big Picture: Scalar Conservation Laws

$$u_t + f(u)_x = 0, \quad |x| < \infty, \quad t = 0, \quad u(x, 0) = u_0(x)$$

• conservative Finite volume scheme

$$v_j^{n+1} = v_j^n - \frac{\Delta t}{\Delta x} (F_{j+1/2}^n - F_{j-1/2}^n)$$

• centered methods

Lax Friedrichs, 1st order

Lax Wendroff, 2nd order

• upwind methods

Godunov, 1st order

1st-order methods \rightarrow dissipation near jumps

2nd-order methods \rightarrow dispersion near jumps

we want to discuss high resolution methods:

in smooth regions of solution
these are methods that are at least second order accurate but avoid oscillations near jumps, and where shock speeds are computed accurately.

High Resolution Methods

There are various approaches to obtain high resolution methods:

- 1) Flux limiters
- 2) slope limiters.

Flux Limiters (see Leveque 16.2 - thin green book)

begin with a conservative scheme

$$v_i^{n+1} = v_i^n - \frac{\Delta t}{\Delta x} \left(G(v_i^n, v_{i+1}^n) - G(v_{i-1}^n, v_i^n) \right)$$

where $G(u_L, u_R)$ is a numerical Flux function.

To develop a high resolution method, we would like

$G \approx G_{\text{HIGH}}$ when u is smooth

and $G \approx G_{\text{LOW}}$ near shocks

Set $G(u_e, u_r) = G_L(u_e, u_r) + \phi [G_H(u_e, u_r) - G_L(u_e, u_r)]$

where G_L is a low order Flux (eg LF)

G_H is a high order Flux (eg LW)

ϕ is a limiter.

want $\phi \approx 1$ when u is smooth
 ≈ 0 near shock

Consider the simple case of linear advection

$$u_t + cu_x = 0, \quad c > 0$$

L.W. : $v_j^{n+1} = v_j^n - \frac{\tau}{2} (v_{j+1}^n - v_{j-1}^n) + \frac{\tau^2}{2} (v_{j+1}^n - 2v_j^n + v_{j-1}^n)$, $\tau = \frac{c\Delta t}{\Delta x}$

This is usual way of writing LW, but rewrite as the following:

$$v_j^{n+1} = v_j^n - \tau (v_j^n - v_{j-1}^n) - \frac{1}{2} \tau (1-\tau) (v_{j+1}^n - 2v_j^n + v_{j-1}^n)$$

1st order
 upwind
 ↓
 stable if
 $0 \leq \tau \leq 1 \rightarrow$

Flux correction to make
 scheme 2nd order

For this range of τ , $\tau(1-\tau) > 0$
 therefore this Flux correction term
 is "anti-diffusive", to remove the
 diffusive error in upwind method

If the solution is smooth, then the anti-diffusive
 correction is effective, ie it makes the scheme 2nd order.
 However, if the solution is not smooth, the correction
 overcorrects to give oscillations.

$$v_j^{n+1} = v_j^n - \sigma (v_j^n - v_{j-1}^n) - \frac{1}{2} \sigma (1 - \sigma) (v_{j+1}^n - 2v_j^n + v_{j-1}^n)$$

The corresponding flux in the conservative scheme

$$v_j^{n+1} = v_j^n - \frac{\Delta t}{\Delta x} \left(G(v_j^n, v_{j+1}^n) - G(v_{j-1}^n, v_j^n) \right)$$

$$\text{is } G(v_j^n, v_{j+1}^n) = \underbrace{+c v_j^n}_{G_L} + \underbrace{\frac{1}{2} c (1 - \sigma) (v_{j+1}^n - v_j^n)}_{G_H - G_L}$$

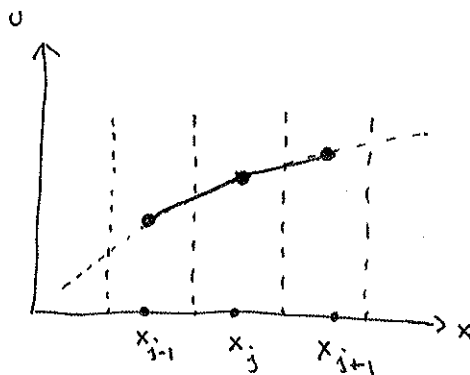
limited Flux: limited

$$G(v_j^n, v_{j+1}^n) = c v_j^n + \frac{1}{2} c (1 - \sigma) (v_{j+1}^n - v_j^n) \phi_j$$

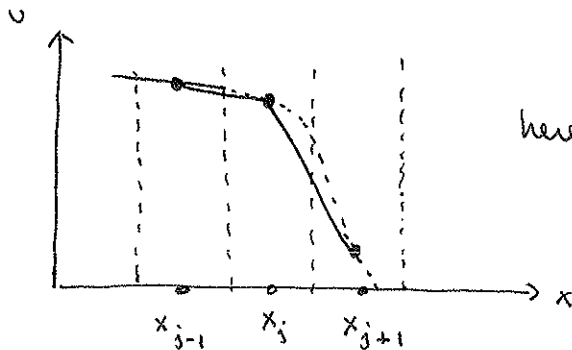
The limiter needs some measure of smoothness of solution

One choice:

$$\phi_j = \frac{v_j^n - v_{j-1}^n}{v_{j+1}^n - v_j^n} = \text{ratio of successive differences}$$



here, $\phi_j \approx 1 \rightarrow$ in smooth, monotone case, $\phi_j \approx 1$
(points of extrema do present problems)



here, $\phi_j \neq 1$

Set $\phi_j = \Phi(\Theta_j)$

ie let ϕ_j be some function of Θ_j

there are many such functions (Bee, Van Leer, etc...)

Need some way to guide the choice of limiters:

\Rightarrow TVD Methods (total variation diminishing)

define the total variation of a grid function:

$$TV(v_j^n) \equiv \sum_{j=-\infty}^{\infty} |v_j^n - v_{j-1}^n| \rightarrow \text{add up all successive differences in absolute value}$$

usually assume that v_j approaches a constant as $j \rightarrow \pm\infty$
so that $TV(v_j)$ is bounded.

A TVD method is one in which

$$TV(v_j^{n+1}) \leq TV(v_j^n) \text{ for all } n$$

It turns out that solutions of the scalar conservation law have the continuous version of ~~the~~ TVD.

$$TV(u) = \int_{-\infty}^{\infty} |u_x| dx$$

so want to reproduce this behavior in numerical scheme.

For our purposes,

TVD will imply no oscillations.

We will construct limiter functions such that TVD occurs.

Flux Limiters (continued)

Consider $u_t + cu_x = 0$

Lax-Wendroff:

$$v_j^{n+1} = \underbrace{v_j^n - \sigma(v_j^n - v_{j-1}^n)}_{\substack{\text{upwind} \\ \text{method} \\ 1^{\text{st}} \text{ order}}} - \underbrace{\frac{\sigma}{2}(1-\sigma)(v_{j+1}^n - 2v_j^n + v_{j-1}^n)}_{\substack{2^{\text{nd}} \text{ order} \\ \text{correction}}}, \quad \sigma = \frac{a\Delta t}{\Delta x}, 0 \leq \sigma \leq 1$$

the Flux with the limiter ϕ is

$$G(v_j^n, v_{j+1}^n) = cv_j^n + \frac{1}{2}c(1-\sigma)(v_{j+1}^n - v_j^n)\phi_j$$

Set the limiter to be some function of a grid function Θ_j^n , where Θ_j^n measures the smoothness of the solution

$$\phi_j^n = \Phi(\Theta_j^n), \quad \Theta_j^n = \frac{v_j^n - v_{j-1}^n}{v_{j+1}^n - v_j^n}$$

want to construct Φ st the resulting method is TVD.

Lax-Wendroff scheme with limiter

$$v_j^{n+1} = v_j^n - \frac{\Delta t}{\Delta x} \left[c(v_j^n - v_{j-1}^n) + \frac{c}{2}(1-\sigma) \left((v_{j+1}^n - v_j^n)\phi_j^n - (v_j^n - v_{j-1}^n)\phi_{j-1}^n \right) \right]$$

May rewrite as

$$v_j^{n+1} = v_j^n - \left[\sigma - \frac{1}{2}\sigma(1-\sigma)\phi_{j-1}^n \right] (v_j^n - v_{j-1}^n) - \left[\frac{\sigma}{2}(1-\sigma)\phi_j^n \right] (v_{j+1}^n - v_j^n)$$

Theorem (by Harten)

In order for a method of the form

$$v_j^{n+1} = v_j^n - \alpha_{j-1} (v_j^n - v_{j-1}^n) + \beta_j (v_{j+1}^n - v_j^n)$$

to be TVD, the following conditions are sufficient

$$\alpha_{j-1} \geq 0, \quad \beta_j \geq 0, \quad \alpha_{j-1} + \beta_j \leq 1 \quad \text{For all } j$$

For a proof, see Loeveque, p. 178.

In our case, we have

$$\left. \begin{aligned} \alpha_{j-1} &= \sigma - \frac{\sigma}{2}(1-\sigma) \phi_{j-1}^n \\ \beta_j &= -\frac{\sigma}{2}(1-\sigma) \phi_j^n \end{aligned} \right\} \text{this selection is not unique.}$$

Notice, if $\phi_j^n \approx 1$, $\beta_j < 0$ because $0 \leq \sigma \leq 1$.

Another possibility is

$$\alpha_{j-1} = \sigma - \frac{1}{2}\sigma(1-\sigma) \phi_{j-1}^n + \frac{1}{2}\sigma(1-\sigma) \phi_j^n \left(\frac{v_{j+1}^n - v_j^n}{v_j^n - v_{j-1}^n} \right)$$

$$\beta_j = 0$$

Then the sufficient conditions for TVD become

$$0 \leq \alpha_{j-1} \leq 1$$

$$\Rightarrow 0 \leq \sigma \left[1 + \frac{1}{2}(1-\sigma) \left(\underbrace{\phi_j^n}_{\Phi(\phi_j^n)} \underbrace{\left(\frac{v_{j+1}^n - v_j^n}{v_j^n - v_{j-1}^n} \right)}_{\frac{1}{\phi_j^n}} - \underbrace{\phi_{j-1}^n}_{\Phi(\phi_{j-1}^n)} \right) \right] \leq 1$$

$$\Rightarrow 0 \leq 1 + \frac{1}{2}(1-\sigma) \left(\frac{\Phi(\phi_j^n)}{\phi_j^n} - \Phi(\phi_{j-1}^n) \right) \leq 1$$

You can show under what conditions this inequality is satisfied

$$0 \leq 1 + \frac{1}{2}(1-\sigma) \left(\frac{\Phi(\theta_j^n)}{\theta_j^n} - \Phi(\theta_{j-1}^n) \right) \leq 1$$

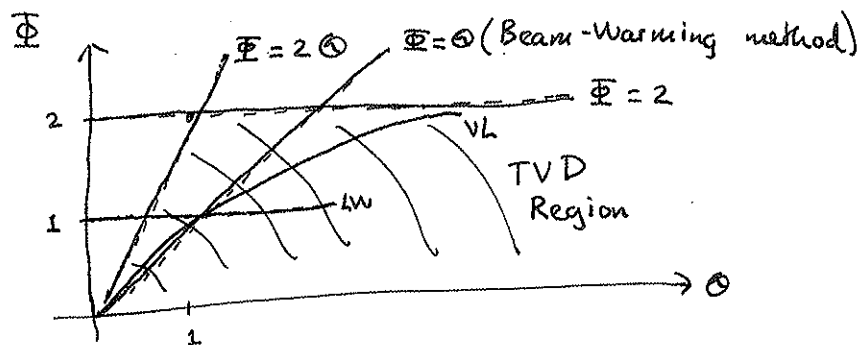
This condition is satisfied if $0 \leq \sigma \leq 1$
and if

$$\left| \frac{\Phi(\theta_j^n)}{\theta_j^n} - \Phi(\theta_{j-1}^n) \right| \leq 2, \quad \theta_j^n = \frac{v_j^n - v_{j-1}^n}{v_{j+1}^n - v_j^n}$$

Notice, if $\theta_j^n < 0$, then neighboring differences in v_j^n have opposite sign $\overbrace{\dots}^{\text{sign}}, \overbrace{\dots}^{\text{sign}}$.
Safest to set $\Phi(\theta_j^n) = 0$ when $\theta_j^n < 0$. This avoids the growth of peaks or dips, but also makes the scheme only first order accurate near smooth local extrema.

