

Numerical Solution of ODEs

$$v''(t) + v(t) = 0 \quad \text{2nd order differential equation}$$

general solution: $v(t) = A \cos t + B \sin t$

where A, B are constants of integration.

This two-parameter family of functions are called integral curves. In this case, the curves oscillate with period 2π . The constants select the amplitude and phase of the oscillation.

In order to obtain a unique solution, need to specify initial or boundary data.

Suppose $v(0) = \alpha, v'(0) = \beta$ (initial data)

$$\Rightarrow A = \alpha, B = \beta$$

unique solution: $v(t) = \alpha \cos(t) + \beta \sin(t)$

Or suppose instead $v(0) = \alpha, v(b) = \beta$ (boundary data)

solution is: $v(t) = \alpha \cos t + \left(\frac{\beta - \alpha \cos b}{\sin b} \right) \sin t$

unique solution if $\sin b \neq 0$, otherwise there are NO solutions or an INFINITE number of solutions.

IVP is a problem consisting of an ODE and initial data. "Solution evolves locally"

BVP is a problem consisting of an ODE and boundary data (usually specified at endpoints of the ~~domain~~ interval of interest.)

BVP are typically more difficult to handle either analytically or numerically, but the solution set is often more interesting.

Initial Value Problems

Focus on IVPs in the form of a system of 1st order eqns.

$$\underline{y}'(t) = \underline{f}(t, \underline{y}), \quad 0 \leq t \leq b, \quad \underline{y}(0) = \underline{c}$$

where $\underline{y}'(t) = \begin{bmatrix} y_1'(t) \\ y_2'(t) \\ \vdots \\ y_m'(t) \end{bmatrix}, \quad \underline{f}(t, \underline{y}) = \begin{bmatrix} f_1(t, \underline{y}) \\ f_2(t, \underline{y}) \\ \vdots \\ f_m(t, \underline{y}) \end{bmatrix}, \quad \underline{c} \in \mathbb{R}^m$

constant vector

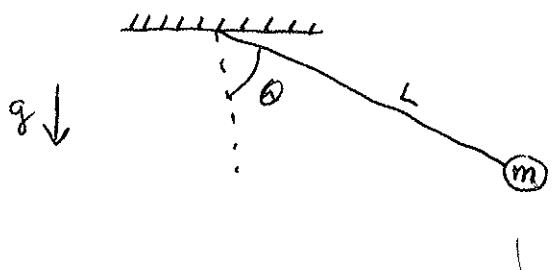
Assume that $\underline{f}(t, \underline{y})$, the "slope function" is a smooth function of t, \underline{y} . Note that if the interval is $[a, b]$, then the change of variables $\tau = a + \frac{b-a}{b} t$

transforms to $[0, b]$.

Examples

(1) Pendulum (dynamical system, governed by autonomous ODEs)

Consider a mass suspended on a rigid arm attached to a pivot:



$\theta(t)$ = angular displacement (radians)

m = mass

L = length

g = acceleration due to gravity

Newton's Second Law: (vertical direction)

$$m L \frac{d^2\theta}{dt^2} = -mg \sin\theta \Rightarrow \frac{d^2\theta}{dt^2} = -\frac{g}{L} \sin\theta$$

Initial Conditions $\theta(0) = \theta_0, \frac{d\theta(0)}{dt} = \omega_0$

First-order system

$$\begin{aligned} y_1 &= \theta(t) \\ y_2 &= \frac{d}{dt}\theta(t) \end{aligned} \Rightarrow \boxed{y_1' = y_2, y_2' = -\frac{g}{L} \sin(y_1)}$$

$$\boxed{y_1(0) = \theta_0, y_2(0) = \omega_0}$$

We expect solutions to oscillate.

IF $|\theta|$ is small then $\sin\theta \sim \theta \Rightarrow \theta'' + \frac{g}{L}\theta = 0$

The general solution is $\theta(t) = A \cos \sqrt{\frac{g}{L}} t + B \sin \sqrt{\frac{g}{L}} t$

$$\begin{matrix} \uparrow \\ \theta_0 \end{matrix} \quad \begin{matrix} \uparrow \\ \omega_0 \sqrt{\frac{g}{L}} \end{matrix}$$

what about finite Θ ?

Consider an integral of the eqn: $\Theta'' + \frac{g}{L} \sin \Theta = 0$

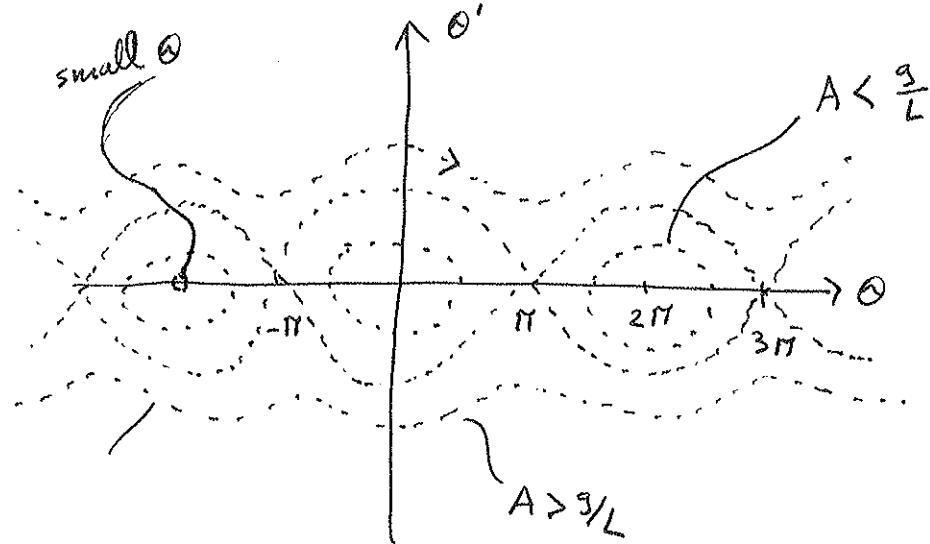
Multiply through by Θ' :

$$\Theta''\Theta' + \frac{g}{L} \sin \Theta \cdot \Theta' = 0$$

$$\Rightarrow \frac{1}{2} \Theta'^2 - \frac{g}{L} \cos \Theta = A, \quad A = \text{constant}$$

the initial conditions specify: $A = \frac{1}{2} \omega_0^2 - \frac{g}{L} \cos \Theta_0$

Plot the curves in the phase plane:



(2) Predator-Prey (dynamical systems)

3/9

Dynamical system that arises in population dynamics.
 Consider the interaction of coupled populations:

$y_1(t)$ = population of a prey species

$y_2(t)$ = population of a predator species

In isolation, assume

$$\begin{aligned} \frac{dy_1}{dt} &= \alpha y_1 \\ \frac{dy_2}{dt} &= \gamma y_2 \end{aligned} \quad \left\{ \begin{array}{l} \text{change in population is proportional} \\ \text{to the population} \end{array} \right.$$

Generally assume that $\alpha > 0$, $\gamma < 0$. Assume that the interaction is proportional to the product of the populations, and that the interaction inhibits the growth of prey and promotes the growth of predator.

With interaction:

$$\frac{dy_1}{dt} = \alpha y_1 - \beta y_1 y_2$$

$$\frac{dy_2}{dt} = \gamma y_2 + \delta y_1 y_2$$

where $\beta, \delta > 0$. Include initial conditions

$$y_1(0) = A, \quad y_2(0) = B.$$

Solutions oscillate. Consider phase plane,

equilibrium points: $(0,0)$ or $(-\frac{\delta}{\beta}, \frac{\alpha}{\beta})$

Consider a linearization about $(-\frac{\delta}{\beta}, \frac{\alpha}{\beta})$:

$$y_1 = -\frac{\gamma}{\beta} + \tilde{y}_1, \quad y_2 = \frac{\alpha}{\beta} + \tilde{y}_2$$

Find that :

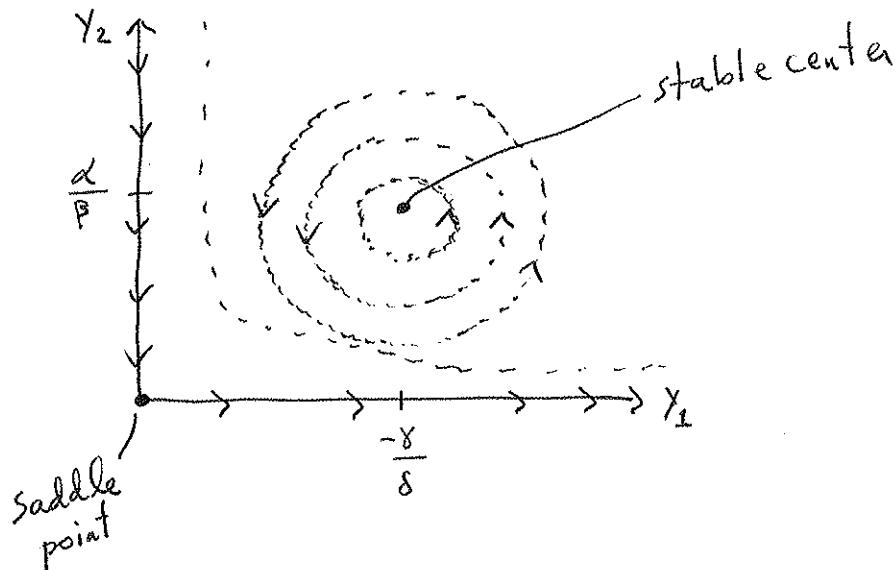
$$\tilde{y}_j'' + \tau_j \tilde{y}_j = 0, \quad \tau_j > 0 \text{ is constant}$$

Integrate :

$$\frac{dy_2}{dx_1} = \frac{\delta y_2 + \delta x_1 y_2}{\alpha y_1 - \beta y_1 y_2} = \frac{y_2(\delta + \delta y_1)}{y_1(\alpha - \beta y_2)}$$

The equation is separable (Boyce Di Prima)

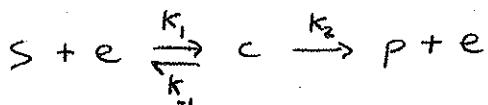
$$\Rightarrow \underbrace{\alpha \ln y_2 - \beta y_2 - \delta \ln y_1 - \delta y_1}_{\text{const. of motion}} = \text{constant}$$



(3) Biochemical Kinetics

4/4

2-step biochemical reaction involving a substrate $s(t)$, an enzyme $e(t)$, a "complex" $c(t)$, and a product $p(t)$.



k_1, k_{-1}, k_2 are rate constants.

Rate laws (laws of mass action):

$$(1) \quad \frac{ds^*}{dt^*} = -k_1 s^* e^* + k_{-1} c^* \quad *-\text{dimensional quantities}$$

$$(2) \quad \frac{de^*}{dt^*} = -k_1 s^* e^* + k_{-1} c^* + k_2 c^*$$

$$(3) \quad \frac{dc^*}{dt} = k_1 s^* e^* - k_{-1} c^* - k_2 c^*$$

$$(4) \quad \frac{dp^*}{dt} = k_2 c^*$$

Initial conditions: $s^*(0) = \bar{s}$, $e^*(0) = \bar{e}$, $c^*(0) = p^*(0) = 0$

1/17/03

notice, add eqns (2) and (3)

$$\Rightarrow \frac{d}{dt^*}(e^* + c^*) = 0 \rightarrow e^* + c^* = \bar{e}$$

Add (1), (3) and (4), again right-hand side vanishes

$$\frac{d}{dt^*}(s^* + c^* + p^*) = 0 \rightarrow s^* + c^* + p^* = \bar{s}$$

In effect we've integrated the system twice, so we are left with two ODEs and two algebraic relations.

Use these 2 algebraic relations to eliminate e^* and p^* .

$$\frac{ds^*}{dt^*} = -k_1 \bar{e} s^* + (k_1 s^* + k_{-1}) c^*, \quad s^*(0) = \bar{s}$$

$$\frac{de^*}{dt^*} = k_1 \bar{e} s^* - (k_1 s^* + k_{-1} + k_2) c^*, \quad c^*(0) = 0$$

Choose scales for t^* , s^* , c^* (make dimensionless)

$$\text{let } s = \frac{s^*}{\bar{s}}, \quad c = \frac{c^*}{\bar{e}}, \quad t = \frac{t^*}{\sqrt{k_1 \bar{e}}}$$

$$\Rightarrow \begin{cases} s' = -s + (s + K - \lambda)c, & s(0) = 1 \\ \varepsilon c' = s - (s + K)c, & c(0) = 0 \end{cases}$$

$$\text{where } \varepsilon = \frac{\bar{e}}{\bar{s}}, \quad K = \frac{k_{-1} + k_2}{k_1 \bar{s}}, \quad \lambda = \frac{k_2}{k_1 \bar{s}}$$

$$\text{Often, } \bar{e} \ll \bar{s} \Rightarrow 0 < \varepsilon \ll 1$$

$$K, \lambda = O(1)$$

Analyze the behaviour of the solution.

$$\text{set } \varepsilon = 0:$$

$$\begin{aligned} s' &= -s + (s + K - \lambda)c \\ 0 &= s - (s + K)c \Rightarrow c = \frac{s}{s + K} \end{aligned}$$

$$\Rightarrow s' = \frac{-\lambda s}{s + K}, \quad s(0) = 1$$

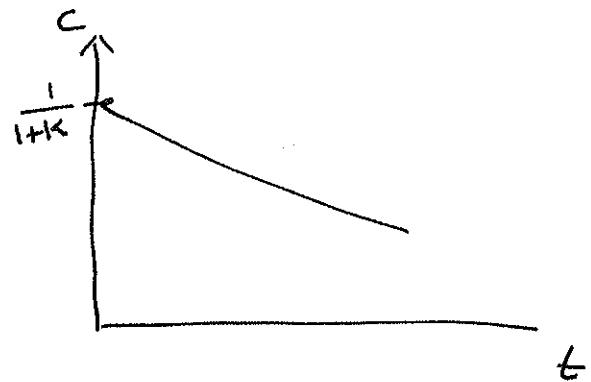
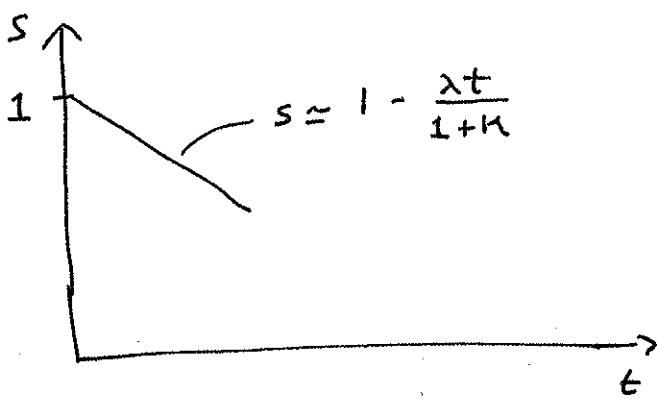
separable equation, integrate

$$\Rightarrow \boxed{s + K \ln s = -\lambda t + 1, \text{ for } \varepsilon \rightarrow 0}$$

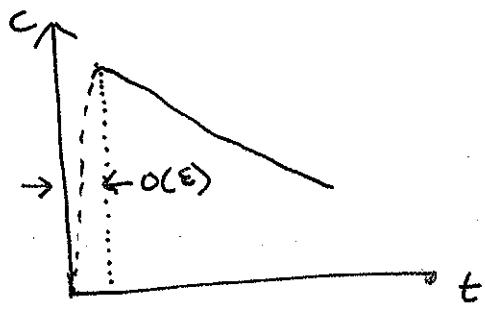
$$s + K \ln s = -\lambda t + L$$

notice, when $s \approx$ near $s = 1$, linear

(5)



Problem with this solution is that c does not satisfy initial condition, $c(0) = 0$.



Consider an "initial layer" solution. This is a small region near $t=0$ where c' is large so that $\epsilon c' = O(1)$

Set $\zeta = \frac{t}{\epsilon} \Rightarrow$ when $t = O(\epsilon)$, $\zeta = O(1)$
(exact change of variables)

$$\Rightarrow \frac{ds}{d\zeta} = \epsilon \left[-s + (s + k - \lambda) c \right]$$

$$\frac{dc}{d\zeta} = s - (s + k) c$$

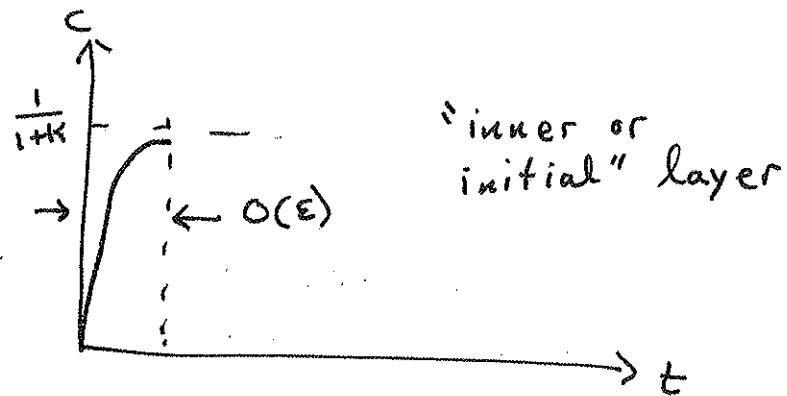
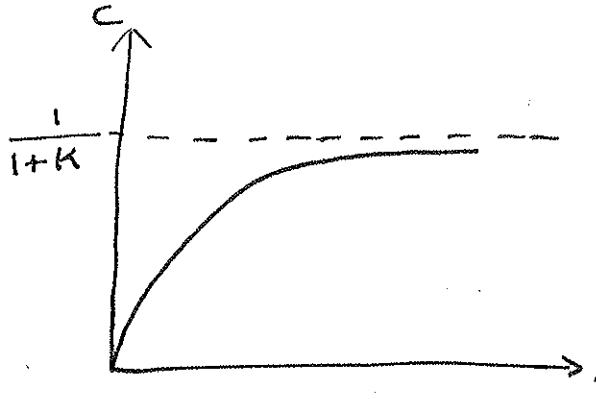
since ϵ is small, first equation goes to zero
so that easily integrated $\Rightarrow s = 1$

$$\Rightarrow \frac{dc}{d\zeta} = 1 - (1+k)c, \quad c(0) = 0$$

$$\text{Integrate: } c = \frac{1}{1+k} \left(1 - e^{-(1+k)\zeta} \right)$$

$$c = \frac{1}{1+k} \left(1 - e^{-(1+k)\varepsilon} \right)$$

$$\Rightarrow c = \frac{1}{1+k} \left(1 - e^{-(1+k)t/\varepsilon} \right)$$



The point is, s is a "slow" variable

whereas c is a "fast" variable.

Numerically, we refer to this rapid behavior
as "stiffness" in the problem.

EXAMPLE : Diffusion Problem

ODEs often arise through some reduction (analytical or numerical) of a PDE

$$u_t = (pu_x)_x + g(x, u) \quad , \quad 0 < x < 1 , \quad t > 0$$

where $u(x, t)$ = temperature

$p(x)$ = diffusivity

$g(x, u)$ = heat source

Boundary conditions: $u(0, t) = \alpha(t)$, $u(1, t) = \beta(t)$

Initial conditions: $u(x, 0) = g(x)$

use Method of Lines

introduce grid: $x_j = jh$, $h = \frac{1}{m+1}$, $1 \leq j \leq m$

$$v_j(t) \approx v(x_j, t)$$

Replace derivatives in x with finite difference approximations

$$\frac{d}{dt} v_j = \frac{1}{h} \left[P(x_j + \frac{h}{2}) \left(\frac{v_{j+1} - v_j}{h} \right) - P(x_j - \frac{h}{2}) \left(\frac{v_j - v_{j-1}}{h} \right) \right]$$

Pv_x at

$$x_j + \frac{h}{2}$$

Pv_x at $x_j = \frac{h}{2}$

set of ODEs, $1 \leq j \leq m$, $v_0(t) = \alpha(t)$

$$v_{m+1}(t) = \beta(t)$$

$$v_j(0) = q(x_j)$$

Remarks

- 1) Unlike previous cases, the number of ODEs is typically large. Might use only a lower order method and match the order of accuracy of the spatial discretization. Can be expensive computationally.
- 2) Equations are stiff

A Regularity Result for IVPs

IVP $\underline{y}' = \underline{f}(t, \underline{y})$, $\underline{y}(0) = \underline{c}$, $0 \leq t \leq b$

Theorem Let $\underline{f}(t, \underline{y})$ be continuous for all (t, \underline{y}) in a region

$$D = \left\{ 0 \leq t \leq b, -\infty < \|\underline{y}\| < \infty \right\}.$$

Moreover, assume that \underline{f} is Lipschitz in the variable \underline{y} , ie a constant L exists such that

$$\|\underline{f}(t, \underline{y}) - \underline{f}(t, \hat{\underline{y}})\| \leq L \|\underline{y} - \hat{\underline{y}}\|$$

for all $(t, \underline{y}), (t, \hat{\underline{y}})$ in D .

Then:

- (1) For any $\underline{c} \in \mathbb{R}^m$ there exists a unique solution $\underline{y}(t)$ throughout the interval $[0, b]$. The solution is differentiable.
- (2) The solution depends continuously on the initial data, ie if $\hat{\underline{y}}(t)$ also satisfies the ODE, then
$$\|\underline{y}(t) - \hat{\underline{y}}(t)\| \leq e^{Lt} \|\underline{y}(0) - \hat{\underline{y}}(0)\|$$
- (3) IF $\hat{\underline{y}}(t)$ satisfies, more generally, a perturbed ODE
$$\hat{\underline{y}}' = \underline{f}(t, \hat{\underline{y}}) + \underline{\varepsilon}(t, \hat{\underline{y}}) \text{ where } \underline{\varepsilon} \text{ is continuous and bounded in } D, \|\underline{\varepsilon}\| < M. \text{ Then}$$

$$\|\underline{y}(t) - \hat{\underline{y}}(t)\| \leq e^{Lt} \|\underline{y}(0) - \hat{\underline{y}}(0)\| + \frac{M}{L} (e^{Lt} - 1)$$

Remarks

1) If $\underline{f}(t, y)$ is smooth then $\underline{\Sigma}$ is Lipschitz if

$$\left\| \frac{\partial \underline{f}}{\partial y}(t, y) \right\| \leq L \quad \forall (t, y) \text{ in } D.$$

2) Often $\underline{\Sigma}$ cannot be shown to be Lipschitz for D if $-\infty < \|y\| < \infty$. A restricted result:
suppose \underline{f} satisfies

$$\|\underline{\Sigma}(t, \hat{y}) - \underline{\Sigma}(t, y)\| \leq L \|\hat{y} - y\| \quad \text{for}$$

$D = \{0 \leq t \leq b, \|y - c\| \leq r\}$. Suppose also that

$\|\underline{\Sigma}\| \leq M$ on D . Then the results above hold

but for $0 \leq t \leq \min\{b, r/M\}$.