So to satisfy the above condition, let

$$0 \leq \frac{\Phi(\theta_j^n)}{\theta_j^n} \leq 2, \quad 0 \leq \Phi(\theta_j^n) \leq 2$$



Note, $\Phi = 1$ For Lax-Wendroff
 $\Phi = 0$ For Beam-Warming method

A requirement for 2nd order is that $\Phi(1) = 1$, smoothly.

we can write down many Functions that stay in TVD region and pass through $\Phi(1) = 1$, smoothly.

van Leer limiter $\Phi = \frac{|\theta| + \theta}{1 + |\theta|} \Rightarrow \text{TVD}$

Example

$$u_t + u_x = 0, \quad t > 0 \quad (c=1)$$

$$u(x,0) = \tanh(\lambda x), \quad \lambda - \text{parameter}$$

$$\text{exact solution is } u(x,t) = \tanh(\lambda(x-t))$$

Choose a grid $x_j = j \Delta x$ and choose large enough range of j st $u \rightarrow \text{constant}$ for $|j|$ large.

Integrate numerically using the Flux-limited Lax-Wendroff method

$$\phi_j = \Phi(\theta_j^n), \quad \Phi(\theta_j^n) = \frac{|\theta_j^n| + \theta_j^n}{1 + |\theta_j^n|}$$

$$\theta_j^n = \frac{v_j^n - v_{j-1}^n}{v_{j+1}^n - v_j^n} = \frac{\delta_{-x} v_j^n}{\delta_{+x} v_j^n}$$

$$\Rightarrow \phi_j^n = \frac{\left| \frac{\delta_{-x} v_j^n}{\delta_{+x} v_j^n} \right| + \frac{\delta_{-x} v_j^n}{\delta_{+x} v_j^n}}{1 + \left| \frac{\delta_{-x} v_j^n}{\delta_{+x} v_j^n} \right|} = \frac{|\delta_{-x} v_j^n| + \delta_{-x} v_j^n \frac{\delta_{+x} v_j^n}{|\delta_{+x} v_j^n|}}{|\delta_{+x} v_j^n| + |\delta_{-x} v_j^n|}$$

$\swarrow \text{sgn}(x)$

$$\rightarrow \phi_j^n = \frac{|\delta_{-x} v_j^n| + \delta_{-x} v_j^n \cdot \text{sgn}(\delta_{+x} v_j^n)}{|\delta_{+x} v_j^n| + |\delta_{-x} v_j^n| + \epsilon}, \quad \epsilon \sim 10^{-10}$$

this avoids singularity as you approach $u \sim \text{const}$ for $|x|$ large

Slope Limiters

- this approach is based on Godunov's method

$$u_t + F(u)_x = 0$$

Godunov's method (1st order upwind method)

$$v_j^{n+1} = v_j^n - \frac{\Delta t}{\Delta x} \left(F_{j+1/2}^n - F_{j-1/2}^n \right)$$

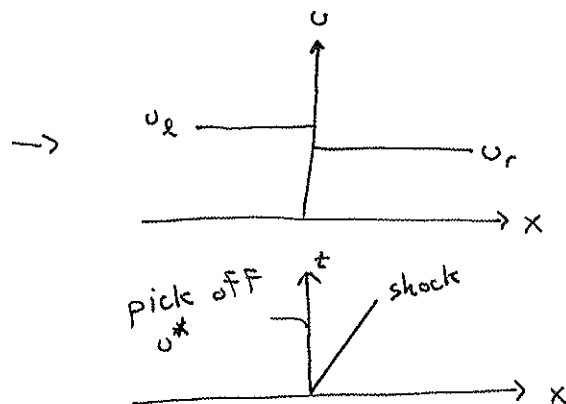
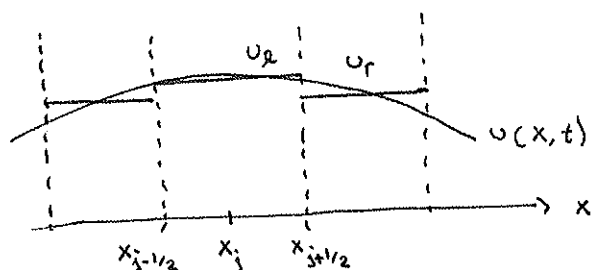
where $F_{j+1/2}^n = F(u^*(v_j^n, v_{j+1}^n))$

the state $u^*(u_l, u_r)$ is found by solving a Riemann problem

$$u_t + F(u)_x = 0, \quad u(x, 0) = \begin{cases} u_l, & x < 0 \\ u_r, & x > 0 \end{cases}$$

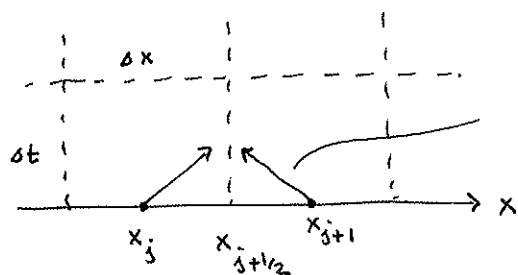
then $u^* = u(x, t)$ for $x=0, t>0$.

To obtain high resolution method, must somehow increase the order of accuracy. In order to do this, you need to include some derivative information into the states u_l, u_r .

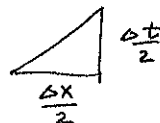


To increase accuracy, instead of using piecewise constant approximations, use piecewise linear approximations, (or quadratic, cubic, etc...).

For second-order approximations, need to include slope information \rightarrow think midpoint rule.



linear projection to center of interfaces



use Taylor series

$$u\left(x_j + \frac{\Delta x}{2}, t_n + \frac{\Delta t}{2}\right) = u(x_j, t_n) + \frac{\Delta x}{2} u_x(x_j, t_n) + \frac{\Delta t}{2} u_t(x_j, t_n) + \dots$$

$$= u + \frac{\Delta x}{2} u_x - \frac{\Delta t}{2} f'(u) u_x + \dots$$

$$= u + \underbrace{\frac{1}{2} \left(1 - \frac{\Delta t}{\Delta x} f'(u) \right) \Delta x u_x}_{\text{slope correction}} + \dots$$

For the Riemann problem about $x_{j+1/2}$, take

$$u_L = v_j^n + \frac{1}{2} \left(1 - \frac{\Delta t}{\Delta x} f'(v_j^n) \right) \delta_{+x} v_j^n$$

$$u_R = v_{j+1}^n - \frac{1}{2} \left(1 + \frac{\Delta t}{\Delta x} f'(v_{j+1}^n) \right) \delta_{-x} v_j^n$$

this results in a 2nd order extension of Godunov's method

There is no limiting in these formulas so that oscillations near shocks would occur. To suppress oscillation we want to include a slope limiter.

One effective choice:

$$\bullet \quad u_L = v_j^n + \frac{1}{2} \left(1 - \frac{\Delta t}{\Delta x} \max(f'(v_j^n), 0) \right) \Delta v_j^n$$

$$\Delta v_j^n = \text{minmod}(\delta_{+x} v_j^n, \delta_{-x} v_j^n)$$

where the minimum modulus function is defined by

$$\text{minmod}(a, b) = \begin{cases} a & \text{if } ab > 0, |a| < |b| \\ b & \text{if } ab > 0, |a| > |b| \\ 0 & \text{otherwise} \end{cases}$$

$$\bullet \quad u_r = v_{j+1}^n - \frac{1}{2} \left(1 + \frac{\Delta t}{\Delta x} \min(f'(v_{j+1}^n), 0) \right) \Delta v_{j+1}^n$$

11/13/03

Systems of Conservation Laws

$$u_t + \underline{f}(u)_x = 0$$

where u is a vector of m state variables and \underline{f} is a vector of m flux functions. The system is hyperbolic if $\underline{f}_u(u)$ is diagonalizable with real eigenvalues, $\lambda_1(u) \leq \lambda_2(u) \leq \dots \leq \lambda_m(u)$

Finite volume method

$$v_j^{n+1} = v_j^n - \frac{\Delta t}{\Delta x} \left(F_{j+1/2}^n - F_{j-1/2}^n \right)$$

Here

$$v_j^n \approx \frac{1}{\Delta x} \int_{x_{j-1/2}}^{x_{j+1/2}} u(x, t_n) dx$$

where $F_{j+1/2}^n$ is a numerical flux function

Standard methods such as Lax-Friedrichs and Lax-Wendroff carry over

$$F_{j+1/2}^n = \frac{1}{2} \left(\bar{F}(v_j^n) + \bar{F}(v_{j+1}^n) \right) - \frac{\Delta x}{2 \Delta t} \left(v_{j+1}^n - v_j^n \right) \quad \text{- Lax-Friedrichs}$$

$$F_{j+1/2}^n = \frac{1}{2} \left(\bar{F}(v_j^n) + \bar{F}(v_{j+1}^n) \right) - \frac{\Delta t}{2 \Delta x} \bar{F}_v(v) \left(\bar{F}(v_{j+1}^n) - \bar{F}(v_j^n) \right) \quad \text{- Lax-Wendroff}$$

\uparrow
 $\frac{1}{2}(v_j^n + v_{j+1}^n)$

Godunov's method

$$F_{j+1/2}^n = \bar{F}(u^*(v_j^n, v_{j+1}^n))$$

where $u^*(u_l, u_r)$ is found by solving the Riemann problem

$$u_t + \bar{F}(u)_x = 0$$

$$u(x, 0) = \begin{cases} u_l & , x < 0 \\ u_r & , x > 0 \end{cases}$$

take $u^* = u(0, t)$, $t > 0$

The solution of the Riemann problem for systems is more complicated and more costly computationally. Often solve an approximate Riemann problem instead. Usually these are based on some linearization,

$$u_t + \underline{A} u_x = 0, \quad u(x, 0) = \begin{cases} u_l & , x < 0 \\ u_r & , x > 0 \end{cases}$$

where \underline{A} approximates the Jacobian \bar{F}_u . You require

$$\underline{A} \rightarrow \bar{F}_u(u) \quad \text{as} \quad u_l \rightarrow u \quad \text{and} \quad u_r \rightarrow u$$

Typically, take \underline{A} to be the Jacobian evaluated at some average state \hat{U} , $\underline{A} = \underline{F}_U(\hat{U})$, where $\hat{U} = \frac{U_L + U_R}{2}$

this is Roe's method.

The solution of the linear problem is straightforward.

Compute the eigenvalues of \underline{A} , (because these are the characteristic speeds of the linear problem) and the corresponding eigenvectors:

$$\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_m$$

$$\Gamma_1, \Gamma_2, \dots, \Gamma_m$$

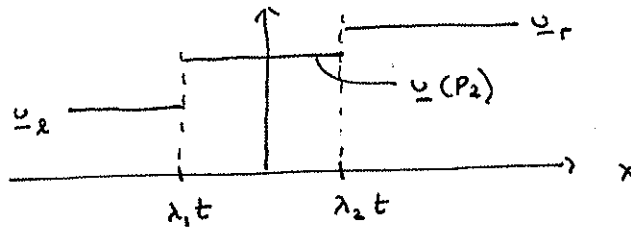
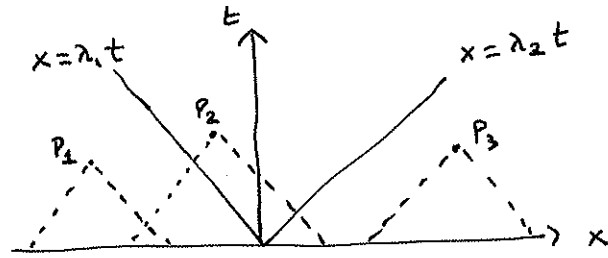
compute $\alpha_1, \alpha_2, \dots, \alpha_m$ such that

$$\alpha_1 \Gamma_1 + \alpha_2 \Gamma_2 + \dots + \alpha_m \Gamma_m = U_R - U_L$$

Then the solution becomes

$$\begin{aligned} U(x, t) &= U_L + \sum_{\lambda_k \leq \frac{x}{t}} \alpha_k \Gamma_k \\ &= U_R - \sum_{\alpha_k > \frac{x}{t}} \alpha_k \Gamma_k \end{aligned}$$

Example, $m=2$, $\lambda_1 < 0$, $\lambda_2 > 0$

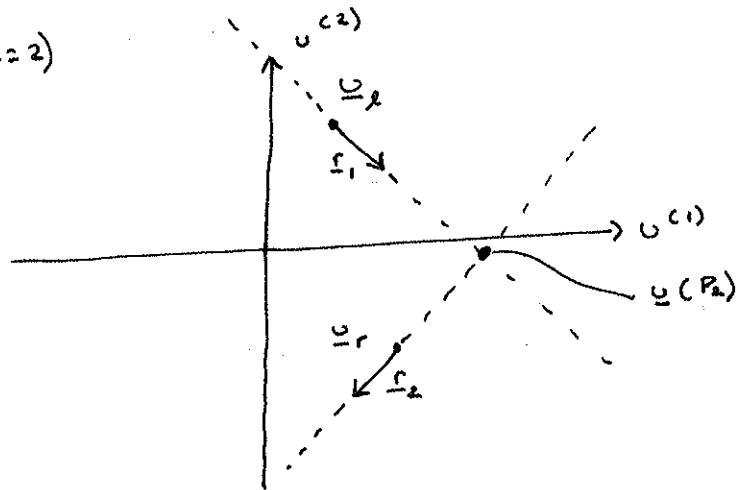


$$u(P_1) = u_l$$

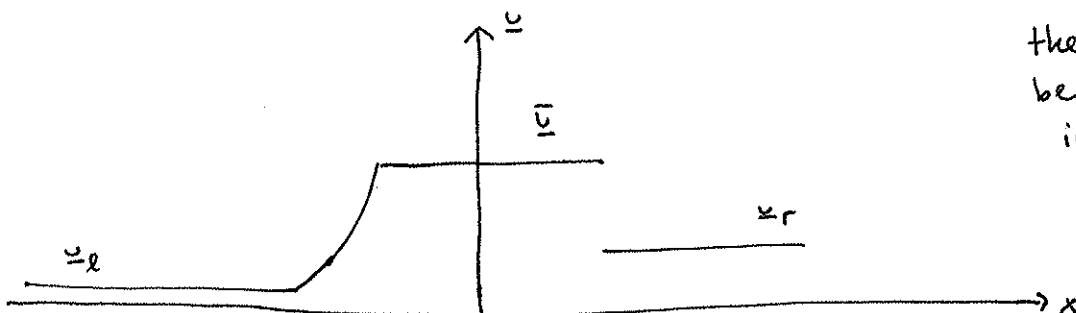
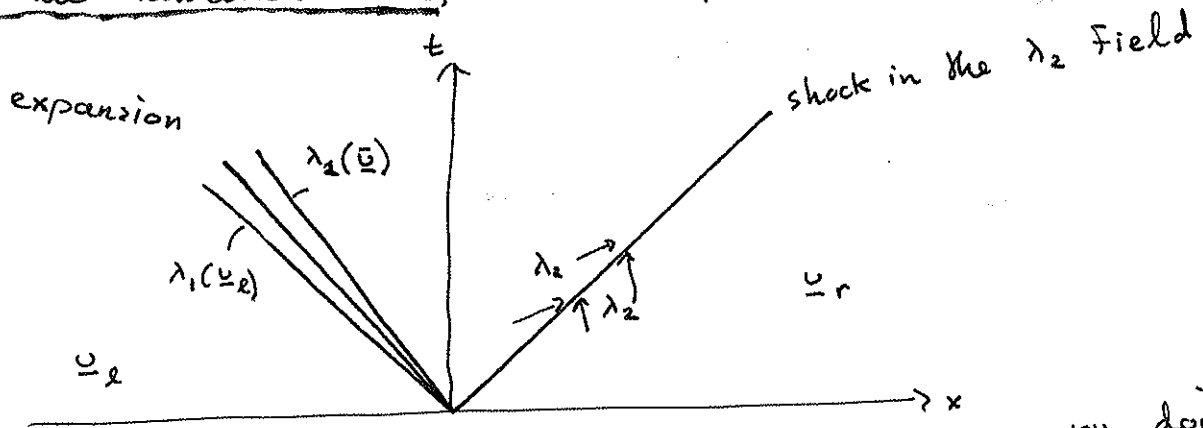
$$u(P_2) = u_l + \alpha_1 \Gamma_1$$

$$u(P_2) = u_l + \alpha_1 \Gamma_1 + \alpha_2 \Gamma_2 = u_r$$

Phase plane: ($m=2$)



For the nonlinear case, an example is the following, for $m=2$



you don't know the structure beforehand, typical involves iteration

Multiple Dimensions

Consider two-dimensions

$$u_t + f(u)x + g(u)y = 0, \quad u(x, y, 0) = u_0(x, y)$$

where $u(x, y, t)$ is the state variable and $f(u), g(u)$ are fluxes in the x, y directions, respectively.

Finite volume scheme - obtained by integrating this equation in x, y, t . Grid: $x_j = j\Delta x$, $y = k\Delta y$

$$v_{j,k}^n = \frac{1}{\Delta x \Delta y} \int_{x_{j-1/2}}^{x_{j+1/2}} \int_{y_{k-1/2}}^{y_{k+1/2}} u(x, y, t_n) dy dx$$

Integrate on a cell to get

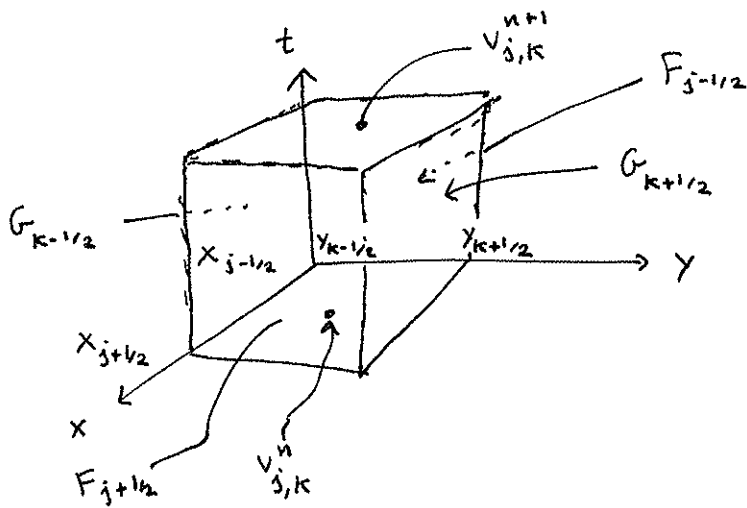
$$v_{j,k}^{n+1} = v_{j,k}^n - \frac{\Delta t}{\Delta x} \left(F_{j+1/2,k}^n - F_{j-1/2,k}^n \right) - \frac{\Delta t}{\Delta y} \left(G_{j,k+1/2}^n - G_{j,k-1/2}^n \right)$$

where

$$F_{j+1/2,k}^n = \frac{1}{\Delta y \Delta t} \int_{t_n}^{t_{n+1}} \int_{y_{k-1/2}}^{y_{k+1/2}} f(u(x_{j+1/2}, y, t)) dy dt$$

$$G_{j,k+1/2}^n = \frac{1}{\Delta x \Delta t} \int_{t_n}^{t_{n+1}} \int_{x_{j-1/2}}^{x_{j+1/2}} g(u(x, y_{k+1/2}, t)) dx dt$$

referred to as an
unsplit scheme - there are
also directional splitting



Lax-Friedrichs Formulas

$$F_{j+1/2,k}^n = \frac{1}{2} \left(F(v_{j,k}^n) + F(v_{j+1,k}^n) \right) - \frac{\Delta x}{2\Delta t} (v_{j+1,k}^n - v_{j,k}^n)$$

$$G_{j,k+1/2}^n = \frac{1}{2} \left(g(v_{j,k}^n) + g(v_{j,k+1}^n) \right) - \frac{\Delta y}{2\Delta t} (v_{j,k+1}^n - v_{j,k}^n)$$

Directional Splitting

General scheme, $v_{j,k}^{n+1} = S_y(\Delta t) S_x(\Delta t) v_{j,k}^n$

where $S_x(\Delta t)$ denotes one time step of a scheme to solve $u_t + f(u)_x = 0$, a 1D scheme. Then $S_y(\Delta t)$ denotes one time step of a scheme to solve $u_t + g(u)_y = 0$. This is only first order accurate in time, no matter how you solve each individual 1D problem. This can be avoided by implementing Strang splitting.

Strang splitting: $v_{j,k}^{n+1} = S_x\left(\frac{\Delta t}{2}\right) S_y(\Delta t) S_x\left(\frac{\Delta t}{2}\right) v_{j,k}^n$

→ 2nd order accurate (at most)

Strang splitting: $v_{j,k}^{n+1} = S_x\left(\frac{\Delta t}{2}\right) S_y(\Delta t) S_x\left(\frac{\Delta t}{2}\right)$

For many steps:

$$v_{j,k}^{n+2} = \underbrace{\left[S_x\left(\frac{\Delta t}{2}\right) S_y(\Delta t) S_x\left(\frac{\Delta t}{2}\right) \right]}_{S_x(\Delta t)} \left[S_x\left(\frac{\Delta t}{2}\right) S_y(\Delta t) S_x\left(\frac{\Delta t}{2}\right) \right]$$

becomes more efficient

CFL condition becomes something like:

$$\frac{\Delta t}{\Delta x} \max(\lambda_p) + \frac{\Delta t}{\Delta y} \max(\nu_p) \leq \text{CFL}$$

where λ_p eigenvalues of F_u
 ν_p eigenvalues of g_u

ELLIPTIC EQUATIONS

Recall the 2nd order PDE

$$A u_{xx} + 2B u_{xy} + C u_{yy} = D$$

The equation is elliptic if $B^2 - AC < 0$

Coordinate transform leads to $u_{\xi\xi} + u_{\eta\eta} = \hat{D}$, Poisson's eqn

or consider, $u_{xx} + u_{yy} = f(x, y, u, u_x, u_y)$

More generally, in K dimensions, the PDE

$$\sum_{p,q=1}^K a_{p,q}(x) \frac{\partial^2 u}{\partial x_p \partial x_q} + \sum_{p=1}^K b_p(x) \frac{\partial u}{\partial x_p} + c(x) u = d(x)$$

with $a_{p,q} = a_{q,p}$, is elliptic if $\underline{A} = [a_{p,q}(x)]$ is positive definite.