Boundary Value Problems

A BVP involves an ODE defined on some interval and data supplied at more than one point on the interval.
Typically, the data is supplied on the boundaries.

2-point BVP

For example: $\underline{y}' = \underline{f}(t, \underline{y})$, $0 \leq t \leq b$

$$g(\underline{y}(0), \underline{y}(b)) = 0$$

\underline{f}

$\underline{y}(t)$ is a vector of m scalar functions, \underline{f} is a given "slope" function and g is a set of m algebraic equations

Consider the solution of an associated initial value problem, $\underline{y}' = \underline{f}(t, \underline{y})$, $\underline{y}(0) = \underline{c}$, $0 \leq t \leq b$ call the solution $\underline{y}(t; \underline{c})$. Then to solve the BVP,

$$g(\underline{c}, \underline{y}(b; \underline{c})) = 0 \quad - \text{a set of } m \text{ nonlinear algebraic equations for } \underline{c}.$$

Even if the IVP is well behaved, there is no guarantee that the solution of the BVP exists and if it exists, whether it is unique.

The solution behavior of a BVP is non-local.

Example

1) A diffusion problem

$$v_t = (pv_x)_x + g(x, v), \quad 0 \leq x \leq 1, \quad t \geq 0$$

$$\text{BCs } v(0, t) = \alpha(t), \quad v(1, t) = \beta(t)$$

$$\text{ICs } v(x, 0) = g(x)$$

Suppose that α, β approach constant values for t sufficiently large, so that as $t \rightarrow \infty$ a steady state is reached.

The steady state solves for $u = u(x)$

$$0 = (p(x)u_x) + g(x, u) , \quad 0 \leq x \leq 1$$

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

write as a first-order system:

$$\text{set } y_1(x) = u \quad (\text{temperature})$$

$$\text{and } y_2(x) = p(x)u' \quad (\text{heat flux})$$

ODEs:

$$y_1' = \frac{y_2}{p(x)} , \quad y_2' = -g(x, u) , \quad 0 \leq x \leq 1$$

$$y_1(0) = \alpha , \quad y_1(1) = \beta$$

no conditions given for y_2

often $p(x) \neq 0$, in fact, $p(x) > 0$. There are cases where $p(x) = 0$ (typically at a boundary) and then some special treatment is required.

2) Singular Perturbation Problem

Consider the second-order, linear BVP:

$$\epsilon v'' + a(x)v' + b(x)v = 0, \quad 0 \leq x \leq L$$

$$v(0) = \alpha, \quad v(L) = \beta$$

Suppose $a(x)$ and $b(x)$ are smooth and that $\epsilon > 0$ and small in magnitude. We want to determine the solution behavior in the limit $\epsilon \rightarrow 0$.

Consider the case $\epsilon = 0$.

$$a(x)v' + b(x)v = 0$$

$$\rightarrow \frac{v'}{v} = -\frac{b(x)}{a(x)} \quad \text{assume } a(x) \neq 0 \text{ for } 0 \leq x \leq L$$

Integrate:

$$\ln v = - \int \frac{b(x)}{a(x)} dx + C$$

$$v = C \exp \left\{ - \int \frac{b(x)}{a(x)} dx \right\}$$

Choose C such that v satisfies the boundary condition:

e.g., if $v(0) = \alpha$ then

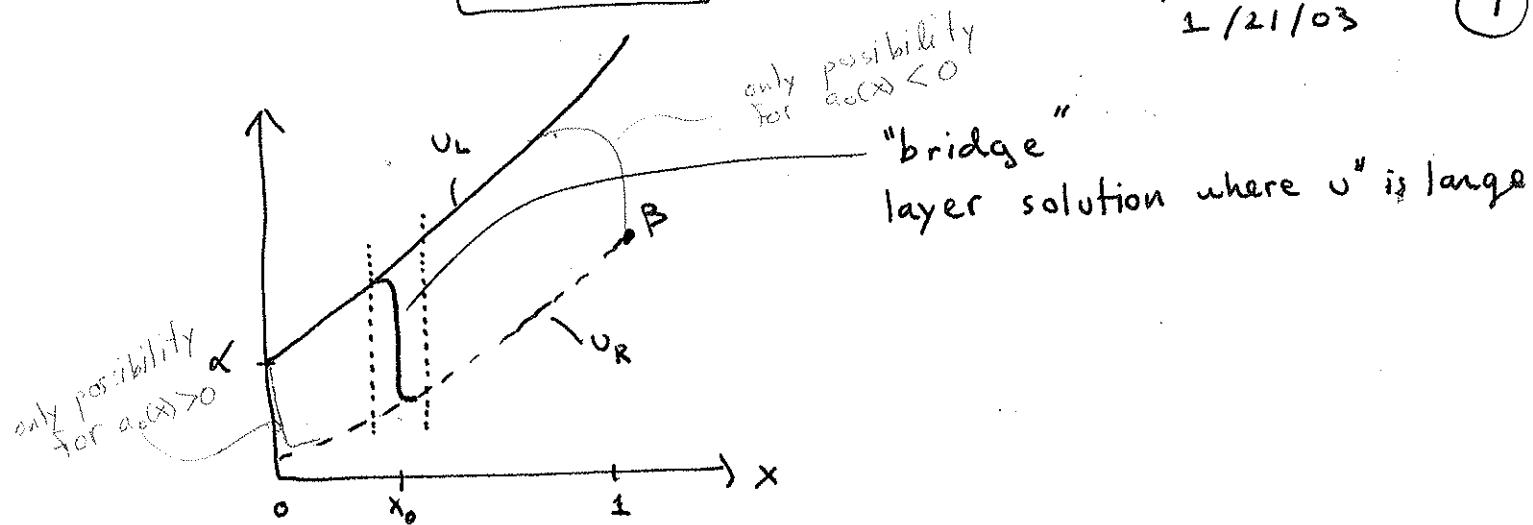
$$v_L = \alpha \exp \left\{ - \int_0^L \frac{b(s)}{a(s)} ds \right\}$$

notice
the negative sign

or if $v(L) = \beta$

$$v_R = \beta \exp \left\{ \int_x^L \frac{b(s)}{a(s)} ds \right\}$$

$$\text{In general, } v_L \neq v_R \text{ unless } \beta = \alpha \exp \left\{ - \int_0^L \frac{b(s)}{a(s)} ds \right\}$$



Consider "layer-type" solutions

$$\xi = \frac{x - x_0}{\varepsilon} \rightarrow x = x_0 + \xi \varepsilon$$

$$\frac{d}{dx} = \frac{d}{d\xi} \cdot \frac{d\xi}{dx} = \frac{1}{\varepsilon} \frac{d}{d\xi}$$

$$\frac{d^2}{dx^2} = \frac{1}{\varepsilon^2} \frac{d^2}{d\xi^2}$$

$$\Rightarrow \varepsilon \left(\frac{1}{\varepsilon^2} \frac{d^2}{d\xi^2} v \right) + a(x_0 + \xi \varepsilon) \frac{1}{\varepsilon} \frac{d}{d\xi} v + b(x_0 + \xi \varepsilon) v = 0$$

third term is small compared to first two terms

Consider $\varepsilon = 0$: $v'' + a(x_0)v' = 0$

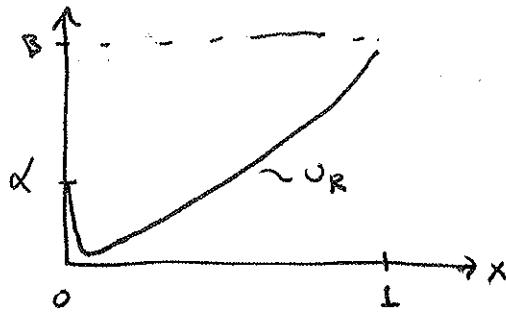
general solution: $v = A + Be^{-a(x_0)\xi}$

solution behavior depends on value of $a(x_0)$

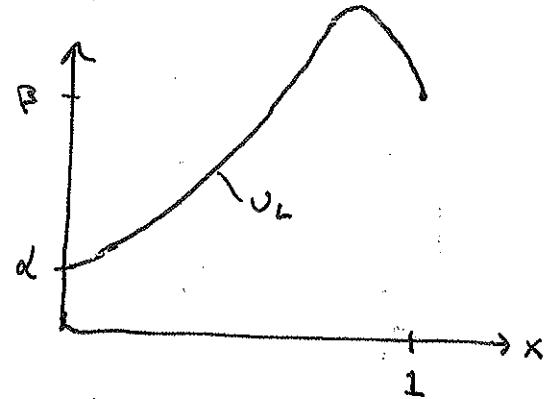
$a(x_0) > 0 \Rightarrow$ decay exponentially as ξ increases

$a(x_0) < 0 \Rightarrow$ decay exponentially as ξ decreases

IF $a > 0$ then $x_0 = 0$



IF $a < 0$ then $x_0 = 1$



No internal layers are possible

Suppose $a_s(x) > 0$

$$u = A + B e^{-ax}$$

$$\text{at } x=0 \quad u = A + B = d$$

$$\text{as } x \rightarrow \infty \quad u = A = u_R(0) = B \exp\left\{\int_0^x \frac{b(s)}{a(s)} ds\right\}$$

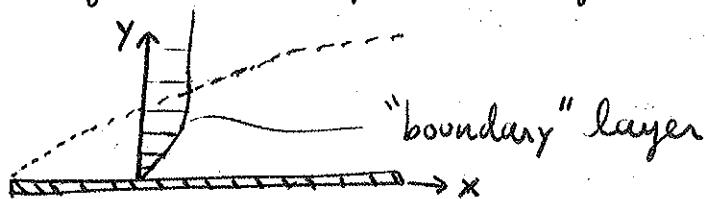
$$\text{then } B = d - A$$

This is the solution, to 1st order of accuracy.

3) Blasius Problem (BOUNDARY LAYER)

Consider the flow of an incompressible fluid over a flat plate.

$$u=1 \quad \overrightarrow{\quad} \quad \overrightarrow{\quad}$$



Eqs for velocity: (u, v)

$$u_x + v_y = 0 \quad (\text{conservation of mass})$$

$$uu_x + vu_y = v u_{yy} \quad (\text{balance of momentum in } x\text{-direction})$$

Boundary conditions

$$u = v = 0 \quad \text{at } y = 0$$

$$u = 1 \quad \text{at } y \rightarrow \infty$$

Reduce the eqns for $u(x,y)$ and $v(x,y)$ to an ODE for a single variable Ψ .

Define $\Psi(x,y)$ such that

$$u = \Psi_y, \quad v = -\Psi_x$$

with this choice, the conservation of mass eqn is satisfied identically. The momentum balance eqn becomes

$$\Psi_y \Psi_{xy} - \Psi_x \Psi_{yy} = v \Psi_{yyy}$$

BCs $\Psi_x = \Psi_y = 0$ on $y=0 \Rightarrow \Psi = 0$ on $y=0$

$$\Psi_y \rightarrow L \text{ as } y \rightarrow \infty$$

Solution for Ψ has the form

$$\Psi(x,y) = \sqrt{2x} f(u), \quad u = \frac{y}{\sqrt{2x}} \quad \text{similarity solution}$$

where u = similarity variable

Need to find a differential equation for $f(u)$

$$\Psi_x = \frac{1}{2} \sqrt{\frac{v}{x}} f(u) + \sqrt{2x} f'(u) \left(-\frac{y x^{-3/2}}{\sqrt{2x}} \right)$$

$$\Psi_x = \frac{\sqrt{v}}{2} f(u) (x^{-1/2}) - \frac{y}{2x} f'(u) \quad y = u \sqrt{2x}$$

$$\Psi_x = \frac{1}{2} \sqrt{\frac{v}{x}} f(u) - \frac{1}{2x} u \sqrt{2x} f'(u)$$

$$\bullet \Psi_x = \frac{1}{2} \sqrt{\frac{v}{x}} (f(u) - u f'(u))$$

$$\Psi_y = \sqrt{2x} f'(u) \frac{1}{\sqrt{2x}} \Rightarrow \bullet \Psi_y = f'(u)$$

$$\Psi_{xy} = \Psi_{yx} = f''(u) \left(-\frac{1}{2} \frac{\sqrt{v}}{\sqrt{2x}} \frac{y}{x^{3/2}} \right) = f''(u) \left(\frac{-u \sqrt{2x}}{2 \sqrt{x^3 v}} \right) = -\frac{u}{2x} f'''(u)$$

$$\Psi_{yy} = \frac{1}{\sqrt{v_x}} f''(n)$$

$$\Psi_{yyy} = \frac{1}{v_x} f'''(n)$$

Substitute partial derivatives into

$$\Psi_y \Psi_{xy} - \Psi_x \Psi_{yy} = v \Psi_{yyy}$$

$$\Rightarrow \frac{1}{2} \left((f'(n))^2 - f(n) f''(n) \right) = f'''(n)$$

Boundary conditions

$$\Psi_y = 0 \text{ on } y=0 \Rightarrow f'(0) = 0$$

$$\Psi_x = 0 \text{ on } y=0 \Rightarrow f(0) = 0$$

$$\Psi_y \rightarrow 1 \text{ as } y \rightarrow \infty \Rightarrow \lim_{n \rightarrow \infty} f'(n) = 1$$

Let $f''(0) = s$, shear stress

pick a value for s , integrate out to infinity, compare $f(n)$ to $\lim_{n \rightarrow \infty} f'(n) = 1$. Adjust s as necessary. This is the shooting method.

Differential Algebraic Equation

$$y' = f(t, y) \quad \text{explicit ODE}$$

$$F(t, y, y') = 0 \quad \text{implicit ODE}$$

If $\frac{\partial F}{\partial y'}$, Jacobian matrix, is non-singular, then the implicit equation can be "solved" for y' . A DAE is one for which $\frac{\partial F}{\partial y'}$ is singular. Then the eqn involves both "differential" variables and "algebraic" variables.

Some cases this is made explicit:

$$\text{DAE} \left\{ \begin{array}{l} y' = f(t, y, z) \\ 0 = g(t, y, z) \end{array} \right.$$

In this case, $y(t)$ = differential variables, $z = z(t)$ = algebraic variables.

One issue is what boundary or initial data is allowed.

eg: $y' = z$ \rightarrow $| \cancel{y(0)} | I = z$
 $0 = y - t \rightarrow y = t$

Hence Here, no data was needed.

IF the equations are simple enough, solve for z in terms of t, y and substitute into y' for ODE.

1/24/03

Example

(1) Mechanics: Let some mechanical system be described by a set of n generalized coordinates.

q_j , $j = 1, \dots, n$ = generalized coordinates

Typically, q_j = positions or possibly angles, areas, etc...

Then define generalized velocities as q'_j , $j = 1, \dots, n$.

Suppose Let $g_i(t, q) = 0$ be a set of m constraints on the motion.

Suppose $U = U(q)$ is the potential energy of the system and $T = T(q, q')$ is the kinetic energy of the system.

1/24/03

A Lagrangian for the system is

$$L = \dot{q} + T - \sum_{i=1}^m \lambda_i q_i, \quad L = L(t)$$

where λ are the Lagrangian Lagrange multipliers.

Hamilton's Principle - The motion of the system from

$t=0$ to $t=b$ is one that minimizes the

integral $I_H = \int_0^b L(t) dt \Rightarrow \underline{\text{Calculus of Variations}}$

\Rightarrow Euler Lagrange equations.

$$\frac{\partial L}{\partial q_j} - \frac{d}{dt} \left(\frac{\partial L}{\partial \dot{q}_j} \right) = 0, \quad j=1, \dots, n$$

\Rightarrow system of n 2nd order ODEs for q_i

let $\underline{v} = \underline{q}'$ be the generalized velocities

$$\Rightarrow \boxed{\underline{M}(t, \underline{q}) \underline{v}' = \underline{f}(t, \underline{q}, \underline{v}) - \underline{G}^T(t, \underline{q}) \lambda}$$

$$\underline{v} = \underline{q}'$$

$$\underline{o} = \underline{q}(t, \underline{q})$$

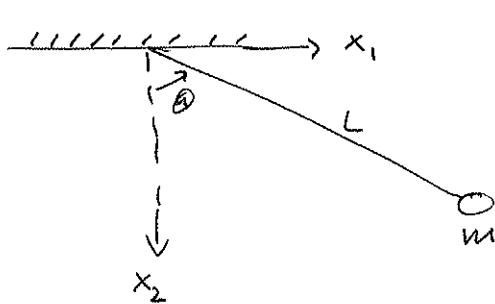
where \underline{v} - generalized velocities

\underline{M} - "mass matrix," (positive definite)

\underline{f} - applied forces

$$\underline{G} = \frac{\partial \underline{q}}{\partial \underline{q}}$$

Now consider a ~~2D~~ case:



(x_1, x_2) - Cartesian coordinates

\hat{g} - gravity

let $q_j = x_j$, $j = 1, 2$

potential energy, $V = -m\hat{g}x_2$ (PE goes to 0 when $x_2 \rightarrow 0$)

kinetic energy, $T = \frac{m}{2}(x_1'^2 + x_2'^2)$

constraint: $g = x_1^2 + x_2^2 - L^2 = 0$

Lagrangian: $L = -V + T - \lambda g$

$$L = \frac{m}{2}(x_1'^2 + x_2'^2) + m\hat{g}x_2 - \lambda(x_1^2 + x_2^2 - L^2) = 0$$

Euler-Lagrange eqns:

$$\frac{\partial L}{\partial x_1} - \frac{d}{dt}\left(\frac{\partial L}{\partial x_1'}\right) = 0 \Rightarrow -2\lambda x_1 - mx_1'' = 0$$

$$\frac{\partial L}{\partial x_2} - \frac{d}{dt}\left(\frac{\partial L}{\partial x_2'}\right) = 0 \Rightarrow m\hat{g} - 2\lambda x_2 - mx_2'' = 0$$

let $x_1' = v_1$ and $x_2' = v_2$

\Rightarrow

$$\boxed{\begin{aligned} mv_1' &= -2\lambda x_1 \\ mv_2' &= m\hat{g} - 2\lambda x_2 \\ x_1^2 + x_2^2 - L^2 &= 0 \end{aligned}}$$

Note, that if you set

$$\begin{aligned} x_1 &= L \sin\theta \\ x_2 &= L \cos\theta \end{aligned} \Rightarrow mL\ddot{\theta} + mg\sin\theta = 0$$

Initial Value Problems

$$y' = f(t, y), \quad t \geq 0, \quad y(0) = c$$

"Problem Stability"

- idea: given a solution $y(t)$ of the IVP, the stability refers to the sensitivity of the solution (output) to perturbations in the data (input)

caution on "stability"

- "stability" is overused. Here stability might better be termed "conditioning," and "stability" referring to a numerical method,
Unfortunately there are many versions of "stability" defined for the numerics.

Test Equation (scalar equation)

$$y' = \lambda y, \quad t \geq 0, \quad y(0) = c$$

Here, λ is a constant, possibly complex. Often λ is an eigenvalue.

solution: $y(t) = y(0) e^{\lambda t}, \quad \lambda = a + ib$

$$\begin{aligned} y(t) &= y(0) e^{at+ibt} \\ &= y(0) e^{at} (\cos bt + i \sin bt) \end{aligned}$$

imaginary part of λ represents oscillatory behavior
and the real part of λ represents exponential growth or decay.

\Rightarrow in terms of stability, we're only interested in the $\operatorname{Re}\lambda$

Suppose $\hat{y}(t)$ solves the test eqn but with $\hat{y}(0)$ as the initial condition.

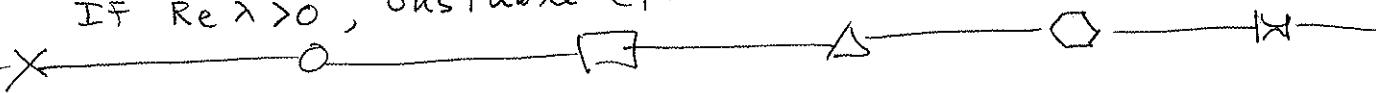
Let $y(t)$ be "the" solution to the test equation and $\tilde{y}(t)$ the perturbed solution, then

$$\begin{aligned} |\hat{y}(t) - y(t)| &= |y(\hat{0}) e^{\lambda t} - y(0) e^{\lambda t}| \\ &= e^{(\operatorname{Re}\lambda)t} |y(\hat{0}) - y(0)| \end{aligned}$$

IF $\operatorname{Re}\lambda \leq 0$, then $|\hat{y}(t) - y(t)|$ remains bounded, ie stable.

IF $\operatorname{Re}\lambda < 0$, then $|\hat{y}(t) - y(t)|$ approaches zero \Rightarrow asymptotically stable.

IF $\operatorname{Re}\lambda > 0$, unstable (perturbations grow).



Full Equation

$$y' = f(t, y), \quad y(0) = \underline{c}, \quad t \geq 0$$

The idea is similar but the details a little blander.

The idea is similar but the details a little blander.

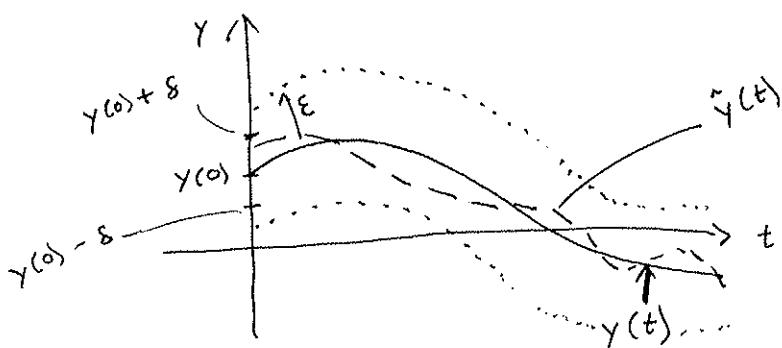
Suppose $y(t)$ is the solution of the ODE for $t \geq 0$.
The solution is stable if given any $\epsilon > 0$ there is a $\delta > 0$ st any other solution $\tilde{y}(t)$ satisfying the ODE and

$$\|\tilde{y}(0) - y(0)\| < \delta$$

also satisfies $\|\tilde{y}(t) - y(t)\| < \epsilon$ for all $t \geq 0$.