Elliptic Eqns

11/17/03

(70)

Canonical Form

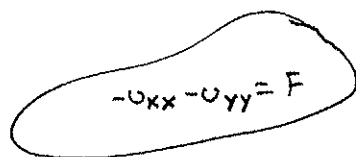
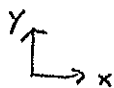
$$-u_{xx} - u_{yy} = F(x, y)$$

Poisson's Eq

Laplace eqn

$$-u_{xx} - u_{yy} = 0$$

boundary value problem



$$du + \beta \frac{\partial u}{\partial n} = g$$

Some properties of Laplace's equation

$$u_{xx} + u_{yy} = 0, \quad x^2 + y^2 < 1$$

$$u = \text{given on } x^2 + y^2 = 1 \quad (\text{Dirichlet})$$

Formulate in polar coordinates

$$\frac{1}{r} (r u_r)_r + \frac{1}{r^2} u_{\theta\theta} = 0, \quad 0 < r < 1, \quad 0 \leq \theta \leq 2\pi$$

$$u(1, \theta) = f(\theta)$$

Solution via separation of variables

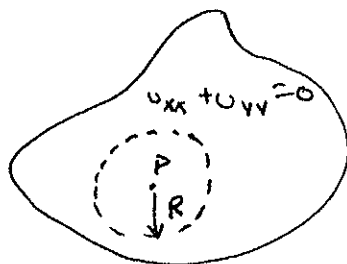
$$u(r, \theta) = a_0 + \sum_{n=1}^{\infty} r^n [a_n \cos n\theta + b_n \sin n\theta]$$

$$a_0 = \frac{1}{2\pi} \int_0^{2\pi} f(\theta) d\theta$$

$$a_n = \frac{1}{\pi} \int_0^{2\pi} f(\theta) \cos n\theta d\theta$$

$$b_n = \frac{1}{\pi} \int_0^{2\pi} f(\theta) \sin n\theta d\theta$$

Notice that at $r=0$, $u(r, \theta) = a_0 = \frac{1}{2\pi} \int_0^{2\pi} f(\theta) d\theta$, which is the mean value of the solution on the perimeter. The same property holds more generally.

Mean-value property for Laplace's Eqn

u_p = mean value of u on circle $r=R$

this holds for any point P and any radius R provided you stay in domain

useful to prove maximum principle, uniqueness, ...

Max / Min Principle

If $u(x,y)$ satisfies Laplace's equation for a domain Ω , then the max (or min) of u must occur on $\partial\Omega$. Proof by contradiction and use the mean value property.

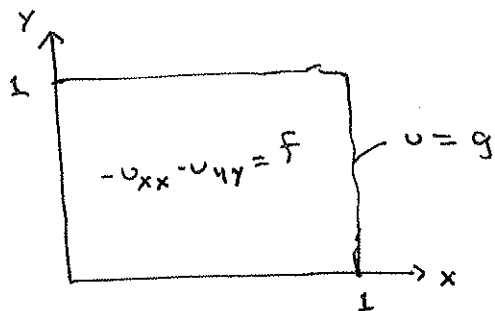
Uniqueness of solution (for Dirichlet problem)

$$\begin{aligned} u_{xx} + u_{yy} &= 0 & x, y \in \Omega \\ u &= \text{given} & x, y \in \partial\Omega \end{aligned}$$

Consider the difference of two solutions u and v and use max/min principle to show that the difference is identically zero.

Numerics

Model problem: $-u_{xx} - u_{yy} = f(x,y)$ $0 < x < 1, 0 < y < 1$
 $u = g(x,y)$ on boundary



Solve using Finite differences:

$$\begin{aligned} \text{grid} \quad x_j &= j\Delta x, \quad \Delta x = 1/N \\ y_k &= k\Delta y, \quad \Delta y = 1/M \end{aligned}$$

$$\text{Set } v_{j,k} \approx u(x_j, y_k)$$

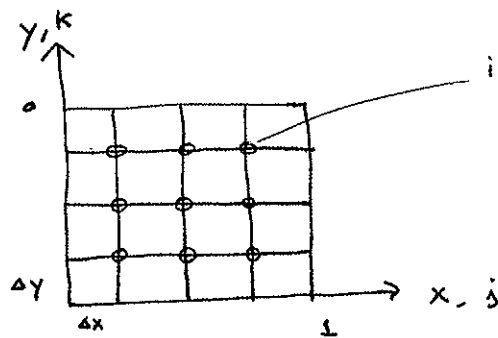
Replace u_{xx} and u_{yy} by centered differences:

$$\left(-\frac{1}{\Delta x^2} \delta_x^2 - \frac{1}{\Delta y^2} \delta_y^2 \right) v_{j,k} = f(x_j, y_k) \quad \text{for } 1 \leq j \leq N-1, 1 \leq k \leq M-1$$

(71)

Set $v_{j,k} = g(x_j, y_k)$ for (j,k) on the boundary

Example suppose $N = M = 4$



interior grid points

we have $(N-1)(M-1)$ linear equations for $v_{j,k}$ in the interior grid points

these linear systems result in a system

$$\underline{A} \underline{v} = \underline{c}$$

we have to decide how to arrange \underline{v} . The question is at what grid point to start and then what direction do you solve in?

$$\text{let } \underline{v} = [v_{11} \ v_{21} \ v_{31} \ v_{32} \ v_{22} \ v_{32} \ v_{13} \ v_{23} \ v_{33}]^T$$

$$\underline{v} = \begin{bmatrix} v_{11} \\ v_{21} \\ v_{31} \\ v_{12} \\ v_{22} \\ v_{32} \\ v_{13} \\ v_{23} \\ v_{33} \end{bmatrix} \quad \underline{A} = \begin{bmatrix} \times & \times & & \times & & & & & \\ & \times & \times & \times & & & & & \\ & & \times & \times & & & & & \\ \times & & & & \times & \times & & \times & \\ & \times & & \times & \times & \times & & \times & \\ & & \times & & \times & \times & & & \times \\ & & & \times & & & \times & \times & \\ & & & & \times & & \times & \times & \times \\ & & & & & \times & & \times & \times \end{bmatrix} \quad \underline{c} = \begin{bmatrix} \times \\ \times \\ \times \\ \times \\ 0 \\ \times \\ \times \\ \times \\ \times \end{bmatrix}$$

where \times = nonzero value

In general

$$\underline{V} = \begin{bmatrix} \begin{matrix} V_{1,1} \\ V_{2,1} \\ \vdots \\ V_{N-1,1} \end{matrix} \\ \begin{matrix} V_{1,2} \\ \vdots \\ V_{N-1,2} \end{matrix} \\ \vdots \\ \begin{matrix} V_{1,M-1} \\ \vdots \\ V_{N-1,M-1} \end{matrix} \end{bmatrix}$$

1st grid line

2nd grid line

$(M-1) \times (N-1)$

$$\underline{A} \text{ is } (N-1)(M-1) \times (N-1)(M-1)$$

$$\underline{A} = \begin{bmatrix} B & D & & & \\ D & B & D & & \\ & D & B & D & \\ & & \ddots & \ddots & \ddots \\ & & & D & B & D \\ & & & & D & B \end{bmatrix}$$

where $\underline{B} =$

tridiagonal $(N-1) \times (N-1)$

$$\begin{bmatrix} \frac{2}{\Delta x^2} + \frac{2}{\Delta y^2} & -\frac{1}{\Delta x^2} & & & \\ -\frac{1}{\Delta x^2} & \frac{2}{\Delta x^2} + \frac{2}{\Delta y^2} & -\frac{1}{\Delta x^2} & & \\ & \ddots & \ddots & \ddots & \\ & & -\frac{1}{\Delta x^2} & \frac{2}{\Delta x^2} + \frac{2}{\Delta y^2} & -\frac{1}{\Delta x^2} \\ & & & -\frac{1}{\Delta x^2} & \frac{2}{\Delta x^2} + \frac{2}{\Delta y^2} \end{bmatrix}$$

$$\underline{D} = \frac{-1}{\Delta y^2} \underline{I}$$

\underline{c} is a vector of boundary contributions

Solvability

Can we solve the linear system $\underline{A} \underline{v} = \underline{c}$?
 For this problem, you can show that \underline{A} is symmetric and positive definite, therefore \underline{A} is nonsingular and the system $\underline{A} \underline{v} = \underline{c}$ has a unique solution for any choice of \underline{c} .

Consider the eigenvalue problem:

$$\underline{A} \underline{w} = \lambda \underline{w}$$

and show that $\lambda > 0$.

$$\rightarrow \left(-\frac{1}{\Delta x^2} \delta_x^2 - \frac{1}{\Delta y^2} \delta_y^2 \right) w_{j,k} = \lambda w_{j,k} \quad \text{for } 1 \leq j \leq N-1, 1 \leq k \leq M-1$$

with $w_{j,k} = 0$ on boundary. The equation is a constant coefficient difference equations, which implies that you can apply separation of variables:

$$w_{j,k} = a^j b^k, \quad a, b = \text{constant}$$

$$\begin{aligned} \rightarrow \delta_x^2 a^j b^k &= b^k \delta_x^2 a^j \\ &= b^k (a^{j-1} - 2a^j + a^{j+1}) \\ &= a^j b^k \left(\frac{1}{a} - 2 + a \right) \end{aligned}$$

$$\Rightarrow -a^j b^k \left[\frac{1}{\Delta x^2} \left(\frac{1}{a} - 2 + a \right) + \frac{1}{\Delta y^2} \left(\frac{1}{b} - 2 + b \right) \right] = \lambda a^j b^k$$

$$\rightarrow \boxed{\lambda = -\frac{1}{\Delta x^2} \left(\frac{1}{a} - 2 + a \right) - \frac{1}{\Delta y^2} \left(\frac{1}{b} - 2 + b \right)}$$

trick: assign $\frac{1}{a} - 2 + a = -2 + 2\cos\theta$
 and $\frac{1}{b} - 2 + b = -2 + 2\cos\phi$

$$\rightarrow \cancel{1/2} \quad \frac{1}{a} + a = 2 \cos \Theta \quad \rightarrow \quad a^2 - 2a \cos \Theta + 1 = 0$$

$$\rightarrow a = \frac{2 \cos \Theta \pm \sqrt{4 \cos^2 \Theta - 4}}{2}$$

$$a = \cos \Theta \pm i \sin \Theta$$

$$\rightarrow a = e^{\pm i \Theta}$$

and likewise, $b = e^{\pm i \phi}$

\Rightarrow solutions are a linear combination of

$$(e^{\pm i \Theta})^j, (e^{\pm i \phi})^k$$

use

$$w_{j,k} = c (e^{i j \Theta} - e^{-i j \Theta}) (e^{i k \phi} - e^{-i k \phi}), \quad c = \text{constant}$$

which satisfies BCs at $j=0$ or $k=0$ ($w_{j,k} = 0$ on $j=0, k=0$)

$$\Rightarrow w_{j,k} = c \sin(j \Theta) \sin(k \phi) \quad \left(\begin{array}{l} \text{constant } c \text{ absorbed} \\ \text{Factor from } \sin = \frac{e^{i\cdot} - e^{-i\cdot}}{2} \end{array} \right)$$

note, $w_{j,k} = 0$ at $j=N \rightarrow \sin N \Theta = 0 \rightarrow N \Theta = p \pi$

$$\rightarrow \boxed{\Theta = \frac{p \pi}{N}, \quad p = 1, \dots, N-1}$$

and $w_{j,k} = 0$ at $k=M \rightarrow$

$$\boxed{\phi = \frac{q \pi}{M}, \quad q = 1, \dots, M-1}$$

Then we have

$$\lambda = \frac{-1}{\Delta x^2} \underbrace{\left(\frac{1}{a} - 2 + a \right)}_{-2 + 2 \cos \Theta} - \frac{1}{\Delta y^2} \underbrace{\left(\frac{1}{b} - 2 + b \right)}_{-2 + 2 \cos \phi}$$

$$\lambda = \frac{2}{\Delta x^2} \underbrace{(1 - \cos \Theta)}_{> 0} + \frac{2}{\Delta y^2} \underbrace{(1 - \cos \phi)}_{> 0}$$

\Rightarrow
 $\lambda > 0$, strictly greater than zero
 therefore A is positive definite

Solvability For More General Cases

Often can show solvability for linear systems obtained via finite differences for elliptic PDE by noting diagonal dominance.

Def: A matrix \underline{A} is diagonally dominant if

$$|a_{ii}| \geq \sum_{\substack{j=1 \\ j \neq i}}^N |a_{ij}| \text{ for all } i$$

and strictly diagonally dominant if

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^N |a_{ij}| \text{ for all } i$$

\underline{A} is an $N \times N$ matrix

Thm: If \underline{A} is strictly diagonally dominant, then it is nonsingular.

consider $\underline{A} \underline{x} = \underline{0} \Rightarrow \underline{x} = \underline{0}$ is the only solution

Example: $-u_{xx} - u_{yy} + \underbrace{au_x + bu_y + cu}_{\text{lower order terms}} = d(x, y)$

Approximate using

$$\left(-\frac{1}{\Delta x^2} \delta_x^2 - \frac{1}{\Delta y^2} \delta_y^2 + \frac{a}{2\Delta x} \delta_{0x} + \frac{b}{2\Delta y} \delta_{0y} + c \right) v_{i,k} = d(x_i, y_k)$$

domain: $0 \leq x \leq 1, 0 \leq y \leq 1$ with $u = g$ on the boundary

Question: Is this problem solvable?

The matrix \underline{A} is not necessarily positive definite or symmetric, so we examine the diagonal and off-diagonal elements.

diagonal element: $\frac{2}{\Delta x^2} + \frac{2}{\Delta y^2} + c > 0$ For $\Delta x^2, \Delta y^2$ sufficiently small

sum of the magnitude off-diagonal elements

$$S = \left| -\frac{1}{\Delta x^2} + \frac{a}{2\Delta x} \right| + \left| -\frac{1}{\Delta x^2} - \frac{a}{2\Delta x} \right| + \left| -\frac{1}{\Delta y^2} + \frac{b}{2\Delta y} \right| + \left| -\frac{1}{\Delta y^2} - \frac{b}{2\Delta y} \right|$$

Suppose Δx and Δy are small enough so that

$$\frac{1}{\Delta x^2} \geq \frac{|a|}{2\Delta x} \Rightarrow \Delta x \leq \frac{2}{|a|} \Rightarrow \boxed{|a|\Delta x \leq 2}$$

$$\frac{1}{\Delta y^2} \geq \frac{|b|}{2\Delta y} \Rightarrow \Delta y \leq \frac{2}{|b|} \Rightarrow \boxed{|b|\Delta y \leq 2}$$

Notice a and b represent convective terms so therefore this is kind of like a CFL condition.

If $|a|\Delta x \leq 2$ and $|b|\Delta y \leq 2$ then

$$S = \left(\frac{1}{\Delta x^2} - \frac{a}{2\Delta x} \right) + \left(\frac{1}{\Delta x^2} + \frac{a}{2\Delta x} \right) + \left(\frac{1}{\Delta y^2} - \frac{b}{2\Delta y} \right) + \left(\frac{1}{\Delta y^2} + \frac{b}{2\Delta y} \right)$$

$$\rightarrow S = \frac{2}{\Delta x^2} + \frac{2}{\Delta y^2}$$

If $c > 0$ then the linear system is strictly diagonally dominant \Rightarrow solvable.

Suppose $c = 0$ then the problem is diagonally dominant.

If $c < 0$, you may be singular.

this is ok, just requires a little more work

If $c = 0$, then the system is only diagonally dominant, assuming $|a| \Delta x \leq 2$, $|b| \Delta y \leq 2$.

Def: An $N \times N$ matrix \underline{A} is reducible if either

a) $N = 1$ and $\underline{A} = 0$

b) $N > 1$ and a $N \times N$ permutation matrix \underline{P} exists such that

$$\underline{P}^T \underline{A} \underline{P} = \left[\begin{array}{c|c} \underline{B} & \underline{C} \\ \hline \underline{0} & \underline{D} \end{array} \right] \left. \begin{array}{l} \text{ } \\ \text{ } \end{array} \right\} \begin{array}{l} r \\ N-r \end{array}$$

this implies that some of the unknowns may be decoupled

and an $N \times N$ matrix is called irreducible if it is not reducible.

Thm If \underline{A} is irreducible and diagonally dominant and if $\star |a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^N |a_{ij}|$ for at least one value for i (strictly diagonally dominant for just one i) then \underline{A} is nonsingular.

For the discretization with $c = 0$, we have that \underline{A} is irreducible (because all equations are coupled) and diagonally dominant and \star holds for grid points near the boundary \Rightarrow the system is solvable.

When $c < 0$ it is possible that c is an eigenvalue of the problem

$$\left(-\frac{1}{\Delta x^2} \delta_x^2 - \frac{1}{\Delta y^2} \delta_y^2 + \frac{a}{2\Delta x} \delta_{0x} + \frac{b}{2\Delta y} \delta_{0y} \right) v_{j,k} = \lambda v_{j,k}$$

so that the matrix would become singular.

Convergence

$$\text{Let } Lu = -u_{xx} - u_{yy}$$

$$\text{Poisson equation: } Lu = f(x, y), \quad 0 < x < 1, \quad 0 < y < 1$$

$$u = g(x, y) \quad \text{on } \partial\Omega$$

Finite difference approximation

$$L_h v_h = \frac{-1}{\Delta x^2} \delta_x^2 v_{j,k} - \frac{1}{\Delta y^2} \delta_y^2 v_{j,k} = f(x_j, y_k) \quad \begin{matrix} 1 \leq j \leq N-1 \\ 1 \leq k \leq M-1 \end{matrix}$$

let G_0 be the interior of grid, i.e. $\{1 \leq j \leq N-1, 1 \leq k \leq M-1\}$

boundary conditions: $v_{j,k} = g(x_j, y_k)$ on ∂G_0

Problem: show convergence, i.e.

$$\max_G |v_{j,k} - u(x_j, y_k)| \rightarrow 0 \quad \text{as } \Delta x, \Delta y \rightarrow 0$$

First consider the truncation error

$$\tau_{j,k} = f(x_j, y_k) - L_h u(x_j, y_k) = O(\Delta x^2, \Delta y^2)$$

by the usual Taylor series method.

$$\text{Therefore, } \max_G |\tau_{j,k}| \rightarrow 0 \quad \text{as } \Delta x, \Delta y \rightarrow 0$$

and the scheme is consistent.

To go from consistency to convergence, require some form of regularity of L_h . (Like stability for a time-dependent problem, but the analysis is different.)

Notation

$$\|v_h\|_G = \max_{(j,k) \in G} |v_{j,k}| \quad \left(\begin{array}{l} \text{on the entire grid} \\ 0 \leq j \leq N, \quad 0 \leq k \leq M \end{array} \right)$$

$$\|v_h\|_{G_0} = \max_{(j,k) \in G_0} |v_{j,k}| \quad \left(\begin{array}{l} \text{on the interior of grid} \\ 1 \leq j \leq N-1, \quad 1 \leq k \leq M-1 \end{array} \right)$$

$$\|v_h\|_{\partial G} = \max_{(j,k) \in \partial G} |v_{j,k}| \quad \left(\begin{array}{l} \text{on the boundary only} \\ j=0, N \quad k=0, M \end{array} \right)$$

Discrete max/min principle

Theorem: If $L_h v_h \leq 0$ ($L_h v_h \geq 0$) on G_0 , then the maximum (minimum) value of $v_{j,k}$ on G occurs on ∂G .

proof: Note that $L_h v_h \leq 0$ implies

$$-\frac{1}{\Delta x^2} \delta_x^2 v_{j,k} - \frac{1}{\Delta y^2} \delta_y^2 v_{j,k} \leq 0$$

$$\left(\frac{1}{\Delta x^2} + \frac{1}{\Delta y^2} \right) v_{j,k} \leq \frac{1}{2} \left[\frac{1}{\Delta x^2} (v_{j+1,k} + v_{j-1,k}) + \frac{1}{\Delta y^2} (v_{j,k+1} + v_{j,k-1}) \right]$$

Suppose that $v_{j,k}$ is a local max on G_0 :

$$\begin{array}{cc} v_{j,k} \geq v_{j+1,k} & v_{j,k} \geq v_{j-1,k} \\ v_{j,k} \geq v_{j,k+1} & v_{j,k} \geq v_{j,k-1} \end{array}$$

use these in the formula above

$$\left(\frac{1}{\Delta x^2} + \frac{1}{\Delta y^2} \right) v_{j,k} \leq \frac{1}{2} \left[\frac{2}{\Delta x^2} v_{j,k} + \frac{1}{\Delta y^2} (v_{j,k} + v_{j,k-1}) \right]$$

if you use all four inequalities

$$\left(\frac{1}{\Delta x^2} + \frac{1}{\Delta y^2} \right) v_{j,k} \leq \frac{1}{2} \left[\frac{1}{\Delta x^2} (v_{j,k} + v_{j,k}) + \frac{1}{\Delta y^2} (v_{j,k} + v_{j,k}) \right] \leq \left(\frac{1}{\Delta x^2} + \frac{1}{\Delta y^2} \right) v_{j,k}$$

since this quantity is bounded above and below, the inequalities are equalities

$$\text{Then } v_{j,k} = v_{j+1,k} = v_{j-1,k} = v_{j,k-1} = v_{j,k+1}$$

and therefore $v_{j,k}$ is not a local max.

We now use the discrete max/min principle to prove a regularity result.

Theorem: Suppose that $w_{i,k}$ is a grid function defined on G with $w_{i,k} = 0$ on ∂G , then $\|w\|_G \leq \frac{1}{8} \|L_h w_{i,k}\|_{G_0}$

Proof: Define $F_{i,k} = L_h w_{i,k}$

and note that $-\|F\|_{G_0} \leq L_h w_{i,k} \leq \|F\|_{G_0}$

Define $z_{i,k} = \frac{1}{4} \left[(x_{i-\frac{1}{2}})^2 + (y_{k-\frac{1}{2}})^2 \right]$ on G

Note that $L_h z_h = -1$

$$\begin{aligned} \text{To see this, } \delta_x^2 (x_{i-\frac{1}{2}})^2 &= (x_{i-1-\frac{1}{2}})^2 - 2(x_{i-\frac{1}{2}})^2 + (x_{i+1-\frac{1}{2}})^2 \\ &= (x_{i-\Delta x-\frac{1}{2}})^2 - 2(x_{i-\frac{1}{2}})^2 + (x_{i+\Delta x-\frac{1}{2}})^2 \\ &= (x_{i-\Delta x})^2 - (x_{i-\Delta x}) + \frac{1}{4} - 2(x_{i-\frac{1}{2}})^2 + (x_{i+\Delta x})^2 - (x_{i+\Delta x}) + \frac{1}{4} \\ &= x_i^2 - 2\Delta x x_i + \Delta x^2 - x_i + \Delta x + \frac{1}{4} - 2(x_i^2 - x_i + \frac{1}{4}) + x_i^2 + 2x_i \Delta x + \Delta x^2 - x_i - \Delta x + \frac{1}{4} \\ &= 2\Delta x^2 \\ \Rightarrow L_h z_h &= \frac{-1}{4} \left[\frac{2\Delta x^2}{\Delta x^2} + \frac{2\Delta y^2}{\Delta y^2} \right] = -1 \end{aligned}$$

$$\text{Next, } L_h (w_h - \|F\|_{G_0} z_h) = L_h w_h + \|F\|_{G_0} \geq 0$$

$$L_h (w_h + \|F\|_{G_0} z_h) = L_h w_h - \|F\|_{G_0} \leq 0$$

Now by the discrete max/min principle

since $L_h (w_h - \|F\|_{G_0} z_h) \geq 0$ then

$\min (w_h - \|F\|_{G_0} z_h)$ occurs on boundary

and since $L_h (w_h + \|F\|_{G_0} z_h) \leq 0$ then

$\max (w_h + \|F\|_{G_0} z_h)$ occurs on boundary

So we can write

$$\max_G (w_h + \|F\|_{G_0} z_h) = \max_{\partial G} (w_h + \|F\|_{G_0} z_h) = \|F\|_{G_0} \max_{\partial G} z_h$$

because we defined $w_{j,k} = 0$ on ∂G

$$\Rightarrow w_h \leq \|F\|_{G_0} \|z_h\|_{\partial G}$$

likewise, $w_h \geq -\|F\|_{G_0} \|z_h\|_{\partial G}$ where $\|z_h\|_{\partial G} = \frac{1}{8}$