The solution is asymptotically stable if in addition to being stable

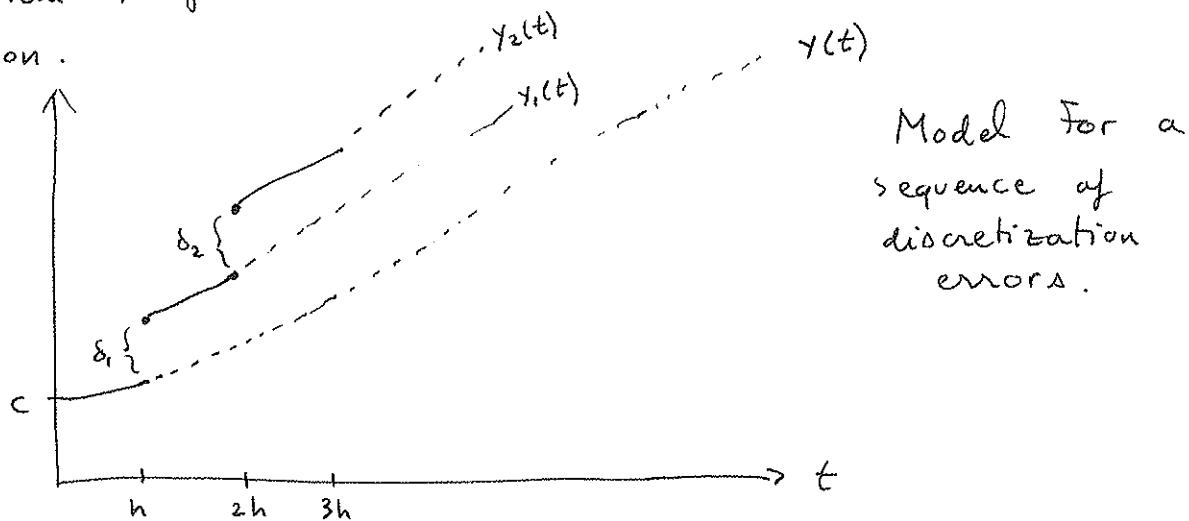
$$\|\tilde{y}(t) - y(t)\| \rightarrow 0 \text{ as } t \rightarrow \infty.$$



Example : Consider the effect of a sequence of perturbations on the solution of the test equation.

We have $y' = \lambda y$, $y(0) = c$, $t \geq 0$

Suppose the solution is perturbed at $t=h$ by an amount δ_1 . The perturbed solution is then itself perturbed at $t=2h$ by δ_2 and so on.



$$y(t) = c e^{\lambda t}$$

$$\begin{aligned} y_1(t) &= (y(h) + \delta_1) e^{\lambda(t-h)} \\ &= (c e^{\lambda h} + \delta_1) e^{\lambda(t-h)} = c e^{\lambda t} + \delta_1 e^{\lambda(t-h)} \end{aligned}$$

$$\begin{aligned} y_2(t) &= (y_1(2h) + \delta_2) e^{\lambda(t-2h)} = (c e^{\lambda 2h} + \delta_1 e^{\lambda h} + \delta_2) e^{\lambda(t-2h)} \\ &= c e^{\lambda t} + \delta_1 e^{\lambda(t-h)} + \delta_2 e^{\lambda(t-2h)} \end{aligned}$$

$$\boxed{y_j(t) = c e^{\lambda t} + \sum_{k=1}^j \delta_k e^{\lambda(t-kh)}} \quad \text{perturbed solution}$$

$$\text{Difference: } y_j(t) - y(t) = \sum_{k=1}^j \delta_k e^{\lambda(t-kh)}$$

Behavior depends on λ :

IF $\operatorname{Re} \lambda \leq 0$, then the perturbations only accumulate linearly with no exponential growth in time (at most).

IF $\operatorname{Re} \lambda < 0 \Rightarrow$ asymptotically stable.

Linear Systems

1) constant coefficient, homogeneous problems,

$\underline{y}' = \underline{A} \underline{y}, \quad t \geq 0, \quad \underline{A}$ is $m \times m$ matrix (real).

$$\boxed{\underline{y}(t) = e^{\underline{A}t} \underline{y}(0)} \quad \text{notice, } e^{\underline{A}t} \text{ is a matrix}$$

$e^{\underline{A}t}$ is defined via power series:

$$e^{\underline{A}t} = I + \underline{A}t + \frac{1}{2}\underline{A}^2t^2 + \frac{1}{3!}\underline{A}^3t^3 + \dots$$

Note, if \underline{A} is diagonalizable then $\underline{T}^{-1}\underline{A}\underline{T} = \underline{\Lambda} = \operatorname{diag}(\lambda_i)_{i=1,\dots,m}$ and column vectors of \underline{T} are the (right) eigenvectors of \underline{A} .

$$\text{IF } \underline{T} \text{ s.t. } \underline{T}^{-1}\underline{A}\underline{T} = \underline{\Lambda} \Rightarrow \underline{A} = \underline{T}\underline{\Lambda}\underline{T}^{-1}$$

$$\underline{A}^2 = (\underline{T}\underline{\Lambda}\underline{T}^{-1})(\underline{T}\underline{\Lambda}\underline{T}^{-1}) = \underline{T}\underline{\Lambda}^2\underline{T}^{-1}$$

$$\Rightarrow e^{\underline{A}t} = \underline{T}(\underline{I} + \underline{\Lambda}t + \frac{1}{2}\underline{\Lambda}^2t^2 + \dots)\underline{T}^{-1}$$

$$e^{\underline{A}t} = \underline{T} \begin{bmatrix} e^{\lambda_1 t} & & \\ & e^{\lambda_2 t} & \\ & \ddots & \\ & & e^{\lambda_m t} \end{bmatrix} \underline{T}^{-1}$$

$$\underline{y}(t) = \underline{T} \begin{pmatrix} e^{\lambda_1 t} \\ \vdots \\ e^{\lambda_m t} \end{pmatrix} \stackrel{T^{-1}}{=} \underline{y}(0)$$

$$\Rightarrow \underline{T}^{-1} \underline{y}(t) = \begin{pmatrix} e^{\lambda_1 t} \\ \vdots \\ e^{\lambda_m t} \end{pmatrix} \stackrel{T^{-1}}{=} \underline{y}(0)$$

If $\underline{w}(t) = \underline{T}^{-1} \underline{y}(t)$ then

$$\underline{w}(t) = \begin{pmatrix} e^{\lambda_1 t} \\ \vdots \\ e^{\lambda_m t} \end{pmatrix} \underline{w}(0) \rightarrow w_j(t) = w_j(0) e^{\lambda_j t}$$

For diagonalizable \underline{A} , if $\operatorname{Re} \lambda_j \leq 0 \forall j \Rightarrow$ problem stability.

if $\operatorname{Re} \lambda_j < 0 \forall j \Rightarrow$ asymptotic stability.

1/28/03

If the difference between \hat{y} and y remains bounded, then the IVP is stable.

If the difference is bounded and approaches zero as $t \rightarrow \infty$, then the IVP is asymptotically stable.

The IVP is unstable if the difference is unbounded.

Suppose \underline{A} is not diagonalizable.

Then \exists T matrix nonsingular, satisfying

$$\underline{T}^{-1} \underline{A} \underline{T} = \begin{bmatrix} \underline{\Lambda}_1 & & \\ & \underline{\Lambda}_2 & \\ & & \ddots & \underline{\Lambda}_n \end{bmatrix}, \text{ being in Jordan Canonical form}$$

where $\begin{cases} \underline{\Lambda}_i & \text{if } \underline{\Lambda}_i = \begin{bmatrix} \lambda_i & 1 & & 0 \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ 0 & & & \lambda_i \end{bmatrix} \text{ — Jordan block} \\ \underline{\Lambda}_i & \text{if } \end{cases}$

For a diagonalizable matrix Λ_i is one by one

$$\underline{y}(t) = \underline{T} \begin{bmatrix} \underline{\Lambda}_1 & & \\ & \underline{\Lambda}_2 & \\ & & \ddots & \underline{\Lambda}_n \end{bmatrix} \underline{T}^{-1} \underline{y}(0)$$

The question of stability depends on behaviour of term $e^{\underline{\Lambda}_i t}$. The corresponding ODE is

$$y_1' = \lambda y_1 + y_2$$

$$y_2' = \lambda y_2 + y_3$$

 \vdots

$$y_{n-1}' = \lambda y_{n-1} + y_n$$

$$y_n' = \lambda y_n$$

For example, suppose $n=2$ then

$$y_1' = \lambda y_1 + y_2, \quad y_1(0) = \alpha, \quad y_2(0) = \beta$$

$$y_2' = \lambda y_2$$

$$\rightarrow y_2 = \beta e^{\lambda t} \rightarrow y_1' = \lambda y_1 + \beta e^{\lambda t}$$

$$\downarrow$$

$$y_1 = (\lambda + \beta t) e^{\lambda t}$$

$$y_1 = (\alpha + \beta t) e^{\lambda t}$$

$$y_2 = \beta e^{\lambda t}$$

- 1) stable if all eigenvalues of \underline{A} satisfy either
 $\operatorname{Re} \lambda < 0$ or $\operatorname{Re} \lambda = 0$ and λ is simple
(λ is simple - belongs to 1×1 Jordan block)
- 2) asymptotic stability if $\operatorname{Re} \lambda < 0$ for
all eigenvalues of A .

Variable Coefficient, Non-homogeneous

$$\underline{y}'(t) = \underline{A}(t) \underline{y}(t) + \underline{q}(t), \quad t \geq 0$$

where \underline{A} is an $m \times m$ real matrix function
and \underline{q} is an m -vector, real function, both
smooth.

The exact solution is found in terms of the
Fundamental solution matrix: $\underline{Y}' = \underline{A}(t) \underline{Y}, t \geq 0$
where $\underline{Y}(0) = \mathbb{I}$. We assume that $\underline{Y}(t)$ is
known.

Set $\underline{y}(t) = \underline{Y}(t) \underline{g}(t)$ vector function to be determined

This approach is variation of parameters.

$$\rightarrow \underline{y}'(t) = \underline{Y}'(t) \underline{g}(t) + \underline{Y}(t) \underline{g}'(t) = \underline{A} \underline{Y} \underline{g} + \underline{Y} \underline{g}' = \underline{A} \underline{y} + \underline{Y} \underline{g}'$$

$$\Rightarrow \underline{A} \underline{y} + \underline{Y} \underline{g}' = \underline{A} \underline{y} + \underline{q} \Rightarrow \underline{Y} \underline{g}' = \underline{q}$$

$$\underline{g}' = \underline{\Upsilon}^{-1} \underline{g}$$

$$\Rightarrow \underline{g}(t) = \int_0^t \underline{\Upsilon}^{-1}(s) \underline{q}(s) ds + \underline{c}$$

$$\underline{y}(t) = \underline{\Upsilon}(t) \left[\underline{c} + \int_0^t \underline{\Upsilon}^{-1}(s) \underline{q}(s) ds \right]$$

If a perturbed solution $\hat{\underline{y}}(t)$ solves the same ODE but with initial conditions $\hat{\underline{y}}(0) = \hat{\underline{c}}$, then the difference is bounded if $\underline{\Upsilon}(t)$ is bounded.

If a perturbed solution $\underline{g}(t)$ solves a perturbed ODE as well, then require $\underline{\Upsilon}(t) \underline{\Upsilon}^{-1}(t)$

Asymptotic stability if $\underline{\Upsilon} \rightarrow 0$ as $t \rightarrow \infty$.

Note that in the variable coefficient case, a solution may grow over certain intervals of time but still be stable or asymptotically stable.

$$\frac{dy}{dt} = \cos t y(t) \rightarrow y(t) = e^{\sin t}$$

The problem is stable (clearly $e^{\sin t}$ is a bounded function)
but NOT asymptotically stable.

Nonlinear Case

$$\underline{y}' = \underline{f}(t, \underline{y}), \quad t \geq 0, \quad \underline{y}(0) = \underline{c}$$

IF $\underline{y}(t)$ is a known trajectory, then is it stable?

One approach is to linearize about $\underline{y}(t)$. Consider a perturbed solution $\hat{\underline{y}}(t)$ solving

$$\hat{\underline{y}}' = \underline{f}(t, \hat{\underline{y}}), \quad \hat{\underline{y}}(0) = \hat{\underline{c}}$$

Expand about $\underline{y}(t)$:

$$\hat{\underline{y}}' = \underline{f}(t, \underline{y}) + \underbrace{\frac{\partial \underline{f}}{\partial \underline{y}}(t, \underline{y})}_{\underline{A}(t)} (\hat{\underline{y}} - \underline{y}) + \dots \text{ higher order terms}$$

IF $\underline{w}(t) = \hat{\underline{y}} - \underline{y}$, then $\underline{w}(t)$ solves

$$\underline{w}'(t) = \underline{A}(t) \underline{w} + \dots$$

IF \underline{w} is small then \underline{w} solves

$$\underline{w}'(t) = \underline{A}(t) \underline{w}(t)$$

to a first approximation.

Can perform an analysis of the linearized problem and this could provide necessary conditions for the nonlinear problem.

Hamiltonian Systems

A Hamiltonian system is a set of ODEs of the

$$\text{form } q_i' = \frac{\partial H}{\partial p_i}, \quad p_i' = -\frac{\partial H}{\partial q_i}, \quad i = 1, 2, \dots, l$$

This is a system of $2l$ equations for the variables (q_i, p_i) , $i = 1, \dots, l$, depending on time t . H is a scalar function, $H = H(q, p)$, and is a smooth function called the Hamiltonian.

In a mechanical system, q = position of a set of particles and p = momentum of particles, and H is the sum of potential and kinetic energies. (The total energy of the system).

$$\frac{dH}{dt} = \frac{\partial H}{\partial p} \cdot \frac{dp}{dt} + \frac{\partial H}{\partial q} \cdot \frac{dq}{dt}$$

$$\downarrow \quad = \frac{\partial H}{\partial p} \cdot \left(-\frac{\partial H}{\partial q} \right) + \frac{\partial H}{\partial q} \cdot \left(\frac{\partial H}{\partial p} \right)$$

$$\downarrow \quad = 0 \quad \Rightarrow \quad H = \text{constant} \Rightarrow \text{total energy is conserved}$$

A geometric implication of this is that the flow conserves area.

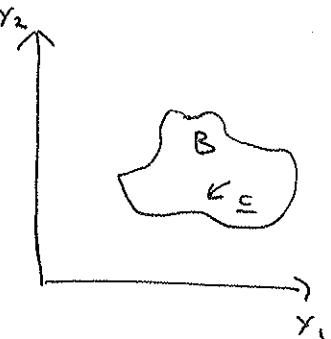
Consider a 2-vector problem

$$\underline{y}'(t) = \Sigma(\underline{y}) \quad (\text{autonomous system})$$

$$\underline{y} = \begin{pmatrix} x_1 \\ y_2 \end{pmatrix}, \quad \underline{y}(0) = \underline{c} = \begin{pmatrix} c_1 \\ c_2 \end{pmatrix}$$

Assume a solution is given by $\underline{y}(t; \underline{c})$.

Consider the set B and the flow of B under the trajectories $\underline{y}(t; \underline{c})$ where, $\underline{c} \in B$.



Define $S(t)B = \{\underline{y}(t; \underline{c}), \underline{c} \in B\}$

If \underline{y} is governed by a constant coefficient linear system then the behavior of the region depends on the stability of the equations. If the problem is asymptotically stable, for example, then the region shrinks to zero. If the $\nabla \cdot \underline{F} = \frac{\partial F_1}{\partial x_1} + \frac{\partial F_2}{\partial x_2} = 0$ then the area remains fixed. (The area of $S(t)B$ remains fixed).

For the Hamiltonian system with $\lambda=1$,

$$F_1 = \frac{\partial H}{\partial p}, \quad F_2 = -\frac{\partial H}{\partial q}, \quad \underline{y} = \begin{bmatrix} q \\ p \end{bmatrix}$$

$$\nabla \cdot \underline{F} = \frac{\partial}{\partial q} \left(\frac{\partial H}{\partial p} \right) + \frac{\partial}{\partial p} \left(-\frac{\partial H}{\partial q} \right) = 0$$

\Rightarrow the area of the flow is preserved

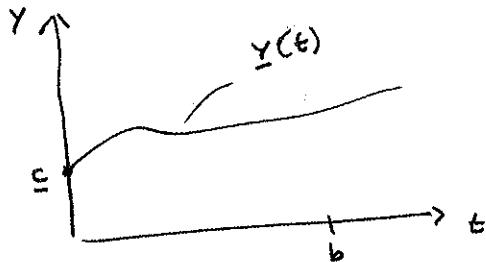
For $\lambda=1$ it turns out that the area is preserved in every (q_i, p_i) plane. The flow is said to be symplectic.

The implication for problem stability is that a Hamiltonian system cannot be asymptotically stable or unstable. A numerical solution may want to mimic this area preserving property \Rightarrow symplectic numerical methods.

Numerical Methods for IVPs

$$\underline{y}' = \underline{f}(t, \underline{y}), \quad t \geq 0, \quad \underline{y}(t) = \underline{c}$$

the solution is a trajectory $\underline{y}(t)$:



Approximate the solution on a mesh or grid.

$$0 = t_0 < t_1 < \dots < t_N = b$$

$$\text{mesh spacing: } h_n = t_n - t_{n-1} > 0$$

The solution is approximated by a set of vectors $\underline{y}_n \in \mathbb{R}^m$ st
 $\underline{y}_n \approx \underline{y}(t_n)$ for $n=0, \dots, N$.

Basic Concepts (consistency, stability, convergence)

I illustrate the basic concepts for Euler's method:

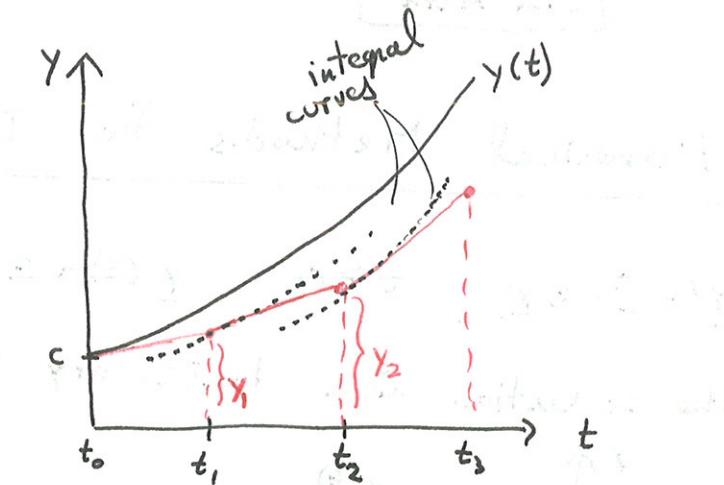
$$\begin{aligned} \underline{y}(t_n) &= \underline{y}(t_{n-1}) + \underline{f}'(t_{n-1}) \frac{(t_n - t_{n-1})}{\cancel{(t_n - t_{n-1})}} + \frac{1}{2} \underline{f}''(t_{n-1}) \frac{\cancel{(t_n - t_{n-1})}^2}{\cancel{(t_n - t_{n-1})}^2} + \dots \\ &= \underline{y}(t_{n-1}) + \underline{f}(t_{n-1}, \underline{y}(t_{n-1})) h_n + O(h_n^2) + \dots \end{aligned}$$

Euler's method ignores $O(h_n^2)$ terms (and higher) to write

$$\boxed{\underline{y}_n = \underline{y}_{n-1} + h_n \underline{f}(t_{n-1}, \underline{y}_{n-1}), \quad \underline{y}_0 = \underline{c}},$$

Geometry

Move along successive tangents to the integral curves.



the smaller the mesh spacing, the more accurate the solution

BASIC CONCEPTS

1/31/03

Euler's Method: $y_n = y_{n-1} + h f(t_{n-1}, y_{n-1})$

$$\frac{y_n - y_{n-1}}{h} = f(t_{n-1}, y_{n-1})$$

Define the difference operator N_h

$$N_h v(t_n) = \frac{v(t_n) - v(t_{n-1})}{h} - f(t_{n-1}, y_{n-1})$$

N_h acts on a grid function $v(t_n)$

N_h is an approximation of N , a differential operator,

where $N_y(t) = y' - f(t, y)$

The local truncation error measures the amount by which the exact solution of the differential equation fails to satisfy the difference eqn.

Define d_n as the local truncation error and set $d_n = N_h y(t_n)$. A numerical method is consistent if $d_n \rightarrow 0$ as $h_n \rightarrow 0$.

A numerical method is consistent

of order P if $|d_n| = O(h_n^P)$ as $h_n \rightarrow 0$

(The numerical method is P^{th} -order accurate.)

$$|d_n| = O(h_n^P) \text{ as } h_n \rightarrow 0 \Rightarrow \lim_{h_n \rightarrow 0} \frac{|d_n|}{h_n^P} = \text{const.}$$

For Euler's Method

$$N_n y(t_n) = \frac{y(t_n) - y(t_{n-1})}{h_n} - f(t_{n-1}, y(t_{n-1}))$$

$$N_n y(t_n) = \frac{(y(t_{n-1}) + h_n y'(t_{n-1}) + \frac{h_n^2}{2} y''(t_{n-1}) + \dots) - y(t_{n-1})}{h_n} - f(t_{n-1}, y(t_{n-1}))$$

$$N_n y(t_n) = \frac{h_n}{2} y''(t_{n-1}) + \dots$$

dominant behavior of d_n

hence, Euler's method is first-order accurate

If a numerical method is consistent then $N_n \rightarrow N$ as $h \rightarrow 0$, but what does this say about the error, $e_n = y_n - y(t_n)$?

Define error, $e_n = y_n - y(t_n)$

A numerical method is convergent if

$\max_{0 \leq n \leq N} |e_n| \rightarrow 0$ as $h \rightarrow 0$ where $h = \max_n h_n$.

A numerical method is convergent of order P if

$\max_{0 \leq n \leq N} |e_n| = O(h^P)$ as $h \rightarrow 0$

Typically, if a numerical method is consistent of order p then it is convergent of order p . To ensure this you must show that the method is stable.

The difference method is "O-stable" (or just "stable") if $h_0 > 0$ and $K > 0$ exist st

$$|x_n - z_n| \leq K \left[|x_0 - z_0| + \max_{1 \leq j \leq N} |N_h x_j - N_h z_j| \right] \text{ for } 0 \leq n \leq N$$

for any mesh functions x_n, z_n and $h \leq h_0$

(Recall our discussion of the regularity of the differential equation - this is similar.)

Thm: If a difference method is consistent of order p with local truncation error d_n and if it is O-stable with constant K , then it is convergent of order p st $|e_n| \leq K \max_{1 \leq j \leq N} |d_j| = O(h^p)$

whenever you talk of numerical methods for differential eqns, the usual theorem is

consistency + stability \Rightarrow convergence

Proof If a method is O-stable then

$$|x_n - z_n| \leq K \left[|x_0 - z_0| + \max_{1 \leq j \leq N} |N_h x_j - N_h z_j| \right]$$

set $x_n = y_n$, set $z_n = y(t_n)$

$$|e_n| \leq K \left[|y_0| + \max_{1 \leq j \leq N} |0 - d_j| \right]$$

$$\Rightarrow |e_n| \leq K \max_{1 \leq j \leq N} |d_j|$$

Examine 0-stability for Euler

$$\text{Set } s_n = x_n - z_n, \quad B = \max_j |N_n x_j - N_n z_j|$$

$$\text{use } N_n x_j = \frac{x_j - x_{j-1}}{h_j} - f(t_{j-1}, x_{j-1})$$

$$N_n z_j = \frac{z_j - z_{j-1}}{h_j} - f(t_{j-1}, z_{j-1})$$

$$\therefore B \geq |N_n x_n - N_n z_n| \quad \text{for } 1 \leq n \leq N$$

$$\Rightarrow B \geq \left| \frac{s_n - s_{n-1}}{h_n} + f(t_{n-1}, x_{n-1}) - f(t_{n-1}, z_{n-1}) \right|$$

$$\geq \left| \frac{s_n}{h_n} - \left(\frac{s_{n-1}}{h_n} + f(t_{n-1}, x_{n-1}) - f(t_{n-1}, z_{n-1}) \right) \right|$$

$$\geq \left| \frac{s_n}{h_n} \right| - \left| \frac{s_{n-1}}{h_n} + f(t_{n-1}, x_{n-1}) - f(t_{n-1}, z_{n-1}) \right|$$

$$\Rightarrow \left| \frac{s_n}{h_n} \right| \leq \left| \frac{s_{n-1}}{h_n} + f(t_{n-1}, x_{n-1}) - f(t_{n-1}, z_{n-1}) \right| + B$$

multiply by h_n , use triangle inequality

$$\Rightarrow |s_n| \leq |s_{n-1}| + h_n |f(t_{n-1}, x_{n-1}) - f(t_{n-1}, z_{n-1})| + h_n B$$

Key step: assume f satisfies Lipschitz condition

$$|f(t_{n-1}, x_{n-1}) - f(t_{n-1}, z_{n-1})| \leq L |s_{n-1}|$$

$$\Rightarrow |s_n| \leq (1 + h_n L) |s_{n-1}| + h_n B$$

$$\leq (1 + h_n L) [(1 + h_{n-1} L) |s_{n-2}| + h_{n-1} B] + h_n B$$

$$= (1 + h_n L)(1 + h_{n-1} L) |s_{n-2}| + (h_n + (1 + h_n L) h_{n-1}) B$$

remember triangle inequality

$$\begin{aligned} |a+b| &\leq |a| + |b| \\ |a+b+c| &\leq |a| + |b| + |c| \\ ||a+b|| &\leq ||a|| + ||b|| \end{aligned}$$

continue with this process:

$$|s_n| \leq (1+h_n L)(1+h_{n-1} L) \dots (1+h_1 L) |s_0| + \\ \left[h_n + (1+h_n L) h_{n-1} + (1+h_n L)(1+h_{n-1} L) h_{n-2} + \dots + \right. \\ \left. (1+h_n L)(1+h_{n-1} L) \dots (1+h_2 L) h_1 \right] B$$

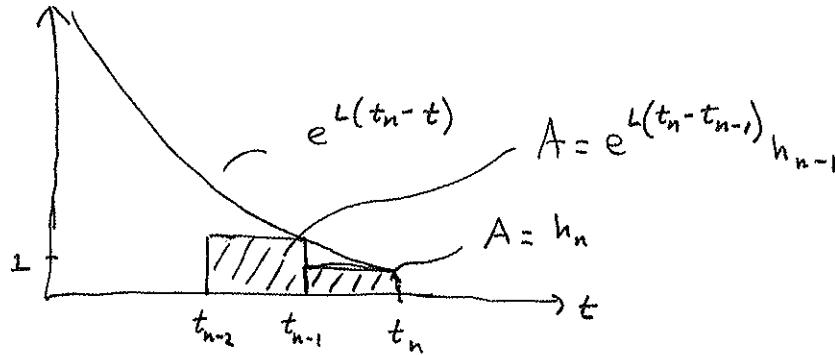
use $1+hL \leq e^{Lh}$

Note $(1+h_n L)(1+h_{n-1} L) \dots (1+h_2 L) L \leq e^{\underset{1}{L}(h_1 + \dots + h_n)} = e^{L t_n}$

$$(1+h_n L)(1+h_{n-1} L) \dots (1+h_{j+1} L) L \leq e^{\underset{j+1}{L}(h_{j+1} + \dots + h_n)} = e^{L(t_n - t_j)}$$

$$\Rightarrow |s_n| \leq e^{L t_n} |s_0| + \left[h_n + e^{L(t_n - t_{n-1})} h_{n-1} + e^{L(t_n - t_{n-2})} h_{n-2} \right. \\ \left. + \dots + e^{L(t_n - t_1)} h_1 \right] B$$

Notice :



$$\Rightarrow \left[h_n + e^{L(t_n - t_{n-1})} h_{n-1} + e^{L(t_n - t_{n-2})} h_{n-2} + \dots + e^{L(t_n - t_1)} h_1 \right] \leq \int_0^{t_n} e^{L(t_n - t)} dt \\ = \frac{1}{L} (e^{L t_n} - 1)$$

$$\Rightarrow |s_n| \leq e^{L t_n} |s_0| + \frac{B}{L} (e^{L t_n} - 1)$$

We want to show that $|s_n| \leq K(|s_0| + B)$

$$|s_n| \leq e^{Lt_n} |s_0| + \frac{B}{L} (e^{Lt_n} - 1)$$

$$\text{Let } K = \max\left\{e^{Lb}, \frac{1}{L}(e^{Lb} - 1)\right\}$$

$$\text{then } |s_n| \leq K(|s_0| + B)$$

Remarks :

- i) Recall the error bound $|e_n| \leq K \max_{1 \leq j \leq N} |d_j|$

where $e_n = y_n - y(t_n)$, $d_j = N_h y(t_j) = \text{truncation error}$

For Euler's method,

$$d_j = \frac{h_j}{2} y''(t_{j-1}) + O(h_j^2)$$

$$\Rightarrow |e_n| \leq \frac{K}{2} \max_{1 \leq j \leq N} |h_j y''(\underset{\substack{\downarrow \\ t_{j-1}}}{\tilde{t}_j})| \quad , \quad t_{j-1} \leq \tilde{t}_j \leq t_j$$

$$\text{where } K = \max\left(e^{Lb}, \frac{1}{L}(e^{Lb} - 1)\right)$$

A direct analysis of the error gives the bound

$$\text{that } |e_n| \leq \frac{Mh}{2L} (e^{Lb} - 1) \text{ where } M = \max_{0 \leq t \leq b} |y''(t)|, h = \max_{1 \leq j \leq N} h_j$$

- 2) The above bound is a very poor estimate for the error.
 One reason is that the bound produces exponential growth in the error, which is not necessarily the case. There are other ways to estimate the error for finite values of h . The purpose of the bound above is to establish the behavior as $h \rightarrow 0$.

(3) The stability analysis for Euler's method may be extended to explicit single-step methods of the form

$$N_h y_n = \frac{y_n - y_{n-1}}{h_n} - \phi(t_{n-1}, y_{n-1}, h_n) = 0$$

where ϕ is the stepping function.
(this includes Runge Kutta methods).

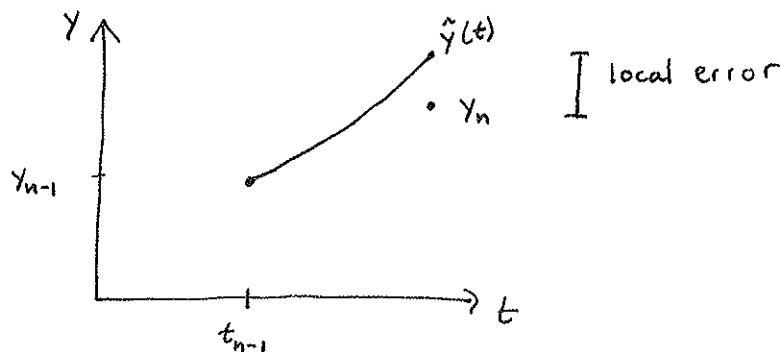
Require that $\phi(t, y, h)$ satisfy a Lipschitz condition in the variable y .

Local Error

This is a measure of discretization error and is NOT the local truncation, though it is related. The local error is a measure of the error committed by one step of the difference formula.

Consider the IVP: $\hat{y}' = f(t, \hat{y})$

$$\hat{y}(t_{n-1}) = y_{n-1}$$



The local error is related to the local truncation error.

For Euler's method

$$\begin{aligned}
 l_n &= \hat{y}(t_n) - (y_{n-1} + h_n f(t_{n-1}, y_{n-1})) \\
 &= \hat{y}(t_n) - [\hat{y}(t_{n-1}) + h_n f(t_{n-1}, \hat{y}(t_{n-1}))] \\
 &= h_n \left[\frac{\hat{y}(t_n) - \hat{y}(t_{n-1})}{h_n} - f(t_{n-1}, \hat{y}(t_{n-1})) \right] \\
 &= h_n N_n \hat{y}(t_n)
 \end{aligned}$$

almost d_n

It can be shown that $l_n = h_n d_n (1 + O(h))$, $h = \max h_n$

If the numerical method is consistent of order p then
 $|l_n| = O(h^{p+1})$ as $h \rightarrow 0$

2/4/03

Absolute Stability

So far we have introduced concepts that concern behavior of a numerical method in the limit $h \rightarrow 0$,

e.g. a method is consistent of order p if

$$d_n = N_n y(t_n) = \text{truncation error} = O(h^p) \text{ as } h \rightarrow 0$$

and if the method is O -stable then it is convergent of order p , ie

$$e_n = y_n - y(t_n) = O(h^p) \text{ as } h \rightarrow 0$$

If we make h sufficiently small then we can achieve a desired accuracy. The question is, how small does h have to be to get an acceptable numerical solution.

Consider the test equation: $y' = \lambda y$, $t \geq 0$, $y(0) = c$
 where λ - complex constant. Apply Euler's method:

$$y_n = y_{n-1} + h_n f(t_{n-1}, y_{n-1}), \quad y_0 = c$$

$$y_n = y_{n-1} + h_n \lambda y_{n-1} \Rightarrow y_n = (1 + h_n \lambda) y_{n-1}$$

Suppose h is fixed, then

$$y_n = (1 + h\lambda) y_{n-1}$$

$$\Rightarrow y_n = (1 + h\lambda)^n c \quad - \text{numerical solution}$$

$$y(t) = c e^{\lambda t} \quad - \text{exact solution}$$

How does the numerical and exact solutions compare for a fixed h ?

Cases:

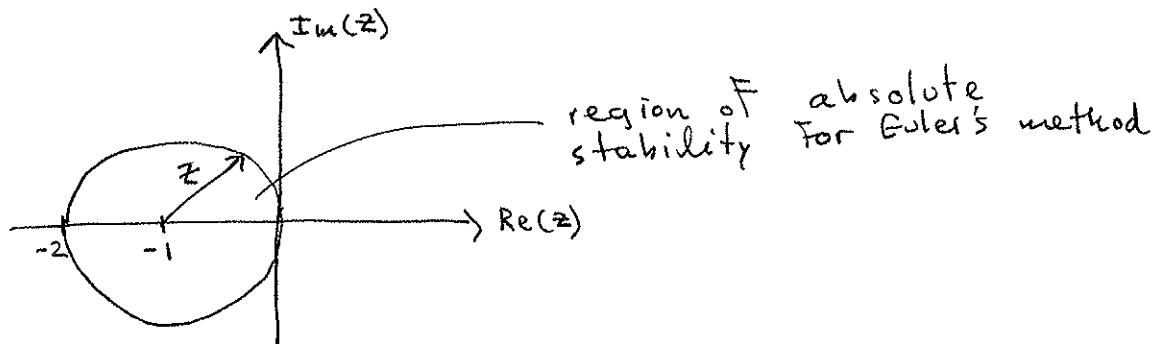
1) $\operatorname{Re}\lambda > 0$. For the unstable case, the exact solution grows exponentially and $(1 + h\lambda)^n$ grows with n . Qualitatively, the solutions have the same growth behavior. Absolute error grows exponentially but the relative error might be ok.

2) $\operatorname{Re}\lambda \leq 0$. The exact solution is bounded and decays exponentially if $\operatorname{Re}\lambda < 0$. The quantity $(1 + h\lambda)$ only decays for a certain h . Therefore we require

$$|1 + h\lambda| \leq 1$$

For a bounded numerical solution. This requirement implies a region of absolute stability.

Let $z = h\lambda$ = complex. Describe the region in the z -plane such that $|1+z| \leq 1$



Suppose that λ is real and negative then we're sitting somewhere on negative real axis, not necessarily restricted to region of absolute stability. For a decaying numerical solution, we require

$$\cancel{-2 < \lambda h < 0} \quad -2 < \lambda h < 0 \Rightarrow h < \frac{2}{-\lambda}$$

This implies that as $-\lambda$ becomes large, h must become small so that the decay is captured.

Extend the analysis to constant coefficient systems.

Constant Coefficient Systems

$$\underline{y}' = \underline{\underline{A}} \underline{y}, \quad t \geq 0, \quad \underline{y}(0) = \underline{c}$$

Apply Euler's method with fixed h .

$$\underline{y}_n = \underline{y}_{n-1} + h \underline{\underline{A}} \underline{y}_{n-1} \Rightarrow \underline{y}_n = (\underline{\underline{I}} + h \underline{\underline{A}}) \underline{y}_{n-1}$$

If $\underline{\underline{A}}$ is diagonalizable, then set $\underline{w}_n = \underline{\underline{T}}^{-1} \underline{y}_n$

$$\Rightarrow \underline{\underline{T}}^{-1} \underline{y}_n = \underline{\underline{T}}^{-1} (\underline{\underline{I}} + h \underline{\underline{A}}) \underline{\underline{T}} \underline{\underline{T}}^{-1} \underline{y}_{n-1}$$

$$\Rightarrow \underline{w}_n = (\underline{\underline{I}} + h \underline{\underline{T}}^{-1} \underline{\underline{A}} \underline{\underline{T}}^{-1}) \underline{w}_{n-1}$$

$$\boxed{\underline{w}_n = (\underline{\underline{I}} + h \underline{\underline{A}}) \underline{w}_{n-1}}$$

$$\underline{w}_n = (\underline{\underline{I}} + h \underline{\underline{A}}) \underline{w}_{n-1}$$

For a bounded solution, require $|1 + h\lambda_j| \leq 1$

For $j=1, \dots, m$ where λ_j are the eigenvalues of $\underline{\underline{A}}$.

Example: spring equation: $m u'' + nu' + ku = 0, t \geq 0$

m -mass, n -damping coeff, $u(0) = \alpha, u'(0) = \beta$

k -spring constant: $m, n, k > 0$

The general solution is: $u(t) = Ae^{\lambda_1 t} + Be^{\lambda_2 t}$

where λ_1, λ_2 are roots of the characteristic polynomial: $m\lambda^2 + n\lambda + k = 0$

$$\lambda = \frac{-n \pm \sqrt{n^2 - 4mk}}{2m}$$

since m, n, k are positive, the real parts of both eigenvalues λ_1, λ_2 are negative \Rightarrow the problem is asymptotically stable.

Consider the behavior of Euler's method for two cases:

1) Overdamped: $n^2 - 4mk > 0 \Rightarrow$ both roots λ_1, λ_2 are real

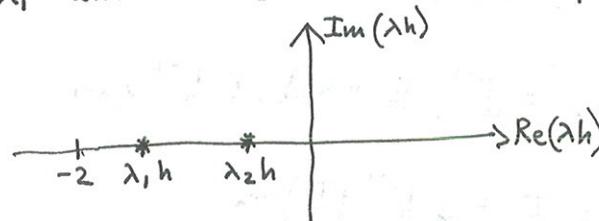
the solution is $u = Ae^{\lambda_1 t} + Be^{\lambda_2 t}$

$$\text{where } \lambda_1 = \frac{-n - \sqrt{n^2 - 4mk}}{2m}, \quad \lambda_2 = \frac{-n + \sqrt{n^2 - 4mk}}{2m}$$

$$(\lambda_1 < \lambda_2)$$

$$A = \frac{d\lambda_2 - \beta}{\lambda_2 - \lambda_1}, \quad B = \frac{\beta - d\lambda_1}{\lambda_2 - \lambda_1}$$

We require that $h\lambda_1$ and $h\lambda_2$ lie in the region of absolute stability



$$\Rightarrow h < \frac{2}{-\lambda_1}$$

we're only worried about $\lambda_1 h$

solution decay-real eigenvalues

$$m'' + \nu u' + k u = 0, \quad t \geq 0, \quad u(0) = \alpha, \quad u'(0) = \beta$$

$$\lambda_1 = \frac{-\nu - \sqrt{\nu^2 - 4mk}}{2m}, \quad \lambda_2 = \frac{-\nu + \sqrt{\nu^2 - 4mk}}{2m}$$

$$m = 0.5, \quad k = 0.5, \quad \nu = 0.2$$

$$\begin{aligned} \lambda_1 &= -0.2 - i0.98... \\ \lambda_2 &= -0.2 + i0.98... \end{aligned}$$

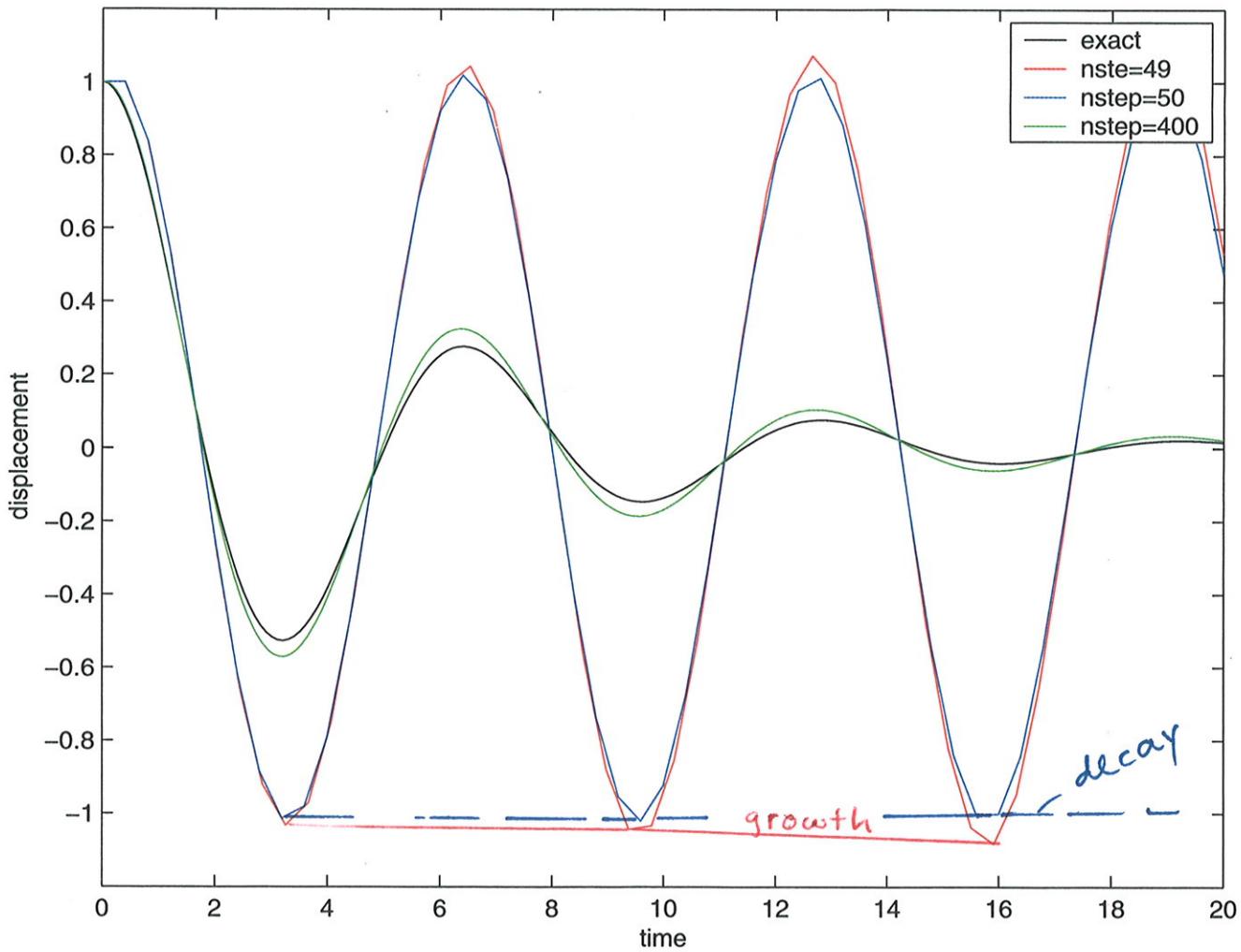
UNDERDAMPED

$$\lambda = \frac{-\nu}{2m} = -0.2$$

Using Euler's method, absolute stability requires that $h < \frac{2}{-\lambda_1}$

$$h \leq \frac{N}{K} = 0.4$$

Region of Absolute Stability



when nstep = 49, $h = \frac{20}{49} = 0.408$, does not obey region, so grows

nstep = 50, $h = \frac{20}{50} = 0.4$, on boundary of disk of stability



(24)

Suppose $m = k = \frac{1}{2}$, $N = 2$

$$\lambda_1 = -2 - \sqrt{3}$$

$$\Rightarrow h < \frac{2}{2 + \sqrt{3}} \approx 0.536$$

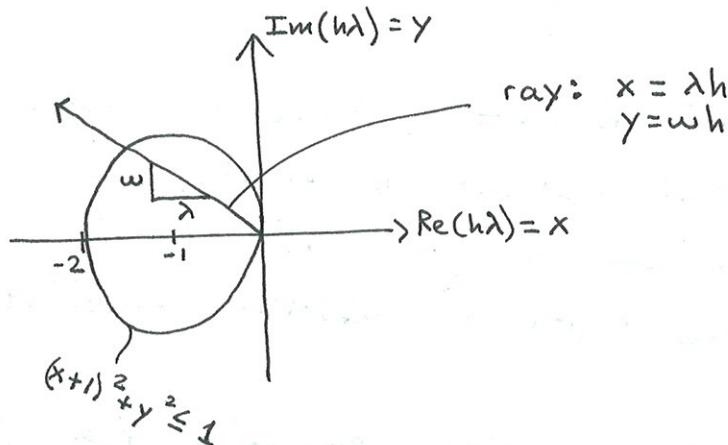
2) Underdamped: $N^2 - 4mk < 0$

Here λ_1, λ_2 are complex conjugates

The solution is: $v = e^{\lambda t} (A \cos(\omega t) + B \sin(\omega t))$

$$\text{where } \lambda = \frac{-N}{2m}, \quad \omega = \frac{1}{2m} \sqrt{4mk - N^2}, \quad A = \alpha, \quad B = \frac{B - \lambda \alpha}{\omega}$$

Absolute stability



$$(x+1)^2 + y^2 = (\lambda h + 1)^2 + (\omega h)^2 \leq 1 \Rightarrow \lambda^2 h^2 + 2\lambda h + 1 + \omega^2 h^2 \leq 1$$

$$\rightarrow h(h\lambda^2 + 2\lambda + \omega^2 h) \leq 0 \rightarrow \lambda^2 h + 2\lambda + \omega^2 h \leq 0$$

$$\rightarrow h \leq \frac{-2\lambda}{\lambda^2 + \omega^2} = \frac{N}{K}$$

$$\text{If } m = k = \frac{1}{2}, \quad N = \frac{1}{5} \quad \text{then } h \leq \frac{2}{5} = 0.4$$

solution is
oscillatory
(imaginary eigenvalues)

Stiffness & Implicit Methods

Region of absolute stability - apply a numerical method to the test equation, $y' = \lambda y$. Plot the region in the complex λh -plane for which the numerical solution is bounded. For Euler's method, the region is a disk of unit radius:



IF $\operatorname{Re} \lambda \leq 0$, then the exact solution of the test equation is stable and the numerical solution is bounded only if h is chosen st $h\lambda$ is in the region of absolute stability.