Therefore we can write

$$-\frac{1}{8} \|L_h w_h\|_{G_0} \leq w_h \leq \frac{1}{8} \|L_h w_h\|_{G_0}$$

END OF PROOF

Back to convergence:

$$\text{set } w_{j,k} = v_{j,k} - u(x_j, y_k)$$

note, $w_{j,k} = 0$ on boundary due to Dirichlet conditions

$$\text{Then } L_h w_h = L_h v_h - L_h u(x_j, y_k) = z_{j,k}$$

The regularity result gives us

$$\|w_h\|_G \leq \frac{1}{8} \|L_h w_h\|_{G_0} = \frac{1}{8} \|z_h\|_{G_0} = O(\Delta x^2, \Delta y^2)$$

$$\rightarrow \boxed{\|w_h\|_G \leq O(\Delta x^2, \Delta y^2) \text{ as } \Delta x, \Delta y \rightarrow 0}$$

Solution Schemes

Model problem:

$$\frac{1}{\Delta x^2} \delta_x^2 v_{j,k} + \frac{1}{\Delta y^2} \delta_y^2 v_{j,k} = F(x_j, y_k) \quad \begin{matrix} 1 \leq j \leq N-1 \\ 1 \leq k \leq M-1 \end{matrix}$$

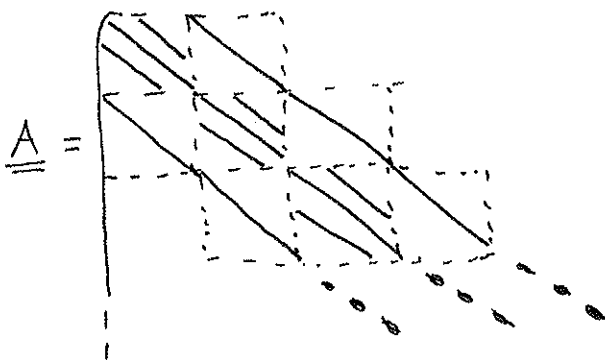
with $v_{j,k} = g(x_j, y_k)$ For j, k on ∂G

This implies a linear system $\underline{A} \underline{v} = \underline{c}$

where \underline{A} was block tri-diagonal

The problem boils down to linear algebra and whether you employ direct or iterative methods.

Direct Methods



The diagonal matrices are tri-diagonal and the sub and super diagonal matrices are diagonal

the matrix \underline{A} is banded with bandwidth $= 2 \min(N-1, M-1) + 1$
ie bandwidth is $O(\min(N, M))$

Use a direct banded matrix solver based on Gaussian elimination with partial pivoting: this results in an operation count $= O(\min(N, M)^2 \cdot NM)$
if $N = M$ then operation count $= O(N^4)$

this is a safe (reliable, stable) method
but slow

The optimal operation count is $O(N^2)$, ie one calculation for every grid point.

Direct Factorization

Since \underline{A} is symmetric, positive definite, you could use Cholesky.

$$\underline{A} = \underline{L} \underline{L}^T$$

where \underline{L} is lower triangular and retains the same band structure as original matrix \underline{A}

$$\underline{L} = \begin{bmatrix} \times & & & \\ & \times & & \\ & & \times & \\ & & & \times \end{bmatrix}$$

the bandwidth is $\frac{1}{2}$ the original bandwidth

operation count = $O(N^4)$

Block Tri diagonal Solvers

$$\underline{A} = \begin{bmatrix} C_1 & D_1 & & & \\ B_2 & C_2 & D_2 & & \\ & B_3 & C_3 & D_3 & \\ & & \ddots & \ddots & \ddots \end{bmatrix}$$

where B, C, D are $(N-1) \times (N-1)$ matrices

write down the augmented matrix

$$\begin{bmatrix} C_1 & D_1 & & \\ B_2 & C_2 & D_2 & \\ & \ddots & \ddots & \ddots \end{bmatrix} \begin{bmatrix} F_1 \\ F_2 \\ \vdots \end{bmatrix}$$

← multiply by $-B_2 C_1^{-1}$ and add to second row

$$\rightarrow \tilde{C}_2 = C_2 - B_2 C_1^{-1} D_1$$

$$\tilde{F}_2 = F_2 - B_2 C_1^{-1} F_1$$

$$\text{let } C_1^{-1} D_1 = Z_1 \rightarrow C_1 Z_1 = D_1$$

where C_1 is tridiagonal
cost to compute Z_1 is $O(N^2)$

ultimately, using the tridiagonal block solvers,
the operational cost remains $O(N^2)$

Iterative Methods (Residual Correction Methods)

Let \underline{w} denote an approximation of \underline{v} .

Then $\underline{e} = \underline{w} - \underline{v} = \text{error}$
and $\underline{r} = \underline{c} - \underline{A}\underline{w} = \text{residual}$

The error and residual are related

$$\underline{A}\underline{e} = \underline{A}(\underline{w} - \underline{v}) = \underline{A}\underline{w} - \underline{c} = -\underline{r} \Rightarrow \underline{e} = -\underline{A}^{-1}\underline{r}$$

Then the solution

$$\underline{v} = \underline{w} - \underline{e} = \underline{w} + \underline{A}^{-1}\underline{r} \rightarrow \underline{v} = \underline{w} + \underline{A}^{-1}\underline{r}$$

we don't have \underline{A}^{-1} (it is expensive to calculate).

However, suppose that \underline{B} approximates \underline{A}^{-1} , then

$$\underline{w}^{n+1} = \underline{w}^n + \underline{B}\underline{r}^n \quad \text{"residual" correction scheme}$$

There are several choices of \underline{B} :

$$\text{Let } \underline{A} = \underline{L} + \underline{D} + \underline{U} \quad \text{where}$$

$$\underline{A} = \underbrace{\begin{bmatrix} \circ & & \\ \text{lower triangular} & & \\ & \circ & \end{bmatrix}}_{\underline{L}} + \underbrace{\begin{bmatrix} \text{diagonal} & & \\ & \circ & \\ & & \circ \end{bmatrix}}_{\underline{D}} + \underbrace{\begin{bmatrix} \text{upper triangular} & & \\ & \circ & \\ & & \circ \end{bmatrix}}_{\underline{U}}$$

lower triangular portion diagonal portion upper triangular portion

Choices for \underline{B}

- ① Jacobi choice: $\underline{B} = \underline{D}^{-1}$
- ② Gauss-Seidel $\underline{B} = (\underline{L} + \underline{D})^{-1}$
- ③ Successive Over Relaxation (SOR)

$$\underline{B} = \omega (\underline{D} + \omega \underline{L})^{-1} \quad \text{where } \omega \text{ is a relaxation parameter} \\ (\omega \text{ is a scalar})$$

Comments:

- the implementation is very simple
- convergence occurs because \underline{A} comes from the discretization of an elliptic PDE

11/24/03

Example:

consider a 1D elliptic problem, $-u_{xx} = F$, $u=0$ on $\partial\Omega$

$$-\frac{1}{\Delta x^2} \delta_x^2 v_j = F_j, \quad 1 \leq j \leq N-1, \quad v_0 = v_N = 0$$

Suppose $N=4$

$$\begin{cases} -v_2 + 2v_1 = \Delta x^2 F_1 \\ -v_3 + 2v_2 - v_1 = \Delta x^2 F_2 \\ 2v_3 - v_2 = \Delta x^2 F_3 \end{cases}$$

This is a 3×3 linear system

$$\underline{A} \underline{v} = \underline{c} \quad \rightarrow \quad \begin{bmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{bmatrix} \begin{pmatrix} v_1 \\ v_2 \\ v_3 \end{pmatrix} = \Delta x^2 \begin{pmatrix} F_1 \\ F_2 \\ F_3 \end{pmatrix}$$

Suppose we instead "solve" for the diagonal elements:

$$v_1 = \frac{1}{2} (\Delta x^2 F_1 + v_2)$$

$$v_2 = \frac{1}{2} (\Delta x^2 F_2 + v_1 + v_3)$$

$$v_3 = \frac{1}{2} (\Delta x^2 F_3 + v_2)$$

add and subtract diagonal

$$v_1 = v_1 + \frac{1}{2}(\Delta x^2 f_1 + v_2 - 2v_1)$$

$$v_2 = v_2 + \frac{1}{2}(\Delta x^2 f_2 + v_1 - 2v_2 + v_3)$$

$$v_3 = v_3 + \frac{1}{2}(\Delta x^2 f_3 + v_2 - 2v_3)$$

this suggests some sort of iterative scheme:

$$\left. \begin{aligned} v_1^{n+1} &= v_1^n + \frac{1}{2}(\Delta x^2 f_1 + v_2^n - 2v_1^n) \\ v_2^{n+1} &= v_2^n + \frac{1}{2}(\Delta x^2 f_2 + v_1^n - 2v_2^n + v_3^n) \\ v_3^{n+1} &= v_3^n + \frac{1}{2}(\Delta x^2 f_3 + v_2^n - 2v_3^n) \end{aligned} \right\} \underline{\text{Jacobi}}$$

residual

If use the $n+1$ term right away:

$$\left. \begin{aligned} v_1^{n+1} &= v_1^n + \frac{1}{2}(\Delta x^2 f_1 + v_2^n - 2v_1^n) \\ v_2^{n+1} &= v_2^n + \frac{1}{2}(\Delta x^2 f_2 + v_1^{n+1} - 2v_2^n + v_3^n) \\ v_3^{n+1} &= v_3^n + \frac{1}{2}(\Delta x^2 f_3 + v_2^{n+1} - 2v_3^n) \end{aligned} \right\} \underline{\text{Gauss-Seidel}}$$

residual

Include a relaxation parameter ω

$$\left. \begin{aligned} v_1^{n+1} &= v_1^n + \frac{\omega}{2}(\Delta x^2 f_1 + v_2^n - 2v_1^n) \\ v_2^{n+1} &= v_2^n + \frac{\omega}{2}(\Delta x^2 f_2 + v_1^{n+1} - 2v_2^n + v_3^n) \\ v_3^{n+1} &= v_3^n + \frac{\omega}{2}(\Delta x^2 f_3 + v_2^{n+1} - 2v_3^n) \end{aligned} \right\} \underline{\text{SOR}}$$

Typical Algorithm

$$\left(\frac{-1}{\Delta x^2} \delta_x^2 - \frac{1}{\Delta y^2} \delta_y^2 \right) v_{j,k} = f_{j,k}, \quad 1 \leq j \leq N-1, \quad 1 \leq k \leq M-1$$

with $v_{j,k} = 0$ on boundary

For SOR:

- one sweep
- 1) Set $v_{j,k} = 0$ everywhere, pick ω
 - 2) For $n = 1, 2, \dots, n_{\max}$
 - 3) For $k = 1, \dots, M-1$
 - 4) For $j = 1, \dots, N-1$
 - 5)
$$v_{j,k} = v_{j,k} + \frac{\omega}{D} \left(f_{j,k} + \frac{1}{\Delta x^2} \delta_x^2 v_{j,k} + \frac{1}{\Delta y^2} \delta_y^2 v_{j,k} \right)$$
 - 6) stop if $\max |r_{j,k}| < \text{tol}$ — $D = \frac{2}{\Delta x^2} + \frac{2}{\Delta y^2}$

cost per sweep: $O(MN)$

total cost to converge: $O(nMN)$ — number of iterations

so the question is, how many iterations are necessary for convergence? $n \sim O(MN)$

so we're back to $O(N^4)$ — these iterative methods haven't bought us much over the direct methods. Need to use multi grid...