In the previous example, the choice of h was determined by accuracy and not stability \Rightarrow "non-stiff".

In other cases, the situation could be reversed and a small h is dictated by stability \Rightarrow "stiff."

Example: spring with a small mass

$$\varepsilon u'' + 2u' + u = 0, \quad 0 \leq t \leq 0(1)$$

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

where ε is a positive constant

25

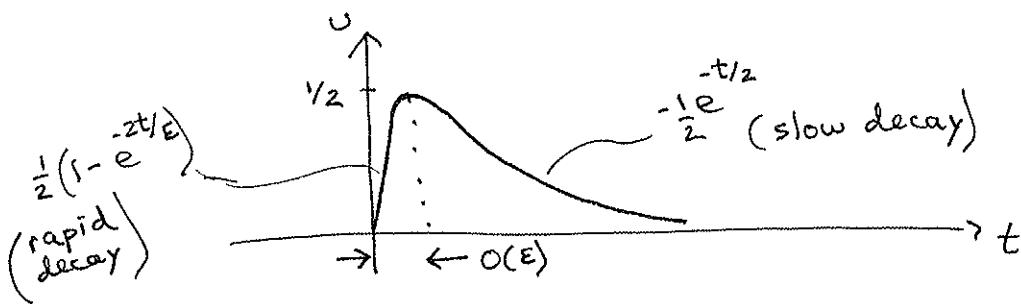
exact solution: $u(t) = A(e^{r_1 t} - e^{r_2 t})$

$$\text{where } r_1 = \frac{-1}{\varepsilon} + \frac{\sqrt{1-\varepsilon}}{\varepsilon}, \quad r_2 = \frac{-1}{\varepsilon} - \frac{\sqrt{1-\varepsilon}}{\varepsilon}, \quad A = \frac{1}{2\sqrt{1-\varepsilon}}$$

in the limit as $\varepsilon \rightarrow 0$, $r_1 \rightarrow -\frac{1}{2}$, $r_2 \rightarrow -\frac{2}{\varepsilon}$, $A \rightarrow \frac{1}{2}$

$$\text{and } u \rightarrow \frac{1}{2}(e^{-t/2} - e^{-2t/\varepsilon})$$

this is a problem of multiple time scales



Accuracy $\Rightarrow h \sim O(\varepsilon)$ for $t \sim O(\varepsilon)$ to resolve rapid decay
and $h \sim O(1)$ for $t \sim O(1)$ to resolve slow decay

Suppose we want to apply Euler's method, then take the equation and write it as a system:

$$x_1 = u, \quad x_2 = u' \Rightarrow \begin{cases} x_1' = x_2 \\ x_2' = -\frac{1}{\varepsilon}(x_1 + 2x_2) \end{cases}$$

This is of the form $y' = Ay$, where $A = \begin{pmatrix} 0 & 1 \\ -1/\varepsilon & -2/\varepsilon \end{pmatrix}$

We'll have: $y_n = y_{n-1} + hAy_{n-1}$ choose h
 $y_0 = \text{given}$

Concerning absolute stability, $h \leq \frac{2}{-\lambda_j}$, $j=1,2$

What are the eigenvalues? $\lambda_1 \sim -1/2$, $\lambda_2 = -2/\varepsilon$

λ_2 is the 'controlling' eigenvalue as its magnitude is greater than that of λ_1

$$\Rightarrow \boxed{h \leq \frac{2}{-(-2/\varepsilon)} = \varepsilon} \quad \text{this is true for all time}$$

The problem of stiffness is essentially an issue of multiple scales, ie, the solution evolves on different time scales.

$$\text{eg, } \dot{\mathbf{y}}' = \mathbf{A} \mathbf{y}, \quad 0 \leq t \leq b$$

\mathbf{A} has eigenvalues $\lambda_j, j=1, \dots, n$

We require all eigenvalues of \mathbf{A} with $\operatorname{Re} \lambda_j \leq 0$ to lie in region of abs stability

The problem is stiff if $b(-\operatorname{Re} \lambda) \gg 1$.

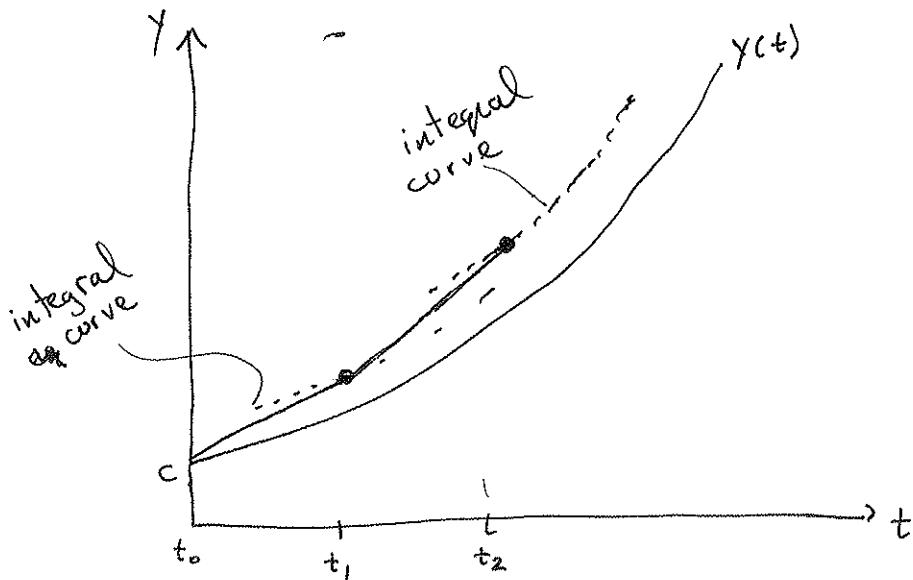
Most methods for stiff equations are implicit.

Standard case: backward Euler method

$$\dot{\mathbf{y}}' = \mathbf{f}(t, \mathbf{y}), \quad t \geq 0, \quad \mathbf{y}(0) = \mathbf{c}$$

$$\text{backward Euler: } \mathbf{y}_n = \mathbf{y}_{n-1} + h \mathbf{f}(t_n, \mathbf{y}_n)$$

The derivation involves a Taylor series centered at (t_n, \mathbf{y}_n) .



Solve

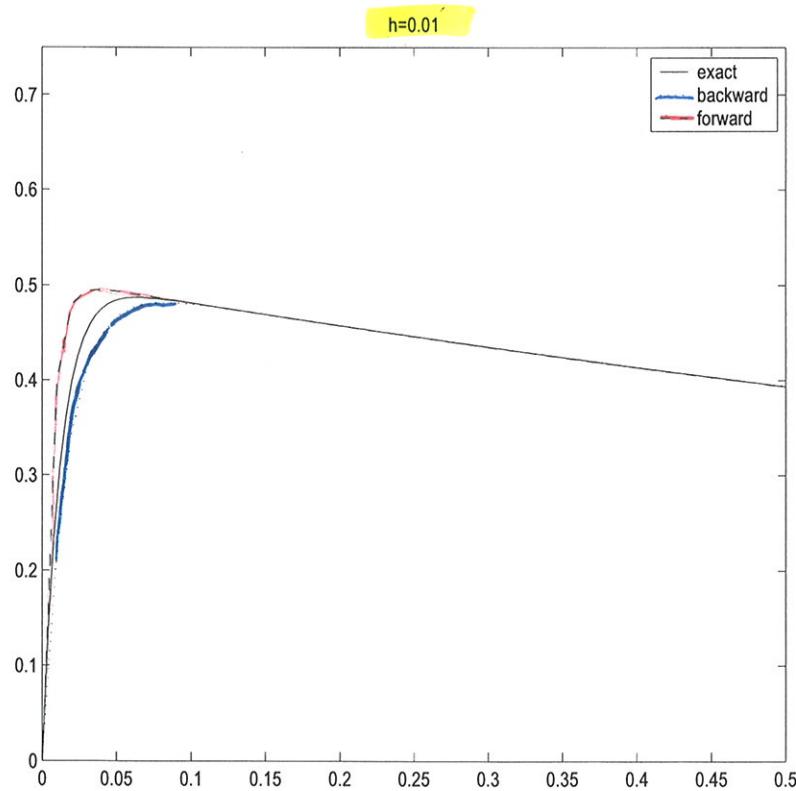
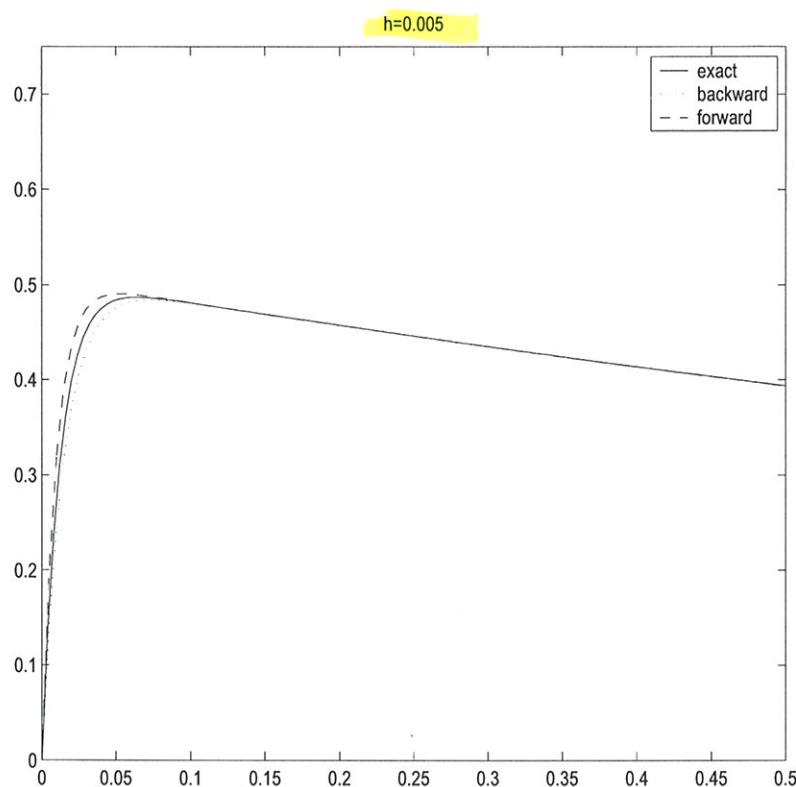
\mathbf{y}_n is defined implicitly \Rightarrow need to solve a set of (nonlinear) algebraic eqns to find \mathbf{y}_n at each step.

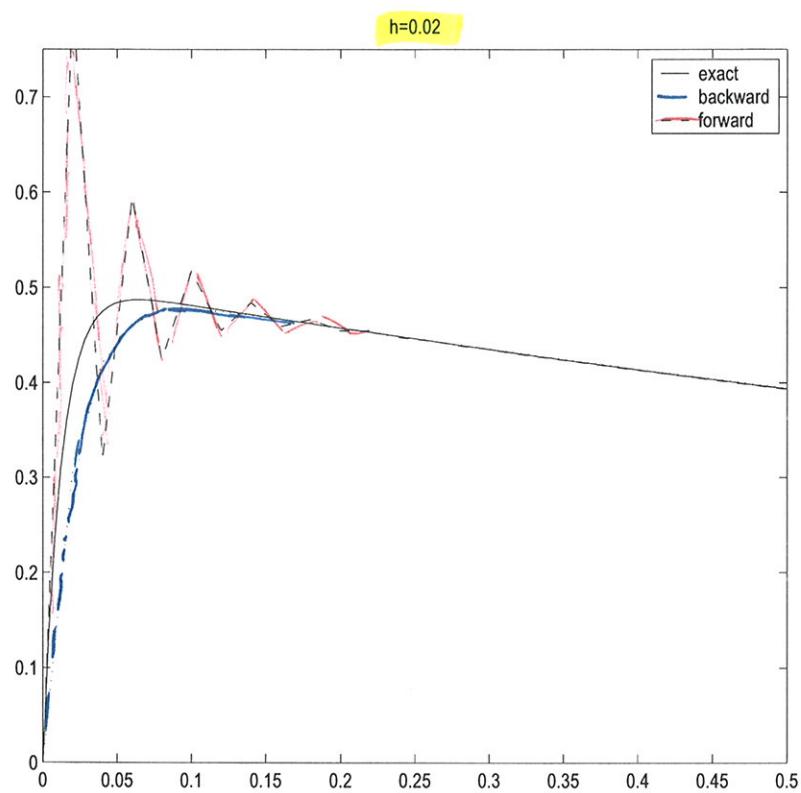
$$\varepsilon v'' + 2v' + v = 0$$

$$v(0) = 0$$

$$\varepsilon v'(0) = 1$$

25a





Apply backward Euler to the test equation:

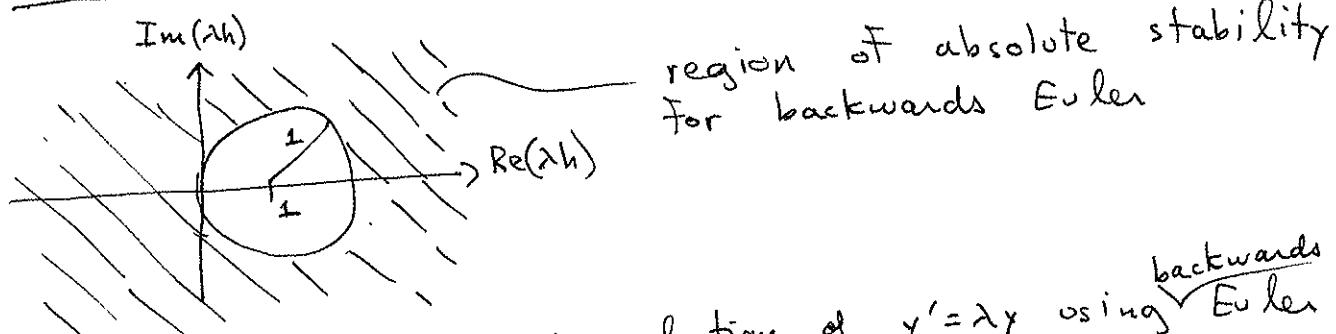
$$\underline{y}' = \lambda y \Rightarrow y_n = y_{n-1} + h f(t_n, y_n)$$

$$y_n = y_{n-1} + h \lambda y_n$$

$$(I - h\lambda) y_n = y_{n-1}$$

$$y_n = \frac{y_{n-1}}{1 - h\lambda}$$

absolute stability: $|\frac{1}{1-h\lambda}| \leq 1 \rightarrow 1 \leq |1-h\lambda|$



Note that a numerical solution of $\underline{y}' = \lambda y$ using backward Euler decays whenever $\text{Re}\lambda < 0$ independent of h . The region of absolute stability covers the whole left half of the λh -plane. Such methods are called A-stable.

Example: small-mass, spring problem

$$\epsilon u'' + 2u' + u = 0, \quad 0 \leq t \leq 2, \quad u(0) = 0, \quad \epsilon u'(0) = 1$$

solve using backward Euler

$$\underline{y}' = \underline{A} y, \quad A = \begin{pmatrix} 0 & 1 \\ -1/\epsilon & -2/\epsilon \end{pmatrix}, \quad y(0) = \begin{pmatrix} 0 \\ y_0 \end{pmatrix}, \quad y_1 = 0, \quad y_2 = u'$$

$$y_n = y_{n-1} + h A y_n \rightarrow (I - hA) y_n = y_{n-1}$$

$$\rightarrow y_n = B y_{n-1}, \quad \text{where } B = (I - hA)^{-1} = \frac{1}{1 + \frac{h}{\epsilon}(2+h)} \begin{bmatrix} 1 + 2h/\epsilon & h \\ -h/\epsilon & 1 \end{bmatrix}$$

nonlinear case

$$y' = f(t, y), \quad 0 \leq t \leq b$$

Now we measure stiffness locally about a given trajectory $\hat{y}(t)$. If we linearize about $\hat{y}(t)$, set $y(t) = \hat{y}(t) + \tilde{y}(t)$, where $\tilde{y}(t)$ is small.

$$\Rightarrow \hat{y}' + \tilde{y}' = f(t, \hat{y} + \tilde{y}) \\ = f(t, \hat{y}) + \frac{\partial f}{\partial y}(t, \hat{y}) \tilde{y} + \dots$$

Suppose $\hat{y}' = f(t, \hat{y})$, then \tilde{y} solves $\tilde{y}' = \frac{\partial f}{\partial y}(t, \hat{y}) \tilde{y}$ to a first approximation. Set $\underline{A}_0 = \frac{\partial f}{\partial y}(t_0, \hat{y}(t_0))$

stiffness depends on the eigenvalues of \underline{A}_0 .

stiffness $\Rightarrow \text{b}(-\text{Re}\lambda_j) \gg 1$.

If the problem were stiff then you might use backward

$$\text{Euler} \Rightarrow y_n = y_{n-1} + h f(t_n, y_n)$$

Need to solve for y_n given y_{n-1} , h , using some functional iteration: choose a function $g(y)$ and iterate: $y^v = g(y^{v-1})$, $v=1, 2, \dots$ such that $y^v \rightarrow y_n$ as $v \rightarrow \infty$. Problem becomes the choice of $g(y)$ and the starting value y^0 .

One choice (poor) is $g(y) = y_{n-1} + h f(t_n, y)$

clearly y_n is a fixed-point of the iteration (ie, $y_n = g(y_n)$).

The question is does this functional iteration converge?

Compute the Jacobian: $g_y = h f_y(t_n, y)$

Require that $\rho(g_y) < 1$, where ρ - spectral radius, the max of the magnitudes of the eigenvalues.

$$\rho(g_y) = \max_j |h\lambda_j| < 1$$

we can always find h small enough so that the iteration converges. However if the problem is stiff then h must be chosen to be "too" small for this to happen. This is not a good iteration scheme.

Newton's method:

$$\underline{g}(\underline{y}) = \underline{y} - \underline{\underline{A}}(\underline{y})(\underline{y}_{n-1} - \underline{y} + h \underline{f}(t_n, \underline{y}))$$

$$\text{where } \underline{\underline{A}}(\underline{y}) = (-I + h \underline{\underline{f}}_y(t_n, \underline{y}))^{-1}$$

$$\begin{aligned} \text{Jacobian: } \underline{g}_y(\underline{y}) &= I - \underbrace{\underline{\underline{A}}(\underline{y})(-I + h \underline{\underline{f}}_y(t_n, \underline{y}))}_{I} - \underline{\underline{A}}_y(\underline{y})(\underline{y}_{n-1} - \underline{y} + h \underline{f}(t_n, \underline{y})) \\ &\rightarrow \underline{g}_y(\underline{y}) = -\underline{\underline{A}}_y(\underline{y})(\underline{y}_{n-1} - \underline{y} + h \underline{f}(t_n, \underline{y})) \end{aligned}$$

$\Rightarrow \underline{g}_y(\underline{y}_n) = 0 \Rightarrow$ quadratic convergence (at least), independent of choice of h .

The asymptotic rate of convergence is quadratic and independent of h but the choice of y^* for convergence does depend on h . Evidently, $y^* = y_{n-1}$ is generally good enough.

Algorithm

- ① select $y^* = y_{n-1}$
- ② for $v = 1, 2, \dots$ do ③ → ⑦
- ③ $\underline{x}^v = \underline{y}^v - \underline{y}_{n-1} - h \underline{f}(t_n, \underline{y}^v)$
- ④ $\underline{\underline{\Sigma}}^v = I - h \underline{\underline{f}}_y(t_n, \underline{y}^v)$ (calculate Jacobian)
- ⑤ solve $\underline{\underline{\Sigma}}^v \underline{\Delta y} = \underline{x}^v$ (solve linear system)
- ⑥ $\underline{y}^{v+1} = \underline{y}^v - \underline{\Delta y}$ (update step)
- ⑦ stop if $\|\underline{\Delta y}\| \leq \text{tol}$

Implicit Methods and Newton's Method

Remember, solving the initial value problem $\dot{y} = f(t, y)$, $y(0) = c$, $t \geq 0$

Numerical solution using backward Euler

$$\underline{y}_n = \underline{y}_{n-1} + h_n \underline{f}(t_n, \underline{y}_n)$$

This turns into a root-finding problem for \underline{y}_n .

$$\text{Define: } \underline{F}(\underline{y}) = \underline{y} - \underline{y}_{n-1} - h_n \underline{f}(t_n, \underline{y})$$

$$\text{Clearly } \underline{F}(\underline{y}_n) = 0$$

Apply Newton's method to $\underline{F}(\underline{y})$:

- 1) $\underline{y}^0 = \underline{y}_{n-1}$, initial guess

- 2) for $v = 0, 1, \dots$,

- 3) compute $\underline{F}' = \underline{F}'(\underline{y}^v)$

- 4) compute $\underline{\Sigma}' = \underline{\Sigma}'(\underline{y}^v)$ (Jacobian)

$$\text{Notice, } \underline{\Sigma}(\underline{y}) = \underline{\Sigma} - h_n \underline{f}_y(t_n, \underline{y})$$

- 5) solve $\underline{\Sigma}' \Delta \underline{y} = \underline{F}'$ ($m \times m$ linear system)

- 6) update, $\underline{y}^{v+1} = \underline{y}^v - \Delta \underline{y}$

- 7) check convergence, $\|\Delta \underline{y}\| \leq \text{tol}$

Remarks:

- 1) The expensive step is ⑤, solving the linear system. Typically use Gaussian elimination, partial pivoting. \Rightarrow cost $\frac{2}{3}m^3$ (LU factorization of $\underline{\Sigma}'$). May want to freeze the Jacobian and its LU-factorization after one or two steps.
- 2) Don't need to set tolerance, ~~$\text{tol} = \|\epsilon_{\text{machine}}\|$~~ $\text{tol} = \|\epsilon_{\text{machine}}\|$, instead only need $\text{tol} \ll \text{local error} = \|\Delta \underline{y}\|$

3) Can avoid the analytic calculation of $\frac{dy}{dt}$ by using finite differences.

$$\text{set } \underline{\alpha}_j = \frac{1}{n} (F(t_n, y^* + ne_j) - F(t_n, y^*))$$

j^{th} column of $n \times n$ identity

$$n \approx \sqrt{\epsilon_{\text{machine}}} \approx 10^{-8}$$

Finite diff approximation of j^{th} column of $F(y)$

Trapezoidal Method

- one parameter family of methods

$$y_n = y_{n-1} + h_n [\theta F(t_n, y_n) + (1-\theta) F(t_{n-1}, y_{n-1})]$$

where θ is a parameter, $\theta \in [0, 1]$

If $\theta = 0 \Rightarrow$ Euler

$\theta = 1 \Rightarrow$ backward Euler

$\theta > 0 \Rightarrow$ implicit method

Consider the order of accuracy:

$$N_n y_n = \frac{y_n - y_{n-1}}{h_n} - [\theta F(t_n, y_n) + (1-\theta) F(t_{n-1}, y_{n-1})]$$

local truncation error is: $d_n = N_n y(t_n) = O(\frac{h_n}{2})$ if $\theta \neq \frac{1}{2}$
 $= O(h_n^2)$ if $\theta = \frac{1}{2}$

If $\theta = \frac{1}{2}$

$$y_n = y_{n-1} + \frac{h_n}{2} [F(t_n, y_n) + F(t_{n-1}, y_{n-1})]$$

→ Trapezoidal Method

You can consider the integral of the ODE:

$$\int_{t_{n-1}}^{t_n} y' dt = \int_{t_{n-1}}^{t_n} f(t, y(t)) dt$$

$$y(t_n) - y(t_{n-1}) = \int_{t_{n-1}}^{t_n} f(t, y(t)) dt \rightarrow \text{approximate using trapezoidal rule}$$

$$\approx \frac{h}{2} \left[f(t_n, y(t_n)) + f(t_{n-1}, y(t_{n-1})) \right]$$

This leads to the Trapezoidal method

Consider the stability properties of the Trapezoidal method

$$y' = \lambda y$$

$$y_n = y_{n-1} + \frac{\lambda h}{2} (y_{n-1} + y_n) \rightarrow \left(1 - \frac{\lambda h}{2}\right) y_n = \left(1 + \frac{\lambda h}{2}\right) y_{n-1}$$

$$\rightarrow y_n = \left(\frac{1 + \frac{\lambda h}{2}}{1 - \frac{\lambda h}{2}} \right) y_{n-1}$$

"growth" factor

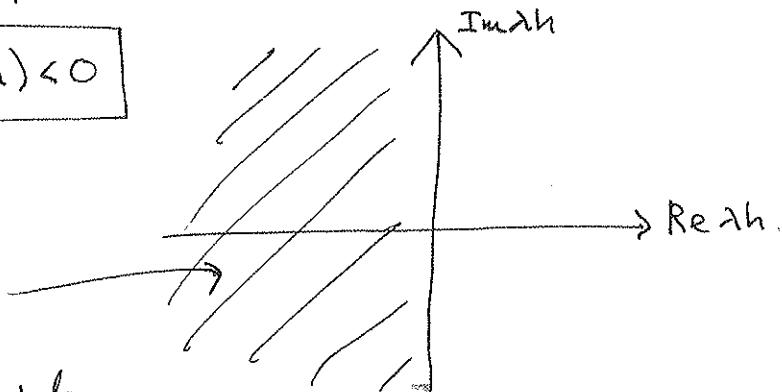
$$\text{Find } \lambda h \text{ complex st } \left| \frac{1 + \frac{\lambda h}{2}}{1 - \frac{\lambda h}{2}} \right| \leq 1$$

$$\text{set } \lambda h = z = a + ib \rightarrow \left| 1 + \frac{a+ib}{2} \right|^2 \leq \left| 1 - \frac{a+ib}{2} \right|^2$$

$$\rightarrow \left(1 + \frac{a}{2}\right)^2 + \frac{b^2}{4} \leq \left(1 - \frac{a}{2}\right)^2 + \frac{b^2}{4} \rightarrow 1 + a + \frac{a^2}{4} \leq 1 - a + \frac{a^2}{4}$$

$$\rightarrow \boxed{a \leq 0} \Rightarrow \boxed{\operatorname{Re}(\lambda h) < 0}$$

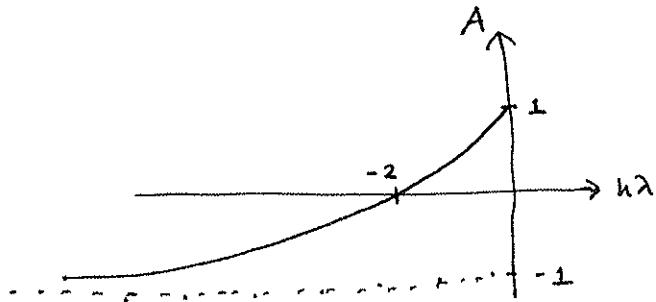
absolute stability
for trap method
hence, A-stable



so the trapezoidal method is A-stable

\Rightarrow if the exact solution of the test equation decays, then the numerical method solution decays independent of h , and if the exact solution grows the numerical soln grows (independent of h)

Consider real λ with $\lambda < 0$, let $A(h\lambda) = \frac{1 + \frac{\lambda h}{2}}{1 - \frac{\lambda h}{2}}$



In the very stiff limit, $-\operatorname{Re}(\lambda h) \gg 1$, then

$$A(h\lambda) = \frac{1 + \frac{\lambda h}{2}}{-1 + \frac{\lambda h}{2}} \sim -1 + \frac{4}{(-\lambda h)} + \dots = -1 + \delta$$

so, $y_n \sim -y_{n-1}$, decay, but very slowly. The solution has a ± 1 oscillation.

Remember, backward Euler, $A = \frac{1}{1-\lambda h} \rightarrow 0$ as $-\operatorname{Re}(\lambda h) \rightarrow \infty$, which is referred to as "stiff" decay.

Example: $\varepsilon u'' + 2u' + u = 0$, $t \geq 0$
 $u(0) = 0$, $\varepsilon u'(0) = 1$, $0 < \varepsilon \ll 1$

$$\Rightarrow \underline{y}' = \underline{A} \underline{y}, \quad \underline{A} = \begin{bmatrix} 0 & 1 \\ -1/\varepsilon & -2/\varepsilon \end{bmatrix}, \quad \underline{y} = \begin{pmatrix} u(t) \\ u'(t) \end{pmatrix}$$

Apply Trapezoidal method:

$$y_n = y_{n-1} + \frac{h}{2}(A y_{n-1} + A y_n) \rightarrow (I - \frac{h}{2}A) y_n = (I + \frac{h}{2}A) y_{n-1}$$

$$y_n = \underbrace{(I - \frac{h}{2}A)^{-1}}_B (I + \frac{h}{2}A) y_{n-1} \Rightarrow y_n = B y_{n-1}$$

$$\rightarrow \text{numerical solution } y_n = B^n y_0, \quad y_0 = \begin{pmatrix} 0 \\ y_0 \end{pmatrix}$$

$$\text{and } B = \begin{bmatrix} 1 + \frac{h}{\varepsilon} - \frac{h^2}{4\varepsilon} & h \\ -\frac{h}{\varepsilon} & 1 - \frac{h}{\varepsilon} - \frac{h^2}{4\varepsilon} \end{bmatrix} \frac{1}{(1 + \frac{h}{\varepsilon} + \frac{h^2}{4\varepsilon})}$$

Single-Step Methods

General form of an explicit single-step method is

$$\underline{y}_n = \underline{y}_{n-1} + h_n \phi(t_{n-1}, \underline{y}_{n-1}, h_n)$$

The function $\phi(t, y, h)$ is called the stepping function.

For example, $\phi(t, y, h) = f(t, y) \Rightarrow$ Euler

Implicit single-step methods have the form

$$\underline{y}_n = \underline{y}_{n-1} + h_n \hat{\phi}(t_{n-1}, \underline{y}_{n-1}, \underline{y}_n, h_n)$$

For example $\hat{\phi}(t, y, wh) = f(t+h, w) \Rightarrow$ backwards Euler

$\hat{\phi}(t, y, w, h) = \frac{1}{2}(f(t, y) + f(t+h, w)) \Rightarrow$ trapezoidal method

Main problem is to construct $\phi(t, y, h)$ to obtain higher-order accuracy. Simplest approach is via Taylor expansions.

$$y(t_n) = y(t_{n-1} + h)$$

$$= y(t_{n-1}) + hy'(t_{n-1}) + \frac{h^2}{2} y''(t_{n-1}) + \dots + \frac{h^p}{p!} y^{(p)}(t_{n-1}) + \dots$$

we know that $y(t)$ solves the ODE $y'(t) = f(t, y(t))$

$$\Rightarrow y'(t_{n-1}) = f(t_{n-1}, y(t_{n-1}))$$

$$y''(t_{n-1}) = \left. \frac{d}{dt} f(t, y(t)) \right|_{t=t_{n-1}} = \left[\frac{\partial f}{\partial t} + \frac{\partial f}{\partial y} f \right]_{t=t_{n-1}}$$

$$y'''(t_{n-1}) = \left[f_{tt} + f_{ty} f + f_{yt} f + f_{yy} f^2 + f_y f_t + f_y^2 f \right]_{t=t_{n-1}}$$

p^{th} order Taylor method

$$y_n = y_{n-1} + h f \Big|_{(t_{n-1}, y_{n-1})} + \frac{h^2}{2} \left. \frac{df}{dt} \right|_{(t_{n-1}, y_{n-1})} + \dots + \frac{h^p}{p!} \left. \frac{d^{p-1} f}{dt^{p-1}} \right|_{(t_{n-1}, y_{n-1})}$$

$$\text{Here, } \phi = f \Big|_{(t_{n-1}, y_{n-1})} + \frac{h}{2} \left. \frac{df}{dt} \right|_{(t_{n-1}, y_{n-1})} + \dots + \frac{h^{p-1}}{p!} \left. \frac{d^{p-1} f}{dt^{p-1}} \right|_{(t_{n-1}, y_{n-1})}$$

Eq., 2nd order Taylor method

$$y_n = y_{n-1} + h \left[f(t_{n-1}, y_{n-1}) + \frac{h}{2} (f_t(t_{n-1}, y_{n-1}) + f_y(t_{n-1}, y_{n-1}) f(t_{n-1}, y_{n-1})) \right]$$

- this is kind of like Euler's method but with a correction

- not hard to show this is 2nd order accurate

~~In principle you can get arbitrarily high order of accuracy~~

$$d_n = \frac{y_n - y_{n-1}}{h} - \left\{ f(t_{n-1}, y_{n-1}) + \frac{h}{2} (f_t(t_{n-1}, y_{n-1}) + f_y(t_{n-1}, y_{n-1}) f(t_{n-1}, y_{n-1})) \right\}$$

$$d_n = y'(t_{n-1}) + \frac{h}{2} y''(t_{n-1}) + \frac{h^2}{6} y'''(t_{n-1}) - \left\{ y'(t_{n-1}) + \frac{h}{2} y''(t_{n-1}) \right\}$$

$$= \frac{h^2}{6} y'''(t_{n-1}) + O(h^3) \Rightarrow d_n = O(h^2) \text{ as } h \rightarrow 0$$

Remarks on Taylor method

- In principle you can get arbitrarily high order of accuracy
- There is an explosion of terms involved in the analytic calculation of the derivatives. This can be avoided by imbedding function evaluations \Rightarrow Runge-Kutta methods.

2/14/03

Runge-Kutta Methods

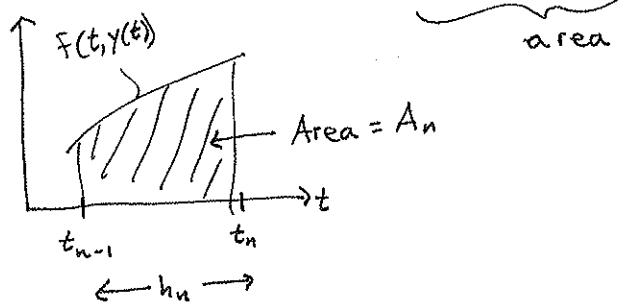
ODE: $y' = f(t, y)$

Main idea: Suppose we know y at t_{n-1} (approximately perhaps) we want to find an approximation for y at $t_n = t_{n-1} + h$

Integrate the ODE:

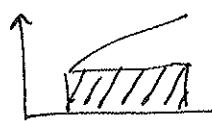
$$y(t_n) = y(t_{n-1}) + \int_{t_{n-1}}^{t_n} f(t, y(t)) dt$$

$\underbrace{\hspace{1cm}}$ area



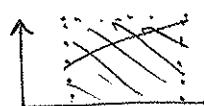
Approximate A_n

$$\text{eg: } A_n \approx h f(t_{n-1}, y(t_{n-1})) \Rightarrow \text{Euler}$$



(EXPLICIT)

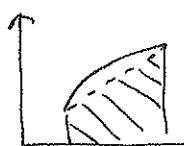
$$A_n \approx h f(t_n, y(t_n)) \Rightarrow \text{Backward Euler}$$



(IMPLICIT)

$$A_n \approx \frac{h}{2} (f(t_{n-1}, y(t_{n-1})) + f(t_n, y(t_n)))$$

\Rightarrow Trapezoidal method



(IMPLICIT)

EXPLICIT MIDPOINT METHOD

Consider the implicit midpoint method:

$$y_n = y_{n-1} + h f\left(t_{n-1/2}, \frac{y_{n-1} + y_n}{2}\right)$$

approximate $y(t_{n-1/2})$ by forward Euler then substitute:

$$\begin{aligned} \hat{y}_{n-1/2} &= y_{n-1} + \frac{h}{2} f(t_{n-1}, y_{n-1}) \\ y_n &= y_{n-1} + h f(t_{n-1/2}, \hat{y}_{n-1/2}) \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{2-stage RK method}$$

local truncation error

$$d_n = \frac{y(t_n) - y(t_{n-1})}{h} - f\left(t_{n-1/2}, y(t_{n-1}) + \frac{h}{2} f(t_{n-1}, y(t_{n-1}))\right)$$

Expand in a Taylor series: (centered about $t_{n-1}, y(t_{n-1})$)
1st term:

$$\begin{aligned} \frac{y(t_n) - y(t_{n-1})}{h} &= \frac{y(t_{n-1}) + hy'(t_{n-1}) + \frac{h^2}{2}y''(t_{n-1}) + \frac{h^3}{6}y'''(t_{n-1}) + \dots - y(t_{n-1})}{h} \\ &= y'(t_{n-1}) + \frac{h}{2}y''(t_{n-1}) + \frac{h^2}{6}y'''(t_{n-1}) + \dots \end{aligned}$$

2nd term

$$\begin{aligned} f\left(t_{n-1/2}, y(t_{n-1}) + \frac{h}{2} f(t_{n-1}, y(t_{n-1}))\right) &= f\left(t_{n-1} + \frac{h}{2}, y(t_{n-1}) + \Delta y\right) \\ &= f + \frac{h}{2}f_t + \Delta y f_y + \frac{1}{2!} \left(\left(\frac{h}{2}\right)^2 f_{tt} + 2f_{ty} \left(\frac{h}{2}\right)(\Delta y) + f_{yy} (\Delta y)^2 \right) + \dots \\ &= f + \frac{h}{2}f_t + \frac{h}{2}f_t \cdot f_y + \frac{1}{2!} \left(\frac{h^2}{4} f_{tt} + 2f_{ty} \frac{h}{2} \cdot \frac{h}{2} f + f_{yy} \left(\frac{h}{2} y\right)^2 \right) + \dots \\ &= f + \frac{h}{2} \left(f_t + f f_y \right) + \frac{h^2}{8} \left(f_{tt} + 2f_{ty} f + f_{yy} f^2 \right) + \dots \end{aligned}$$

$$\Rightarrow d_n = \cancel{y' + \frac{h}{2}y'' + \frac{h^2}{6}y'''} - \left(f + \frac{h}{2} \cancel{(f_t + f f_y)} + \frac{h^2}{8} (f_{tt} + 2f_{ty} f + f_{yy} f^2) \right) + O(h^3)$$

$$\text{Notice: } y' = f(t, y) \text{ so } y'' = \frac{\partial f}{\partial t} \frac{\partial t}{\partial t} + \frac{\partial f}{\partial y} \frac{\partial y}{\partial t} = f_t + f_y f$$

$$\Rightarrow \frac{h}{2}y'' - \frac{h}{2}(f_t + f f_y) = 0$$

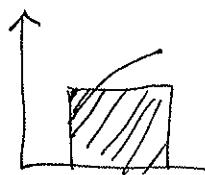
$$\Rightarrow d_n = \frac{h^2}{6}y''' - \frac{h^2}{8}(f_{tt} + 2f_{ty} f + f_{yy} f^2) + O(h^3)$$

$$\Rightarrow \lim_{h \rightarrow 0} \frac{|d_n|}{h^2} = c \Rightarrow$$

the method is consistent
of order 2

$$\Delta_n \approx h f(t_{n-1}, \frac{y(t_{n-1}) + y(t_n)}{2})$$

\Rightarrow Midpoint rule



(IMPLICIT)

Convert the implicit methods to explicit methods using a suitable Euler step:

$$\text{eg: } \hat{y} = y_{n-1} + \frac{h}{2} f(t_{n-1}, y_{n-1}) \approx y(t_{n-1/2}) \Rightarrow \text{Euler half step}$$

use this in Midpoint rule

$$y_n = y_{n-1} + h f(t_{n-1/2}, \hat{y}) \quad \left. \begin{array}{l} \text{"two stage" Runge Kutta (RK) Method} \\ \text{- explicit method} \\ \text{midpoint} \end{array} \right\}$$

Consider:

$$\hat{y} = y_{n-1} + h f(t_{n-1}, y_{n-1}) \approx y(t_n)$$

substitute into trapezoidal method

$$y_n = y_{n-1} + \frac{h}{2} (f(t_{n-1}, y_{n-1}) + f(t_n, \hat{y}))$$

Classical 4th order Runge Kutta (RK4)

$$y_1 = y_{n-1}$$

$$y_2 = y_{n-1} + \frac{h}{2} f(t_{n-1}, y_1)$$

$$y_3 = y_{n-1} + \frac{h}{2} f(t_{n-1/2}, y_2)$$

$$y_4 = y_{n-1} + h f(t_{n-1/2}, y_3)$$

$$y_n = y_{n-1} + \frac{h}{6} [f(t_{n-1}, y_1) + 2f(t_{n-1/2}, y_2) + 2f(t_{n-1/2}, y_3) + 4f(t_n, y_4)]$$

Example: $y' = e^{-t} - y$, $y(0) = 0$

Calculate the solution for $0 \leq t \leq 4$ using various values for h and various RK methods (all explicit) and compare to the exact solution, $y(t) = t e^{-t}$.

error: p^{th} order method

$$e(h) \approx Ch^p \rightarrow \text{global error}$$

$$e\left(\frac{h}{2}\right) \approx C\left(\frac{h}{2}\right)^p$$

$$\rightarrow \frac{e\left(\frac{h}{2}\right)}{e(h)} = 2^{-p} \rightarrow \log\left(\frac{e\left(\frac{h}{2}\right)}{e(h)}\right) = -p \log 2$$

$$\rightarrow -p \approx \frac{\log\left(\frac{e(h_2)}{e(h)}\right)}{\log 2} \rightarrow p \approx \frac{\log\left(\frac{e(h)}{e(h_2)}\right)}{\log 2}$$

General forms

General s -stage RK methods

$$\Sigma_i = y_{n-1} + h \sum_{j=1}^s a_{ij} f(t_{n-1} + c_j h, \Sigma_j), \quad i=1, \dots, s$$

Then

$$y_n = y_{n-1} + h \underbrace{\sum_{i=1}^s b_i f(t_{n-1} + c_i h, \Sigma_i)}_{\text{quadrature rule}}$$

Note that Σ_i as an approximation for y at $t_{n-1} + c_i h$, $0 \leq c_i \leq 1$. The constants $\{a_{ij}\}$, $\{c_i\}$, $\{b_i\}$ are to choose.

Consistency will place some constraints on the constants

Let us consider consistency:

Expand for h small

$$y_{n-1} + h \sum_{j=1}^s a_{ij} f(t_{n-1} + c_j h, \bar{Y}_j)$$

$$= y_{n-1} + h \sum_{j=1}^s a_{ij} f(t_{n-1}, y_{n-1}) + O(h^2)$$

$$= y_{n-1} + h \left(\sum_{j=1}^s a_{ij} \right) f(t_{n-1}, y_{n-1}) + O(h^2) \approx y(t_{n-1} + c_i h)$$

↪ like Euler's method, so we require $\boxed{\sum_{j=1}^s a_{ij} = c_i, i=1, \dots, s}$

Also $y_{n-1} + h \sum_{i=1}^s b_i f(t_{n-1} + c_i h, \bar{x}_i)$

$$= y_{n-1} + h \left(\sum_{i=1}^s b_i \right) f(t_{n-1}, y_{n-1}) + O(h^2) \approx x(t_n)$$

like Euler's method, so require $\boxed{\sum_{i=1}^s b_i = 1}$

Convenient way to express the coefficients

$$\begin{array}{c|ccccc} c_1 & a_{11} & \cdots & a_{1s} & \leftarrow \text{sums to } c_1 \\ c_2 & a_{21} & \cdots & a_{2s} & \\ \vdots & \vdots & \ddots & \vdots & \\ c_s & a_{s1} & \cdots & a_{ss} & \\ \hline b_1 & \cdots & b_s & \leftarrow \text{sums to 1} & \end{array}$$

} most general form which includes implicit methods

If you restrict to explicit methods, the matrix becomes triangular
 $a_{ij} = 0$ for $j > i$

$$\begin{array}{c|cccccc} 0 & 0 & & & & & & \\ c_2 & a_{21} & 0 & & & & & \\ c_3 & a_{31} & a_{32} & 0 & & & & \\ \vdots & \vdots & & \ddots & & & & \\ c_s & a_{s1} & a_{s2} & \cdots & a_{s,s-1} & 0 & & \\ \hline b_1 & b_2 & \cdots & \cdots & \cdots & b_s & & \end{array}$$

$$\frac{c}{\Delta t} \frac{\Delta}{b^T}$$

Example

- Euler $\begin{array}{|c|c|} \hline 0 & 0 \\ \hline 1 & \\ \hline \end{array}$

backward Euler: $\begin{array}{|c|c|} \hline 1 & 1 \\ \hline 1 & \\ \hline \end{array}$

- Midpoint (explicit) $\begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 1/2 & & \\ \hline 0 & 1 & \\ \hline \end{array}$

Midpoint (implicit)

$$\begin{array}{|c|c|} \hline 1/2 & 1/2 \\ \hline & 1 \\ \hline \end{array}$$

- Trapezoidal Rule $\begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 1 & 1 & 0 \\ \hline & y_2 & 1/2 \\ \hline \end{array}$

$$\frac{s}{p} = \frac{1}{2}$$

- RK4

$$\begin{array}{|c|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline 1/2 & 1/2 & 0 & 0 \\ \hline 1/2 & 0 & 1/2 & 0 \\ \hline 1 & 0 & 0 & 1/2 \\ \hline & 1/6 & 2/6 & 2/6 & 1/6 \\ \hline \end{array}$$

Order of Accuracy

For ~~given~~ an s -stage RK method we have many parameters to choose

Impose: $\sum_{i=1}^s b_i = 1$, $\sum_{j=1}^s a_{ij} = c_i$, $i=1, \dots, s$ for consistency.