Analysis of Residual Correction Schemes

Iteration: $\underline{w}^{n+1} = \underline{w}^n + \underline{B} \underline{r}^n$

error eqn: $\underline{A} \underline{e}^n = -\underline{r}^n$ where $\underline{e}^n = \underline{w}^n - \underline{v}$

eliminate \underline{r}^n : $\underline{w}^{n+1} = \underline{w}^n - \underline{B} \underline{A} \underline{e}^n$

subtract \underline{v} from both sides

$$\underline{e}^{n+1} = \underline{e}^n - \underline{B} \underline{A} \underline{e}^n$$

$$\underline{e}^{n+1} = (\underline{I} - \underline{B} \underline{A}) \underline{e}^n$$

Suppose \underline{e}^0 is the initial error, then

$$\underline{e}^n = \underbrace{(\underline{I} - \underline{B} \underline{A})}_{\underline{R}}^n \underline{e}^0$$

← power, not time step

the iteration converges if

$$\lim_{n \rightarrow \infty} \underbrace{\underline{R}^n}_{\text{power}} \underline{e}^0 = 0$$

Theorem: convergence occurs for any \underline{e}^0 if \underline{R} is a convergent ~~matrix~~ matrix, i.e. iff spectral radius of \underline{R} less than one.

So the rate of convergence depends on the eigenvalues of \underline{R} . Suppose \underline{R} has eigenvalues

$$|\lambda_1| > |\lambda_2| > |\lambda_3| > \dots$$

↑
dominant
eigenvalue
(controls
convergence)

Define $\{\alpha_j\}_{j=1}^N$ such that

$$\underline{e}^0 = \sum_{j=1}^N \alpha_j \underline{\xi}_j \quad \text{where } \underline{\xi}_j \text{ is the } j^{\text{th}} \text{ eigenvector of } \underline{R},$$

which is $N \times N$ matrix

$$\rightarrow \underline{e}^1 = \underline{R} \underline{e}^0 = \sum_{j=1}^N \alpha_j \underline{R} \underline{\xi}_j = \sum_{j=1}^N \alpha_j \lambda_j \underline{\xi}_j$$

$$\rightarrow \underline{e}^2 = \sum_{j=1}^N \alpha_j \lambda_j^2 \underline{\xi}_j$$

⋮

$$\underline{e}^n = \sum_{j=1}^N \alpha_j \lambda_j^n \underline{\xi}_j$$

$$\underline{e}^n = \lambda_1^n \left(\alpha_1 \underline{\xi}_1 + \underbrace{\sum_{j=2}^N \alpha_j \left(\frac{\lambda_j}{\lambda_1} \right)^n \underline{\xi}_j}_{\rightarrow 0} \right)$$

$$\rightarrow \underline{e}^n \sim \lambda_1^n \alpha_1 \underline{\xi}_1 \rightarrow 0$$

$$\Rightarrow \frac{\|\underline{e}^{n+1}\|}{\|\underline{e}^n\|} \sim |\lambda_1| \rightarrow \frac{\|\underline{e}^n\|}{\|\underline{e}^0\|} \sim |\lambda_1|^n < \text{tol}$$

$$\Rightarrow n \log(\rho(\underline{R})) < \log(\text{tol}) \rightarrow \boxed{n > \frac{\log(\text{tol})}{\log(\rho(\underline{R}))} = O\left(\frac{-1}{\log(\rho(\underline{R}))}\right)}$$

IF $\rho(\underline{R}) \sim 0 \Rightarrow$ Fast convergence

IF $\rho(\underline{R}) \sim 1 - \epsilon \Rightarrow$ slow convergence

$\rho(\underline{R})$ - spectral radius
of \underline{R} - magnitude
of largest λ_j

Suppose $-\delta_x^2 v_j = \Delta x^2 f_j$, $1 \leq j \leq N-1$

$$v_0, v_N = 0$$

This leads to a system $\underline{A} \underline{v} = \underline{c}$ whose eigenvalues are $\mu_p = 2(1 - \cos(\frac{p\pi}{N}))$, $p = 1, \dots, N-1$.

Note - these are eigenvalues of \underline{A} , not \underline{R} , but they are related.

$$\underline{D} = 2 \underline{I}$$

For the Jacobi method, $\underline{R} = \underline{I} - \underline{D}^{-1} \underline{A}$

$$\rightarrow \underline{R} = \underline{I} - \frac{1}{2} \underline{A}$$

eigenvalue problem for \underline{R} : $\underline{R} \underline{\xi} = \lambda \underline{\xi}$

$$(\underline{I} - \frac{1}{2} \underline{A}) \underline{\xi} = \lambda \underline{\xi} \rightarrow -\frac{1}{2} \underline{A} \underline{\xi} = (\lambda - 1) \underline{\xi}$$

$$\underline{A} \underline{\xi} = \underbrace{-2(\lambda - 1)}_{\mu} \underline{\xi} \rightarrow \mu = -2(\lambda - 1) \rightarrow \lambda = 1 - \frac{\mu}{2}$$

$$\rightarrow \lambda_p = \cos\left(\frac{p\pi}{N}\right) \Rightarrow \rho(\underline{R}) = \cos \frac{\pi}{N} = 1 - \frac{1}{2} \left(\frac{\pi}{N}\right)^2 + \dots$$

really interested in logarithm of spectral radius

$$\ln(\rho(\underline{R})) \approx -\frac{1}{2} \left(\frac{\pi}{N}\right)^2 + \dots$$

$$\rightarrow n_{\text{iterate}} \sim \frac{-1}{\ln(\rho(\underline{R}))} \sim C N^2$$

Multigrid Methods

Multigrid is a technique to accelerate the convergence of a residual correction method.

$$\text{Elliptic PDE: } \begin{aligned} Lu &= f, & x \in \Omega \\ u &= g, & x \in \partial\Omega \end{aligned}$$

(Linear)

$$\text{Finite difference discretization: } \begin{aligned} L_h v_h &= f_h, & x \in \Omega_h \\ v_h &= g_h, & x \in \partial\Omega_h \end{aligned}$$

assume that $v_h \in \mathbb{R}^n$ - (v_h is a grid function in \mathbb{R}^n)

$$\text{test problem: } Lu = u_{xx} = f(x), \quad 0 \leq x \leq 1$$

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

$$\text{approximation: } \begin{aligned} v_{j+1} - 2v_j + v_{j-1} &= h^2 f(x_j), & 1 \leq j \leq N-1 \\ v_0 &= v_N = 0 \end{aligned}$$

residual correction method based on Jacobi

$$v_j^{n+1} = v_j^n + \frac{\omega}{2} (v_{j+1}^n - 2v_j^n + v_{j-1}^n - h^2 f(x_j)) \quad - (\omega\text{-Jacobi})$$

NOT SOP

if $\omega = 1$, this is classical Jacobi }
 if $0 < \omega < 1$, under-relaxed Jacobi } converges take $v_j^0 = 0$
 if $\omega > 1$, over-relaxed Jacobi } diverges

in matrix form

$$\underline{v}^{n+1} = \left(\underline{I} + \frac{\omega}{2} \underline{A} \right) \underline{v}^n - \frac{\omega h^2}{2} \underline{f}$$

$$\underline{v}^n = \begin{bmatrix} v_1^n \\ \vdots \\ v_{N-1}^n \end{bmatrix}, \quad \underline{A} = \begin{bmatrix} -2 & 1 & & \\ 1 & -2 & 1 & \\ & \ddots & \ddots & \ddots \end{bmatrix}, \quad \underline{f} = \begin{bmatrix} f(x_1) \\ \vdots \\ f(x_{N-1}) \end{bmatrix}$$

Error in the n^{th} iterate

$$\underline{e}^n = \underline{v}^n - \underline{v} \quad (\underline{v} \text{ is the true steady-state error})$$

by a similar procedure as just done, you find

$$\underline{e}^{n+1} = \left(\underline{I} + \frac{\omega}{2} \underline{A} \right) \underline{e}^n = \underline{R} \underline{e}^n$$

turns into eigenvalue problem,

$$\underline{R} \underline{\Sigma} = \lambda \underline{\Sigma}$$

write \underline{e}^0 as an eigenvalue expansion

$$\underline{e}^0 = \alpha_1 \underline{\Sigma}_1 + \alpha_2 \underline{\Sigma}_2 + \dots + \alpha_{N-1} \underline{\Sigma}_{N-1}$$

where $\underline{\Sigma}_j$ is the j^{th} eigenvector of $\underline{I} + \frac{\omega}{2} \underline{A}$

and whose eigenvalues are λ_j

From before, $\underline{e}^n = \alpha_1 \lambda_1^n \underline{\Sigma}_1 + \alpha_2 \lambda_2^n \underline{\Sigma}_2 + \dots$

We want to study the behavior of this expansion

Have $\lambda_p = 1 - \omega \left(1 - \cos\left(\frac{p\pi}{N}\right) \right)$ - eigenvalues of \underline{R} , $p=1, \dots, N-1$

j^{th} component
of p^{th} eigenvector

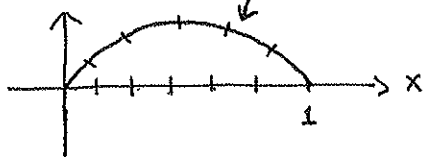
$(\underline{\Sigma}_p)_j = \sin\left(\frac{j p \pi}{N}\right)$ - the j^{th} component of p^{th} eigenvector

We examine the "modes" of the eigenvector expansion

Consider a specific case, $N=6$

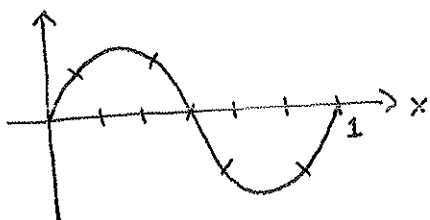
$$(\underline{\Sigma}_p)_j = \sin\left(\frac{j p \pi}{6}\right) = \sin(p \pi x_j), \quad x_j = \frac{j}{6}$$

• $p=1$: $(\underline{\Sigma}_1)_j = \sin(\pi x_j)$, $\lambda_1 = 1 - \omega \left(1 - \frac{\sqrt{3}}{2} \right)$

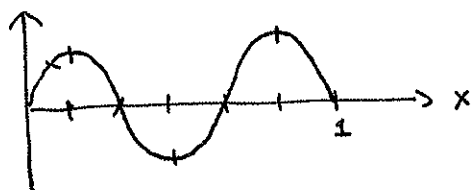


slowest mode
to converge

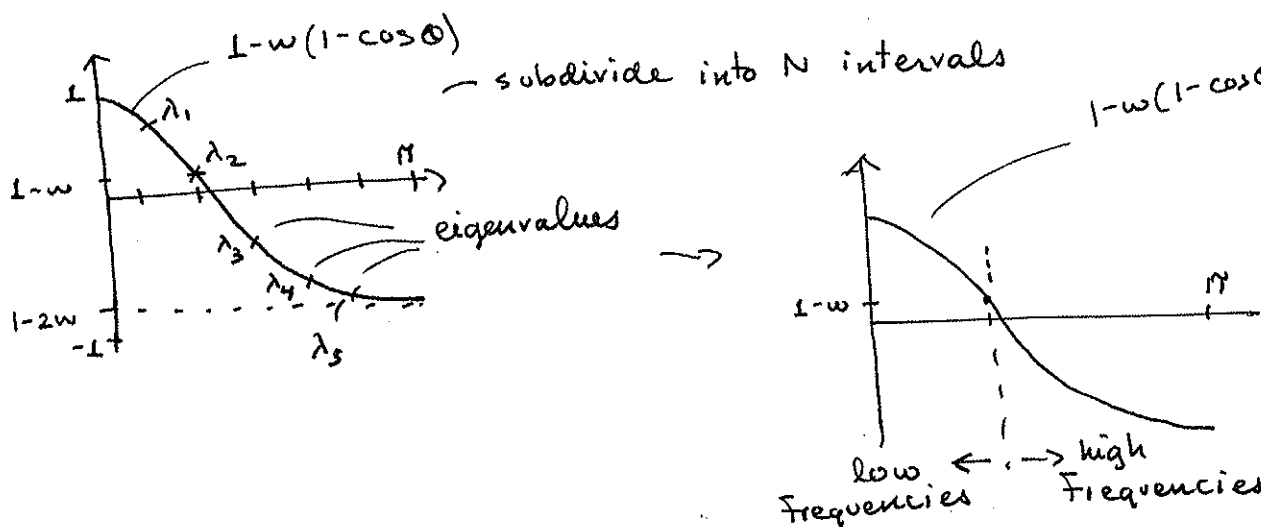
• $p=2$: $(\underline{\Sigma}_2)_j = \sin(2\pi x_j)$, $\lambda_2 = 1 - \frac{\omega}{2}$



$$p=3, \sin(3\pi x_j), \lambda_3 = 1-w$$



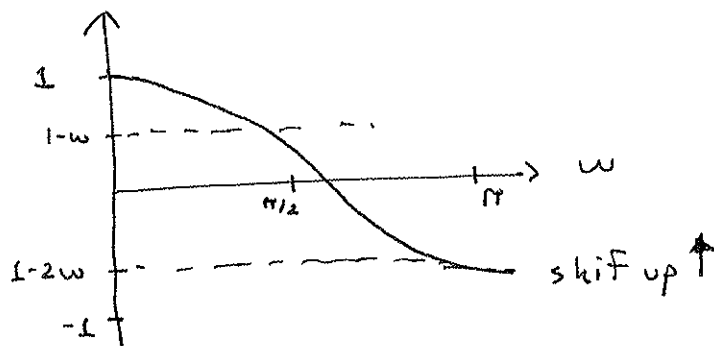
the point - the frequency represented by eigenvectors increases with p



Key observations that lead us to Multigrid

- 1) Convergence occurs for $0 < w < 1$
- 2) If you are in range of convergence, $0 < w < 1$
what is spectral radius? $\rho(R) = \lambda_1 = 1 - w(1 - \cos \frac{\pi}{N})$
 $\sim 1 - \frac{K}{N^2} \Rightarrow$ slow convergence for any choice of w
- 3) Low Frequency modes always converge slowly but the high Frequency modes can be made to converge quickly with a good choice of w .