Beyond that, choose coeffs to:

(1) explicit vs implicit

(2) order of accuracy (want high order methods)

(3) ~~order of accuracy~~ favorable stability properties

(4) error control

Focus on the choice for higher order methods. (accuracy)

For single-step methods, obtaining a higher order method becomes choosing parameters st

truncation error $\sim O(h^p)$

for p^{th} order method

All RK methods are O -stable if f satisfies a Lipschitz condition.

$$\Rightarrow |e_n| \leq k \max |d_i| = O(h^p)$$

truncation error

$$\text{so, } d_n = \frac{y(t_n) - y(t_{n-1})}{h} - h \sum_{i=1}^s b_i \tilde{f}(t_{n-1} + c_i h, \Sigma_i)$$

$$\Sigma_i = y(t_{n-1}) + h \sum_{j=1}^s a_{ij} \tilde{f}(t_{n-1} + c_j h, \Sigma_j)$$

$$\frac{y(t_n) - y(t_{n-1})}{h} = \frac{y(t_{n-1}) + h y'(t_{n-1}) + \frac{h^2}{2} y''(t_{n-1}) + \dots - y(t_{n-1})}{h}$$

$$\Rightarrow d_n = y'(t_{n-1}) + \frac{h}{2} y''(t_{n-1}) + \dots - h \sum_{i=1}^s b_i \tilde{f}(t_{n-1} + c_i h, \Sigma_i) = \cancel{\text{O}(h^p)} \text{ O}(h^p)$$

For this to $\text{O}(h^p)$ the sum from $i=1 \dots s$ has to cancel the first p terms of the expansion

$$\Rightarrow \underbrace{\sum_{i=1}^s b_i \tilde{f}(t_{n-1} + c_i h, \Sigma_i)}_{\text{Expand in Taylor series - this is an "enormous task"} \atop \text{Taylor series within Taylor series} \Rightarrow \text{"an explosion of terms"} \atop \text{- Taylor series within Taylor series for } p\text{-order accuracy}} = y'(t_{n-1}) + \frac{h}{2} y''(t_{n-1}) + \dots + \frac{h^{p-1}}{p!} y^{(p)}(t_{n-1}) + \cancel{\text{O}(h^p)}$$

- Instead, develop certain necessary conditions for p -order accuracy by applying the previous formula for "smart" choices of $\tilde{f}(t, y)$.

Choose: $\tilde{f}(t, y) = y + t^{l-1}$ where l is a positive integer
(linear in y , polynomial in t)

set $t_{n-1} = 0, y_{n-1} = 0$

we're only interested in one step:

$$\Sigma_i = h \sum_{j=1}^s a_{ij} \tilde{f}(t_{n-1} + c_j h, \Sigma_i) = h \sum_{j=1}^s a_{ij} (\Sigma_j + (c_j h)^{l-1})$$

$$\text{Set } \underline{\Sigma} = \begin{bmatrix} \Sigma_1 \\ \Sigma_2 \\ \vdots \\ \Sigma_s \end{bmatrix}, \underline{A} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1s} \\ a_{21} & a_{22} & \dots & a_{2s} \\ \vdots & \vdots & \ddots & \vdots \\ a_{s1} & a_{s2} & \dots & a_{ss} \end{bmatrix}, \underline{c} = \begin{bmatrix} c_1^{l-1} \\ c_2^{l-1} \\ \vdots \\ c_s^{l-1} \end{bmatrix}$$

$$\underline{\Sigma} = h \underline{A} \left(\underline{\Sigma} + h^{\frac{l-1}{2}} \underline{z} \right) = h \underline{A} \underline{\Sigma} + h^{\frac{l-1}{2}} \underline{z}$$

$$\Rightarrow (I - h \underline{A}) \underline{\Sigma} = h^{\frac{l}{2}} \underline{z} \Rightarrow \underline{\Sigma} = h^{\frac{l}{2}} (I - h \underline{A})^{-1} \underline{z}$$

Then: $y_n = y_{n-1} + h \sum_{i=1}^s b_i f(t_{n-1} + c_i h, \underline{\Sigma}_i)$

let $\underline{b} = \begin{pmatrix} b_1 \\ \vdots \\ b_s \end{pmatrix}$

$\underline{\Sigma}_i + (c_i h)$

$$\rightarrow y_n = h \underline{b}^T \underbrace{(\underline{\Sigma} + h^{\frac{l-1}{2}} \underline{z})}_{h^{\frac{l}{2}} \underline{A}^{-1} \underline{\Sigma}} \quad (\text{from top line of page})$$

$$\rightarrow y_n = \underline{b}^T \underline{A}^{-1} \underline{\Sigma}$$

$y_n = h \underline{b}^T \underline{A}^{-1} (I - h \underline{A})^{-1} \underline{A} \underline{z}$

① this is an approximation
for $y(t_n)$ where $y' = t^{\frac{l-1}{2}} + y$
with $y(t_{n-1}) = 0, t_{n-1} = 0$

The exact solution: $y(t) = \int_0^t e^{t-\tau} \tau^{\frac{l-1}{2}} d\tau$

so $y(t_n) = y(h) = \int_0^h e^{h-\tau} \tau^{\frac{l-1}{2}} d\tau$ ②

NOTE,
 y_n is a scalar

General s-stage Runge Kutta method

(34)

$$y_i = y_{n-1} + h \sum_{j=1}^s a_{ij} F(t_j, y_j) , \quad i=1, \dots, s$$

then $y_n = y_{n-1} + h \sum_{i=1}^s b_i F(t_i, y_i)$

nodes: $t_i = t_{n-1} + c_i h , \quad 0 \leq c_i \leq 1$

$$y_i \approx y(t_i)$$

constraints that enforce consistency: $\sum_{i=1}^s b_i = 1 , \sum_{j=1}^s a_{ij} = c_i , \quad i=1, \dots, s$

Matrix form:
$$\begin{array}{c|cc} c_i & a_{ij} \\ \hline & b_i \end{array}$$
 explicit then $a_{ij}=0$ if $j > i$

constraints that enforce order of accuracy:

- Find constraints on the a 's, b 's, c 's st
the method is p^{th} order accurate.

choose a specific ODE \Rightarrow necessary order conditions

let $y' = y + t^{l-1} , \quad l = 1, 2, \dots$

$$y(t_{n-1}) = y_{n-1}$$

without loss you can take $t_{n-1} = y_{n-1} < 0$

apply an s-stage RK to the IVP

$$\textcircled{1} \Rightarrow \boxed{y_n = h^l b^T \bar{A}^{-1} (\bar{I} - h \bar{A})^{-1} \bar{A} \bar{z}} , \quad \text{where } \bar{z} = \begin{bmatrix} c_1^{l-1} \\ \vdots \\ c_s^{l-1} \end{bmatrix}$$

notice, y_n is a scalar

$$\text{the exact solution is } y(t) = \int_0^t e^{t-x} x^{l-1} dx$$

$$\textcircled{2} \Rightarrow \boxed{y(t_n) = y(h) = \int_0^h e^{h-x} x^{l-1} dx}$$

Expand $\textcircled{1}$ and $\textcircled{2}$ in Taylor series about $h=0$, then
match terms up to and including h^p .

$$\Rightarrow \text{local error} = O(h^{p+1})$$

$$\Rightarrow \text{truncation error} = O(h^p)$$

$\Rightarrow p^{\text{th}}$ order method

Note, $(I - h \underline{A})^{-1}$ can be expanded as geometric series

$$(I - h \underline{A})^{-1} = I + h \underline{A} + h^2 \underline{A}^2 + h^3 \underline{A}^3 + \dots$$

$$\underline{A}^{-1}(I - h \underline{A}^{-1})A = I + h \underline{A} + h^2 \underline{A}^2 + h^3 \underline{A}^3 + \dots$$

$$\Rightarrow y_n = h^l \underline{b}^T (I + h \underline{A} + h^2 \underline{A}^2 + \dots) \underline{z} \quad (1e)$$

$$y(h) = y(0) + hy'(0) + \frac{h^2}{2} y''(0) + \dots$$

$$= (l-1)! \left[\frac{h^l}{l!} + \frac{h^{l+1}}{(l+1)!} + \dots \right] \quad (2e)$$

Match terms in (1e) and (2e)

$$\text{coeffs of } h^l: \frac{(l-1)!}{l!} = \underline{b}^T \underline{z} \Rightarrow \boxed{\underline{b}^T \underline{z} = \frac{1}{l}} \quad l = 1, 2, \dots, p$$

$$h^{l+1}: \frac{(l-1)!}{(l+1)!} = \underline{b}^T \underline{A} \underline{z} \Rightarrow \boxed{\underline{b}^T \underline{A} \underline{z} = \frac{1}{l(l+1)}}, \quad l = 1, 2, \dots, p-1$$

.

:

$$h^{l+k}: \boxed{\underline{b}^T \underline{A}^k \underline{z} = \frac{1}{l(l+1)\dots(l+k)}}, \quad l = 1, \dots, p-k, \quad 0 \leq k \leq p-1$$

These are necessary conditions for a p^{th} order method

SPECIAL CASES

$$1) \quad \underline{b}^T \underline{z} = \frac{1}{l}, \quad l = 1, 2, \dots, p \quad \Rightarrow \quad \boxed{\sum_{i=1}^s b_i c_i^{l-1} = \frac{1}{l}} \quad \Rightarrow \text{"pure" quadrature order condition}$$

This condition obtained by considering the ODE $y' = t^{l-1}$

$$y(h) = y(0) + \underbrace{\int_0^h t^{l-1} dt}_{\text{approximating this quadrature}}$$

$\sum_{i=1}^s b_i c_i^{l-1} = \frac{1}{l} \Rightarrow$ the integral is exact for polynomials of degree p or less

SPECIAL CASES

35

(2) Set $\underline{l} = 1$, replace k by $k-1$ for general case

$$\Rightarrow \underline{b}^T \underline{A}^k \underline{z} = \frac{1}{\underline{l}(\underline{l}+1)\cdots(\underline{l}+k)}$$

$$\Rightarrow \underline{b}^T \underline{A}^{k-1} \begin{bmatrix} \underline{l} \\ 1 \\ \vdots \\ 1 \end{bmatrix} = \frac{1}{k!}$$

This is the order condition for $y' = y$

Notice: if $k=1 \Rightarrow \sum_{i=1}^s b_i c_i = 1$ (already established this constraint)

$$\text{if } k=2 \Rightarrow \underline{b}^T \underline{A}^1 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \underline{b}^T \underline{c} = \frac{1}{2}$$

this overlaps with the other special case (1) when $\underline{l}=2$

If you set $\underline{l}=1$, the leading term in the truncation error is

$$h^p \left(\underline{b}^T \underline{A}^p \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} - \frac{1}{(p+1)!} \right) \quad \text{IF the method is explicit, then}$$

$\underline{A}^i = 0$ if $i \geq s$. What this implies is that the maximal order of an explicit s-stage method is $p=4$.

Example: Consider the classical RK4 method

0	0	0	0
$\frac{1}{2}$	$\frac{1}{2}$	0	0
$\frac{1}{2}$	0	$\frac{1}{2}$	0
1	0	0	1
$\frac{1}{6}$	$\frac{2}{6}$	$\frac{2}{6}$	$\frac{1}{6}$

check order conditions:

$$\textcircled{1} \quad \sum_{i=1}^s b_i c_i^{l-1} = \frac{1}{l}, \quad l=1, \dots, p$$

$$\Rightarrow \sum_{i=1}^4 b_i c_i^{l-1} = \frac{1}{l}, \quad l=1, \dots, 4$$

examine $l=1 \Rightarrow \sum_{i=1}^4 b_i c_i = 1$ (satisfied)

$$l=2 \Rightarrow \sum_{i=1}^4 b_i c_i = \frac{2}{6} \left(\frac{1}{2}\right) + \frac{2}{6} \left(\frac{1}{2}\right) + \frac{1}{6} (1) = \frac{1}{2} \quad (\text{satisfied})$$

$$l=3 \Rightarrow \sum_{i=1}^4 b_i c_i^2 = \frac{2}{6} \left(\frac{1}{2}\right)^2 + \frac{2}{6} \left(\frac{1}{2}\right)^2 + \frac{1}{6} (1)^2 = \frac{1}{3} \quad (\text{satisfied})$$

$$l=4 \Rightarrow \sum_{i=1}^4 b_i c_i^3 = \frac{2}{6} \left(\frac{1}{2}\right)^3 + \frac{2}{6} \left(\frac{1}{2}\right)^3 + \frac{1}{6} (1)^3 = \frac{1}{4} \quad (\text{satisfied})$$

$$\textcircled{2} \quad \text{check: } b^T A^{k-1} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \frac{1}{k!}, \quad k=3,4$$

$$A = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 1/2 & 0 & 0 & 0 \\ 0 & 1/2 & 0 & 0 \\ 0 & 0 & 1/2 & 0 \end{bmatrix}, \quad A^2 = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1/4 & 0 & 0 & 0 \\ 0 & 1/2 & 0 & 0 \end{bmatrix}$$

$$A^2 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1/4 \end{bmatrix}, \quad b^T A^2 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \frac{1}{4} b_3 + \frac{1}{2} b_4 = \frac{1}{4} \left(\frac{1}{2}\right) + \frac{1}{2} \left(\frac{1}{4}\right) = \frac{1}{6} = \frac{1}{3!} \quad (\text{satisfied})$$

$$\overbrace{A^3 = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1/4 & 0 & 0 & 0 \end{bmatrix}}, \quad b^T A^3 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \frac{1}{4} b_4 = \frac{1}{4} \cdot \frac{1}{6} = \frac{1}{24} = \frac{1}{4!} \quad (\text{satisfied})$$

necessary but not sufficient

(RK4 is a candidate for 4th order, but it isn't necessarily, ie we haven't shown this yet)

REMARKS

- ① The order conditions given by the two special cases are necessary and sufficient for $s=p \leq 3$.
- ② For $p=s=4$ also require
 $b^T A \begin{bmatrix} c_1^2 \\ c_2^2 \\ c_3^2 \end{bmatrix} = \frac{1}{12}$ (covered by general case)
- ③ and $\{b, c_1, b_2c_2, \dots, b_sc_s\} \triangleq A \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_s \end{bmatrix} = \frac{2}{4!}$ (not covered by the general case)
IVP must consider more general to derive Kutta's condition

- ④ see text for $p=5$

s	1	2	3	4	5	6	7	8	9	10
P_{\max}	1	2	3	4	4	5	6	6	7	8+

Region of Absolute stability

Test equation: $y' = \lambda y$

Apply an s -stage RK method:

$$y_i = y_{n-1} + h\lambda \sum_{j=1}^s a_{ij} y_j, \quad i=1, \dots, s$$

$$y_n = y_{n-1} + h\lambda \sum_{i=1}^s b_i y_i = y_{n-1} + h\lambda b^T y$$

$$\begin{aligned} \underline{y} &= \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} y_{n-1} + h\lambda \underline{A} \underline{y} \rightarrow [\underline{I} - h\lambda \underline{A}] \underline{y} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} y_{n-1} \\ &\rightarrow \underline{y} = [\underline{I} - h\lambda \underline{A}]^{-1} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} y_{n-1} \end{aligned}$$

$$\Rightarrow y_n = y_{n-1} + h\lambda b^T [\underline{I} - h\lambda \underline{A}]^{-1} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} y_{n-1}$$

$$y_n = (\underline{I} + h\lambda b^T (\underline{I} - h\lambda \underline{A})^{-1} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}) y_{n-1}$$

$$\text{use } (\underline{I} - h\lambda \underline{A})^{-1} = \underline{I} + h\lambda \underline{A} + (h\lambda)^2 \underline{A}^2 + \dots$$

$$\text{remember, } b^T \underline{A}^k \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = \frac{1}{k!}, \quad k=1, \dots, p$$

$$\Rightarrow y_n = \left[\underline{I} + \lambda h + \frac{(\lambda h)^2}{2!} + \dots + \frac{(\lambda h)^p}{p!} \right] y_{n-1} + \underbrace{\sum_{j>p} (\lambda h)^j b^T \underline{A}^{j-1} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}}_{\text{only if } s>p, \text{ typically, for } p>4}$$

If explicit, then $\underline{A}^{j-1} = \underline{0}$ for $j > s$

$$\text{let } \lambda h = z$$

$$\Rightarrow y_n = \left[\underline{I} + z + \frac{z^2}{2!} + \dots + \frac{z^p}{p!} + \sum_{j=p+1}^s z^j b^T \underline{A}^{j-1} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \right] y_{n-1}$$

G

only if $s>p$, typically, for $p>4$

$$\text{If } s=p \leq 4 \text{ then } G = \underline{I} + z + \frac{z^2}{2!} + \dots + \frac{z^p}{p!}$$

then $z + |G| \leq 1 \Rightarrow \text{region of absolute stability}$

to plot region of stability, select ω , find $z + \omega(t) = e^{i\theta}$
 OR create grid in z -plane, compute $|G(z)|$, plot contour
 where $|G(z)| = 1$

NO EXPLICIT RUNGE-KUTTA METHOD IS A-STABLE

Error Control

choice of the step size is done by some error control scheme. The main task is to estimate the error and pick \tilde{h} st local error is less than some user defined tolerance.

consider a pair of single-step methods

$$y_n = y_{n-1} + h \phi(t_{n-1}, y_{n-1}, h), \text{ } p^{\text{th}} \text{ order method}$$

and

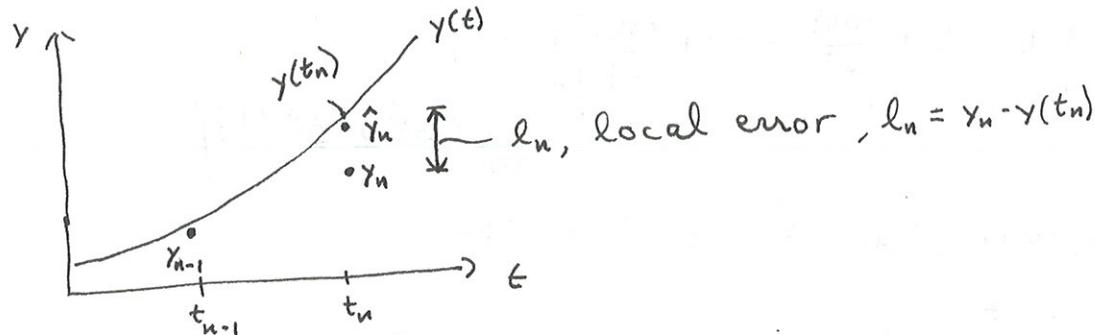
$$\hat{y}_n = y_{n-1} + h \hat{\phi}(t_{n-1}, y_{n-1}, h), \text{ } p+1 \text{ order method}$$

Then the local errors are $|y_n - y(t_n)| \leq \dots = O(h^{p+1})$

$$|\hat{y}_n - y(t_n)| = O(h^{p+2})$$

where $y(t)$ solves the ODE $y' = f(t, y)$, $y(t_{n-1}) = y_{n-1}$

(see pg 22 of notes - if numerical method consistent of order p then $|l_n| = O(h^{p+1})$ as $h \rightarrow 0$)



use $|l_n| \approx |y_n - \hat{y}_n|$

IF $|y_n - \hat{y}_n| \leq \delta$ (some chosen tolerance, or $\delta/|y_n|$ - relative tol)
then the step to t_n would be accepted. If not then
 y_n is rejected and a new h (smaller) is needed.

$$|y_n - \hat{y}_n| = Ch^{p+1} + \dots$$

want \tilde{h} st $|y_n - \hat{y}_n| = \delta = C \tilde{h}^{p+1}$

$$\text{Eliminate } C \Rightarrow \left(\frac{\tilde{h}}{h}\right)^{p+1} = \frac{\delta}{|y_n - \hat{y}_n|} \Rightarrow \tilde{h} = h \left(\frac{\delta}{|y_n - \hat{y}_n|}\right)^{\frac{1}{p+1}}$$

Remarks

- 1) May include a "Fudge" Factor θ to be conservative. Also include limits on h . $\hat{h} = h \left(\frac{\theta \delta}{|y_n - \hat{y}_n|} \right)^{\frac{1}{p+1}}$
- 2) Even if the step is accepted, calculate \hat{h} and use it for the next step.

For Runge Kutta Methods, the computational work to compute y_n and \hat{y}_n can be reduced if ϕ , $\hat{\phi}$ share stage calculations.

For ϕ , might have

$$\begin{array}{c|cc} c & \underline{A} \\ \hline & b^T \end{array}$$

\Rightarrow so ϕ and $\hat{\phi}$ share c and \underline{A}

For $\hat{\phi}$, want

$$\begin{array}{c|cc} c & \underline{A} \\ \hline & \hat{b}^T \end{array}$$

Combine to write :

$$\begin{array}{c|cc} c & \underline{A\phi} \\ \hline & b^T \\ & \hat{b}^T \end{array}$$

Simplest Example: $y_n = y_{n-1} + h f(t_{n-1}, y_{n-1})$ Euler $p=1$

$$\hat{y}_n = y_{n-1} + \frac{h}{2} \left(f(t_{n-1}, y_{n-1}) + f(t_n, y_n) \right), \text{ Modified Euler}$$

$$\begin{array}{c|cc} 0 & 0 & 0 \\ 1 & 1 & 0 \\ \hline & 1 & 0 \\ & \frac{1}{2} & \frac{1}{2} \end{array}$$

Two-stage
order 1/2 pair.

Very popular example is the
Runge-kutta Fehlberg 4/5 pair. (

0	0					
1/4	x	0				
3/8	x	x	0			
12/13	x	x	x	0		
1	x	x	x	x	0	
1/2	x	x	x	x	x	0
	25	0	x	x	x	x
	216					
	x	0	x	x	x	x

see text for values

Implicit Runge Kutta Methods

implicit methods are used for stiff equations. Many are based on quadrature formulas. $y' = f(t, y)$

$$y_n = y_{n-1} + \int_{t_{n-1}}^{t_n} f(t, y) dt$$

Families of Methods

- Gaussian Methods: select c_i 's for a given s to maximize P ,
 $P = 2s$

e.g., midpoint method, $s=1$, $P=2$

- Radau methods: select c_i 's for a given s with $c_s = 1$ to maximize P . $P = 2s - 1$

e.g., backward euler, $s=1$, $P=1$

- Lebatto methods: select c_i 's for a given s with $c_1 = 0$, $c_s = 1$ to maximize P . $P = 2s - P$

e.g., trapezoidal method, $s=2$, $P=2$

These families are collocation methods,

Pick a set of collocation points

$$0 \leq c_1 < c_2 < \dots < c_s \leq 1$$

Construct a polynomial $\psi(t)$ of degree s or less that "collocates" the ODE. i.e,

$$\begin{aligned} \psi(t_{n-1}) &= y_{n-1} \\ \psi'(t_i) &= f(t_i, \psi(t_i)), \quad i = 1, \dots, s \end{aligned} \quad \left\{ \begin{array}{l} s+1 \text{ constraints on } \psi \\ \text{ } \end{array} \right.$$

where $t_i = t_{n-1} + c_i h$
These constraints define ψ uniquely, then $y_n = \psi(t_n)$.