12/1/03



Choose w such
that

$$(1-w) = -(1-2w)$$

$$\Rightarrow w = \frac{2}{3}$$

when $w = \frac{2}{3}$, the spectral radius only for the high frequency modes: $\rho(R)_{\text{high Freq}} = \frac{1}{3}$ - bounded away from 1 independently of N .

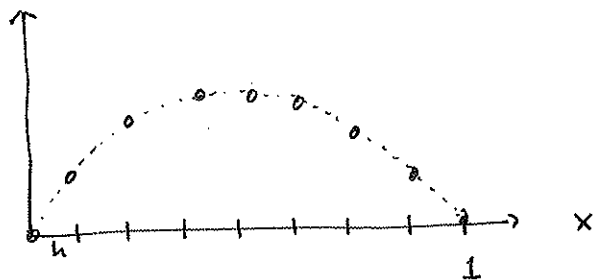
The multigrid idea is to use coarser grids to resolve low frequency components of the error.

Two-Grid Algorithm

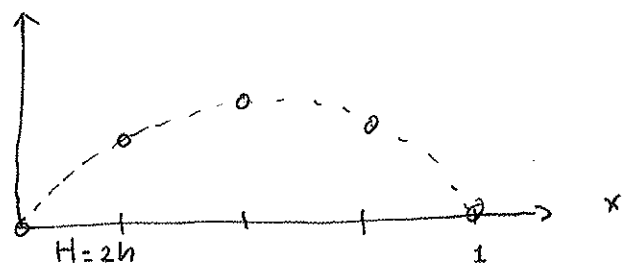
Use grids with $h = \frac{1}{N}$ (Fine grid), ultimately the grid on which we want the solution.

and use grids with $H = 2h$ (coarse grid)

Fine
grid
($N=8$)



coarse
grid
($N=4$)



Two-Grid Algorithm (continued)

step 0: initial guess, \underline{v}^0 given on Fine grid

step ①: Apply ν_1 steps of w -Jacobi (smoothing step)

$$\underline{v}^k = (\underline{I} + \frac{w}{2} \underline{A}) \underline{v}^{k-1} - \frac{w}{2} h^2 \underline{f} \quad k = 1, \dots, \nu$$

$\tilde{\underline{v}} = \underline{v}^{\nu_1}$, the new vector after ν_1 steps of w -Jacobi

the number ν_1 is chosen adaptively, but for simple code, pick $\nu_1 \approx 3$ or 4 - a small number then $\tilde{\underline{v}}$ only has low frequency components of error

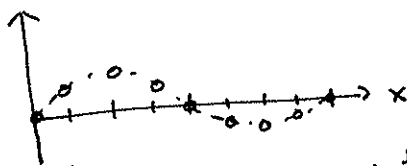
step ②: compute residual and restrict to the coarse grid

$$\underline{r} = \underline{f} - \underline{A} \tilde{\underline{v}} \quad (\text{smooth})$$

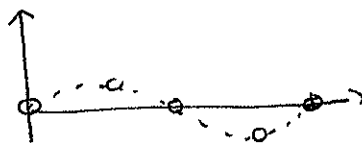
then

$$\underline{r}_H = \underline{Q} \underline{r}$$

suppose this is residual on Fine grid



simply take every other point for the coarse grid.



then for this choice, \underline{Q} is (not square)

$$\begin{bmatrix} \underline{r}_H \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & 1 & 0 \end{bmatrix}}_{\underline{Q}} \begin{bmatrix} \underline{r} \end{bmatrix}$$

instead of creating matrix and operating on \underline{r} , just pick out the values on the coarse grid for \underline{r}_H

step 3

Solve the error equation on coarse grid

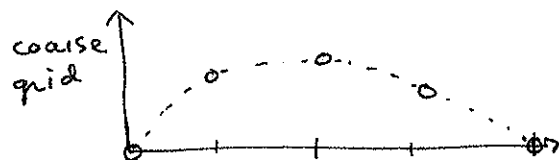
$$\underline{A}_H \underline{e}_H = \underline{r}_H$$

(this is equivalent to solving the discrete problem on coarse grid)

step 4

interpolate \underline{e}_H on to Fine grid error, \underline{e}

$$\underline{e} = \underline{P} \underline{e}_H$$



you could use piecewise linear interpolation of coarse grid values to find fine grid values

you don't necessarily need fancy interpolation because error found in high frequency can be smoothed out by a few Jacobi iterates

$$\begin{bmatrix} \underline{e} \end{bmatrix} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 1/2 & 1/2 & \dots & \\ 0 & 1 & \dots & \\ & \underline{P} & & \end{bmatrix} \begin{bmatrix} \underline{e}_H \end{bmatrix}$$

step 5

update and smooth

$$\underline{\bar{v}} = \underline{\tilde{v}} + \underline{e}$$

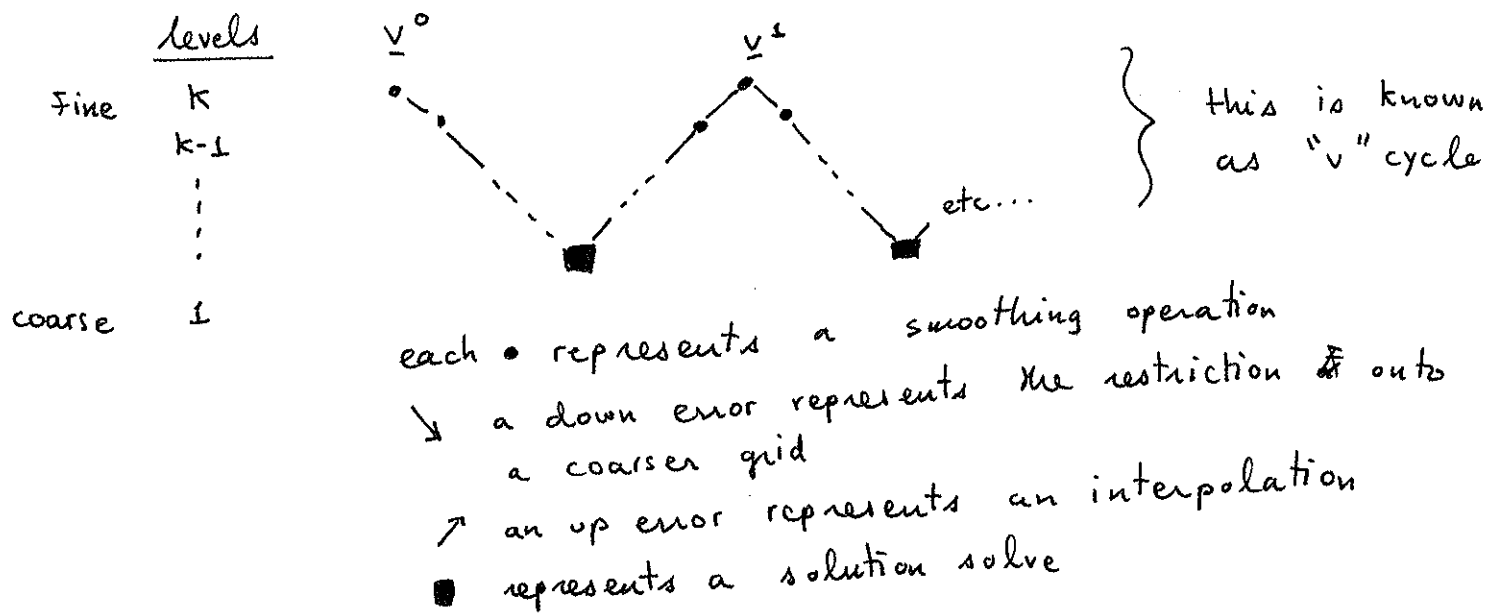
← here we have filtered out the low frequency error in $\underline{\tilde{v}}$ by adding \underline{e} we've introduced lower frequency error due to interpolation, but we now smooth that out

$$\text{set } \underline{v}^0 = \underline{\bar{v}}$$

$$\underline{v}^k = \left(\underline{I} + \frac{\omega}{2} \underline{A} \right) \underline{v}^{k-1} - \frac{\omega}{2} \underline{h}^2 \underline{F}, \quad k = 1, \dots, \nu_2$$

then $\underline{\hat{v}} = \underline{v}^{\nu_2}$ ← this is the next iterate

at this point we've taken one full two-grid algorithm iteration

Multigrids

total operation count: \propto # grid points

Final Exam

Part 1 - in class, thursday

OPEN NOTES, CLOSED BOOKS

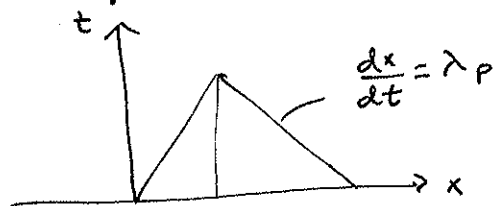
$\approx 75\%$ of Final exam

Part 2: take home, available starting saturday morning

$\approx 25\%$ email for copy, 48 hrs to complete

Topics since last midterm

1) Linear hyperbolic PDEs, $u_t + A u_x = 0$
 hyperbolic if A is diagonalizable with real eigenvalues



Basic schemes

Lax-Friedrichs 1st order

Lax-Wendroff 2nd order

upwind methods 1st order

both are considered centered methods

one-sided scheme

behavior near discontinuities
(dissipative versus dispersive)
determined by examining modified equation
stability and CFL condition
numerical vs exact domain of dependence
non-constant coeff linear PDEs
boundary conditions (inflow vs outflow)

2) Hyperbolic conservation Laws

$$u_t + f(u)_x = 0$$

integral Form (integral conservation Form)

$$\frac{d}{dt} \int_a^b u dx = -f(u) \Big|_a^b$$

shock conditions, method of characteristics

Riemann problems

numerical methods

- conservative finite volume schemes
- numerical Flux Functions (for LF, LW, etc...)
- Lax-Wendroff theorem
- Godunov's method (nonlinear upwind method, 1st order accurate)
- high resolution methods (2nd order at least, but 1st order near shocks)
 - Flux limiters, slope limiters

3) Elliptic PDEs

canonical Forms

properties of Laplace's eqn

- mean value, maximum principle

Finite difference methods

discretizations, solvability of $Ax=b$, convergence

solution schemes (direct vs iterative methods)

multigrid