Show that this collocation method leads to implicit Runge K methods.

$\psi'(t)$ = polynomial of degree $s-1$ or less

$$\psi'(t) = \sum_{j=1}^s \psi'(t_j) L_j(t) \quad \left\{ \begin{array}{l} \text{polynomial in Lagrange interpolating form} \\ (\text{note, the choices for } \\ \subseteq \text{ determine } L) \end{array} \right.$$

$$L_j(t) = \prod_{\substack{k=1 \\ k \neq j}}^s \frac{t - t_k}{t_j - t_k}$$

$$\begin{aligned} \text{Now integrate: } \psi(t) &= \underbrace{y_{n-1}}_{Y_{n-1}} + \int_{t_{n-1}}^t \psi'(z) dz \\ &= y_{n-1} + \int_{t_{n-1}}^t \sum_{j=1}^s \psi'(t_j) L_j(z) dz \\ &= y_{n-1} + \sum_{j=1}^s \psi'(t_j) \int_{t_{n-1}}^t L_j(z) dz \end{aligned}$$

$$\Rightarrow y_n = y_{n-1} + h \sum_{j=1}^s f(t_j, \psi(t_j)) b_j, \quad \text{where } b_j = \frac{1}{h} \int_{t_{n-1}}^{t_n} L_j(z) dz$$

$$\text{Set } \cancel{\partial \psi / \partial t}(\psi(t_i)) \quad \Sigma_i = \psi(t_i) \Rightarrow y_i = y_{n-1} + h \sum_{j=1}^s f(t_j, \Sigma_j) a_{ij}$$

$$a_{ij} = \frac{1}{h} \int_{t_{n-1}}^{t_i} L_j(z) dz$$

Examples : if $c_1 = \frac{1}{2} - \frac{\sqrt{3}}{6}$, $c_2 = \frac{1}{2} + \frac{\sqrt{3}}{6}$

(these are the Gauss nodes for the interval $[0, 1]$)

Find :

$$\begin{array}{c|cc} c_1 & \frac{1}{4} & \frac{3-2\sqrt{3}}{12} \\ c_2 & \frac{3+2\sqrt{3}}{12} & \frac{1}{4} \\ \hline & \frac{1}{2} & \frac{1}{2} \end{array}$$

} 2-stage Gauss method, $P=4$

Behavior for stiff Equations

Test equation, $y' = \lambda y$

apply implicit RK method

Generally, $\underline{\underline{C}} \underline{\underline{A}}$

we had, $y_n = R(z)y_{n-1}$, $z = \lambda h$, where $R(z) = I + z \underline{\underline{b}}^T \left[I - z \underline{\underline{A}} \right]^{-1} (I)$

For an explicit method then R is a polynomial in z .

This implies that no explicit method is A -stable, because regions of absolute stability have to be bounded. For implicit methods, R = rational function of z , $R(z) = \frac{P(z)}{Q(z)}$

\Rightarrow the region of absolute stability may be unbounded.

It turns out that the families of collocation methods (Gauss, Radau, Lebatto) are all A -stable. Additionally, if the degree of Q is greater than the degree of P , then $R \rightarrow 0$ as $z \rightarrow \infty \Rightarrow$ method has stiff decay, (like backward Euler). Radau methods have stiff decay.

One way to obtain stiff decay is if $\underline{\underline{A}}$ is nonsingular and the last row of $\underline{\underline{A}}$ is equal to $\underline{\underline{b}}^T$, then this will have stiff decay.

$$A = \begin{bmatrix} \dots & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \\ \hline & & \underline{\underline{b}}^T \end{bmatrix}, \quad \tilde{A} = \underbrace{\begin{bmatrix} \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}}_{s-1 \text{ columns orthogonal to } \underline{\underline{b}}} \begin{bmatrix} \underline{\underline{b}} \\ \|\underline{\underline{b}}\|^2 \end{bmatrix}$$

$$R = I + z b^T (\underline{I} - z \underline{A})^{-1} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

$$\rightarrow I - z b^T \left(\frac{1}{z} \underline{A}^{-1} \right) \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = I - b^T \underline{A}^{-1} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \quad (\text{assuming } z \underline{A} \gg I)$$

(this is in the limit as $z \rightarrow \infty$)

Notice, $b^T \underline{A}^{-1}$ is a row vector, $[0, 0, \dots, 1]$

$$\Rightarrow \lim_{z \rightarrow \infty} = I - b^T \underline{A}^{-1} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = I - (0, 0, \dots, 0, 1) \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = 0$$

backward Euler: $\frac{s+1}{s} \quad (s=1)$

some Radav method

$$\begin{array}{c|cc} y_3 & 3/12 & -1/12 \\ \hline 1 & 3/4 & 1/4 \\ \hline & 1/4 & 1/4 \end{array} \quad (s=2, p=3)$$

Implicit Runge-Kutta Methods

3/4/03

Families of Implicit Runge-Kutta methods,
derived using polynomial collocation:

Gauss: $p=2s$ (symmetric) mid.

Radav: $p=2s-1$ (stiff decay) b.euler

Labatto: $p=2s-2$ (symmetric) trap } a-stable

Best for stiff eqns.

$$\begin{array}{c|c} c & A \\ \hline & b^T \end{array}$$

◆ test eqn, $y' = \lambda y$

Apply RK method: $y_n = R(z)y_{n-1}, z = \lambda h$

$$R(z) = I + z b^T (\underline{I} - z \underline{A})^{-1} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

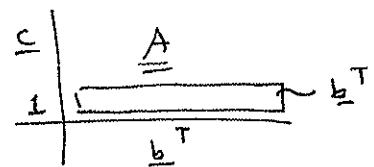
explicit $\rightarrow R = \text{polynomial in } z$

implicit $\rightarrow R = \text{rational fct of } z$

\Rightarrow possible to have an unbounded
region of absolute stability

All members of the families above (Gauss, Radav,
Labatto) are A stable.

Also may have stiff decay if last row of A is b^T



Diagonally Implicit RK Methods (DIRK)

main difficulty - must solve a system of nonlinear algebraic equations at each time step.

General Form:

$$Y_i = y_{n-1} + h \sum_{j=1}^s a_{ij} F(t_j, Y_j), \quad i=1, \dots, s$$

$$\text{where } t_j = t_{n-1} + c_j h$$

then

$$y_n = y_{n-1} + h \sum_{i=1}^s b_i F(t_i, Y_i)$$

Need to solve for the Y_i 's. Consider some

linearization: $Y_i = \hat{Y}_i + \delta_i$ \leftarrow correct to final
↑
current
guess

$$\begin{aligned} \hat{Y}_i + \delta_i &= y_{n-1} + h \sum_{j=1}^s a_{ij} F(t_j, \hat{Y}_j + \delta_j) \\ &= y_{n-1} + h \sum_{j=1}^s a_{ij} \left[F(t_j, \hat{Y}_j) + F_y(t_j, \hat{Y}_j) \delta_j + \dots \right] \end{aligned}$$

I ignore dots, solve for δ_i :

$$\underbrace{\delta_i - h \sum_{j=1}^s a_{ij} F_y(t_j, \hat{Y}_j) \delta_j}_{\text{linear system for } \delta} = y_{n-1} - \hat{Y}_i + h \sum_{j=1}^s a_{ij} F(t_j, \hat{Y}_j)$$

$$\left[\begin{array}{cccc} I - a_{11} h F_y(t_1, \hat{Y}_1) & -a_{12} h F_y(t_2, \hat{Y}_2) & \cdots & \\ -a_{21} h F_y(t_1, \hat{Y}_1) & I - a_{22} h F_y(t_2, \hat{Y}_2) & \cdots & \\ \vdots & \vdots & \ddots & \end{array} \right] \left[\begin{array}{c} \delta_1 \\ \delta_2 \\ \vdots \\ \delta_s \end{array} \right] = \text{mx1 vector}$$

each element an $m \times m$ matrix

so $m \times m$ system

Consider various approaches to solving

For the s 's.

- 1) Freeze calculation and LU-factorization of the Jacobian
- 2) Design its method and its a_{ij} 's s.t. the Jacobian has a simpler structure.

One choice: $a_{ij} = 0$ for $j > i$

$$a_{ii} = \alpha \quad \text{for } i = i'$$

\Rightarrow diagonally implicit runge kutta methods (DIRK)

$$y_i = y_{n-1} + h \alpha f(t_i, y_i) + h \sum_{j=1}^{i-1} a_{ij} f(t_j, y_j)$$

explicit terms

$$\text{set } Y_i = \tilde{Y}_i + \delta_i$$

$$\tilde{Y}_i + \delta_i = y_{n-1} + h \left(f(t_i, \tilde{Y}_i) + f_y(t_i, \tilde{Y}_i) \delta_i + \dots \right) + h \sum_{j=1}^{i-1} a_{ij} f(t_j, Y_j)$$

$$\Rightarrow \underbrace{\left(I - h f_y(t_i, \tilde{Y}_i) \right) \delta_i}_{\text{new matrix}} = \text{known}$$

some "structure" for each i

use $(I - dh f_y(t_{n-1}, y_{n-1})) \delta_i = \text{known}$, for $i = 1, \dots, s$

minimize updates and recalc's of LU

DIRK Methods:

$$p = s + 1$$

$$\begin{array}{c|cc} & v_2 & \\ \hline & 1/2 & \\ & & 1 \end{array} \rightarrow \text{midpoint, } s = 1, p = 2$$

$$\begin{array}{c|cc} & \gamma & \\ \hline 1-\gamma & 1-2\gamma & \gamma \\ & \hline & \gamma & 1/2 \end{array}$$

$$s = 2, p = 3, \gamma = \frac{3 + \sqrt{3}}{6}$$

A-stable

Stiff Decay, $p = 5$

$$\text{eg } \begin{array}{c|cc} & 1 & \\ \hline & 1 & \\ & 1 & \end{array} \quad s = p = 1, \text{ b. euler}$$

$$\begin{array}{c|cc} & \gamma & \\ \hline 1-\gamma & 1-\gamma & \gamma \\ & \hline & 1-\gamma & \gamma \end{array} \quad s = p = 2, \gamma = \frac{2 - \sqrt{2}}{2} \quad (\text{notice last row } A \text{ and } b^T)$$

Midterm : Friday, in class, 2 hrs
open notes
no book
covers everything through RK

Topics:

1) Intro

- a) IVP, BVP, BAE
- b) examples - phase plane, initial layers
- c) regularity of IVPs, Lipschitz
- d) BVP: boundary layer problem

2) IVP

- a) problem stability
 - scalar test
 - linear coeff (constant) system
 - variable coeff sys
 - non-linear

b) Hamiltonian systems

- c) Basic numerical methods and concepts
 - illustrated with Euler's method

- local truncation error
 - consistency \rightarrow
 - order of accuracy
 - global error
 - convergence
 - O-stability

Thm: consistent + O-stability \Rightarrow convergent

- local error vs local truncation error
- absolute stability ($y' = \lambda y$) \rightarrow regions of absolute stability in the complex λh plane

d) stiffness, implicit methods

- backward Euler
- A stability,
- stiff decay
- solving non-linear eqns
 - Newton's method

Trapezoidal rule - A-stable

e) single-step methods

$$y_n = y_{n-1} + h \phi(t_{n-1}, y_{n-1}, h)$$

- Taylor's method

- Runge Kutta Methods

- s-stage method $\Rightarrow \frac{c}{\Delta}$

$$\text{- consistency} \Rightarrow \sum_{i=1}^s a_{ij} = c_i, \quad \sum_{i=1}^s b_i = 1$$

- order conditions \leftarrow how obtained, special forms

- regions of absolute stability

- Error control - For explicit methods

local error estimates

step size selection

- Implicit Methods

Gauss, Radau, Lobatto

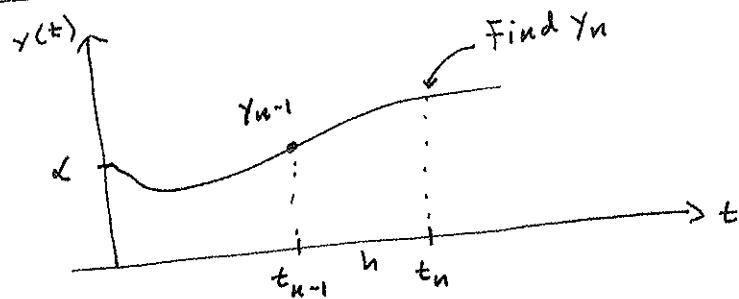
collocation methods

- DIRK Methods

explicit
methods

Linear Multistep Methods

IVP: $y' = f(t, y)$, $0 \leq t \leq b$, $y(0) = \alpha$



single step method uses data at y_{n-1} to calculate y_n
 whereas multi-step method might use data at
 y_{n-2}, y_{n-3}, \dots

\Rightarrow A k -step multistep method uses previous k solutions
 values to compute y_n .

General form

$$\sum_{j=0}^k \alpha_j y_{n-j} = h \sum_{j=0}^k \beta_j f_{n-j}$$

"linear" k -step
multistep method

$\{\alpha_j\}, \{\beta_j\}$ coefficients for the method

set $\alpha_0 = 1 \Rightarrow y_n = \text{stuff}$

Assume $|\alpha_j| + |\beta_j| \neq 0$

define $f_j \equiv f(t_j, y_j)$

Observations

1) "linear" merely implies that y 's, f 's appear

as a linear combination

2) Fixed h (traditionally). Typically, methods are derived using fixed h so that α 's, β 's independent of h . (Though we will discuss varying h for error control.)

3) The method is explicit if $\beta_0 = 0$ and implicit otherwise.

Adams Family

The Adams family is popular for non-stiff ODEs.

Form:
$$y_n = y_{n-1} + h \sum_{j=0}^k \beta_j f_{n-j}$$
, $\alpha_0 = -\alpha_1 = 1$, $\alpha_j = 0$ for $j \geq 2$

The β_j 's are chosen for accuracy.

Amongst this Family of methods there are 2 subclasses:

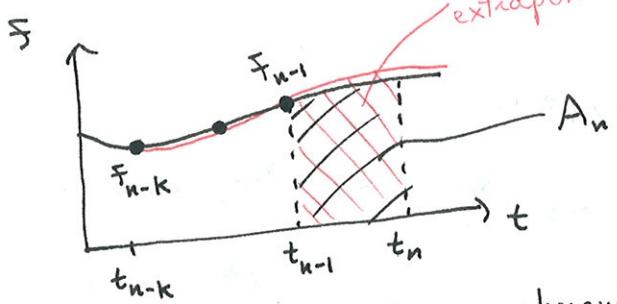
1) Adams-Basforth - explicit methods ($\beta_0 = 0$)
order of accuracy $\equiv p = k$

2) Adams-Moulton - implicit methods ($\beta_0 \neq 0$)
order of accuracy $\equiv p = k+1$

derivation

$$y' = f(t, y), \text{ integrate: } y(t_n) = y(t_{n-1}) + \int_{t_{n-1}}^{t_n} f(t, y(t)) dt$$

Need to approximate the integral



area under curve
extrapolated

interpolate data and
extrapolate for t_{n-1} to t_n

\Rightarrow Fit an interpolating polynomial to the data
 $(t_{n-j}, f_{n-j}), j = 1, \dots, k$

Lagrange polynomials

$$\text{Interpolant: } \tilde{f}(t) = \sum_{j=1}^k f_{n-k} \tilde{L}_j(t)$$

$$\text{where } \tilde{L}_j(t) = \prod_{\substack{i=1 \\ i \neq j}}^k \frac{(t - t_{n-i})}{(t_{n-j} - t_{n-i})} \quad \leftarrow (k-1) \text{ degree polynomial}$$

$$A_n \approx \int_{t_{n-1}}^{t_n} \tilde{f}(t) dt = \sum_{j=1}^k f_{n-k} \underbrace{\int_{t_{n-1}}^{t_n} \tilde{L}_j(t) dt}_{h \beta_j}$$

$$\Rightarrow \beta_j = \frac{1}{h} \int_{t_{n-1}}^{t_n} \tilde{L}_j(t) dt$$

For example (remember $\beta_0 = 0$)

AB1: , $\beta_1 = 1$ (this is just Euler's method)

AB2: , $\beta_1 = \frac{3}{2}$, $\beta_2 = -\frac{1}{2}$

AB3: $\beta_1 = \frac{23}{12}$, $\beta_2 = -\frac{16}{12}$, $\beta_3 = \frac{5}{12}$

For example, the AB3 method becomes

$$y_n = y_{n-1} + \frac{h}{12} [23f_{n-1} - 16f_{n-2} + 5f_{n-3}]$$

Accuracy (more later)

$$d_n = C_{p+1} h^p y^{(p+1)}(t_n) + O(h^{p+1})$$

$$\text{AB: } p=k \text{ so } C_2 = \frac{1}{2}, C_3 = \frac{5}{12}, C_4 = \frac{3}{8}, \dots$$

\uparrow
3rd order

Adams - Moulton fits a k^{th} degree polynomial to the data (t_{n-j}, f_{n-j}) , $j = 0, \dots, k$

$$\hat{f}(t) = \sum_{j=0}^k f_{n-k} \hat{L}_j(t)$$

$$\hat{L}_j(t) = \prod_{\substack{i=0 \\ i \neq j}}^k \frac{(t - t_{n-i})}{(t_{n-j} - t_{n-i})}$$

$$\text{so the AM } \beta_j \text{'s are } \beta_j = \frac{1}{h} \int_{t_{n-1}}^{t_n} \hat{L}_j(t) dt$$

For example: (trapezoidal rule)

$$\text{AM1 : } \beta_0 = \frac{1}{2}, \quad \beta_1 = \frac{1}{2}$$

$$\text{AM2 : } \beta_0 = \frac{5}{12}, \quad \beta_1 = \frac{8}{12}, \quad \beta_2 = \frac{-1}{12}$$

(see text for others)

Accuracy $d_n = C_{p+1} h^p y^{(p+1)}(t_n) + \dots$

$$p = k+1$$

$$C_3 = \frac{-1}{12}, \quad C_4 = \frac{-1}{24}$$

\uparrow
3rd order

Remember,
Adams Family
For nons stiff eqns...

Backward Difference Formulas

Linear Multistep Methods

$$y' = f(t, y), \quad y(0) = c, \quad 0 \leq t \leq b$$

form: $\sum_{j=0}^k \alpha_j y_{n-j} = h \sum_{j=0}^k \beta_j f_{n-j}$, where $f_{n-j} = f(t_{n-j}, y_{n-j})$

for the Adams Method: $\alpha_0 = 1, \alpha_1 = -1, \alpha_j = 0, j \geq 2$

$$\Rightarrow y_n = y_{n-1} + h \sum_{j=0}^k \beta_j f_{n-j}$$

if $\beta_0 = 0 \Rightarrow$ Adams-Basforth (explicit) $p=k$ (order)

if $\beta_0 \neq 0 \Rightarrow$ Adams-Moulton (implicit) $p=k+1$

BDF Methods - Backwards Difference Formulas (stiff problems)

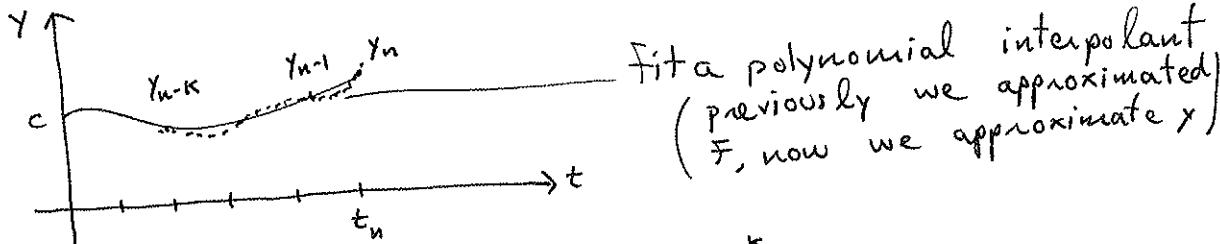
form: $\sum_{j=0}^k \alpha_j y_{n-j} = h \beta_0 f_n, \quad \alpha_0 = 1$

start with differential eqn evaluated at t_n

$$y'(t_n) = f(t_n, y(t_n))$$

now replace $y'(t_n)$ with a backward difference formula
(this naturally derives an implicit method)

Derivation: use the data $(t_{n-i}, y_{n-i}), i=0, \dots, k$



$$\Rightarrow y(t) \approx \sum_{j=0}^k y_{n-j} \hat{L}_j(t), \quad \hat{L}_j(t) = \prod_{\substack{i=0 \\ i \neq j}}^k \frac{t - t_{n-i}}{t_{n-j} - t_{n-i}}$$

differentiate: $y'(t_n) \approx \sum_{j=0}^k y_{n-j} \hat{L}'_j(t_n)$

$$\Rightarrow \sum_{j=0}^k y_{n-j} \hat{L}'_j(t_n) = f(t_n, y_n), \text{ normalize by } \hat{L}'_0(t_n)$$

$$\Rightarrow \alpha_j = \frac{\hat{L}_j'(t_n)}{\hat{L}_0'(t_n)}, \quad \beta_0 = \frac{1}{h} \frac{1}{\hat{L}_0'(t_n)}$$

(44)

Suppose $k=2$

$$\hat{L}_0(t) = \frac{t-t_{n-1}}{t_n-t_{n-1}} \cdot \frac{t-t_{n-2}}{t_n-t_{n-2}} = \frac{1}{2h^2}(t-t_{n-1})(t-t_{n-2})$$

$$\hat{L}_0'(t) = \frac{1}{2h^2}(2t-t_{n-1}-t_{n-2})$$

$$\hat{L}_0'(t_n) = \frac{1}{2h^2}(2t_n-t_{n-1}-t_{n-2}) = \frac{1}{2h^2}(t_n-t_{n-1}+t_n-t_{n-2}) = \cancel{\frac{3}{2}h} \quad \frac{3}{2h}$$

$$\hat{L}_1(t_n) = \frac{t-t_n}{t_{n-1}-t_n} \cdot \frac{t-t_{n-2}}{t_{n-1}-t_{n-2}} \rightarrow \hat{L}_1'(t_n) = -\frac{2}{h}$$

$$\hat{L}_2(t) = \frac{t-t_n}{t_{n-2}-t_n} \cdot \frac{t-t_{n-1}}{t_{n-2}-t_{n-1}} \rightarrow \hat{L}_2'(t_n) = \frac{1}{2h}$$

$$\Rightarrow \alpha_1 = -4/3, \quad \alpha_2 = 1/3, \quad \beta_0 = 2/3$$

$$\Rightarrow \text{BDF2: } y_n = \frac{4}{3}y_{n-1} - \frac{1}{3}y_{n-2} + \frac{2}{3}h f(t_n, y_n)$$

Notice, this is an implicit method
see text for others...

• Order of accuracy is $p=k$

Example:

consider the initial value problem

$$y' = e^{-t} - y, \quad y(0) = 0, \quad 0 \leq t \leq 4$$

Compare solutions of linear multistep methods

for this IVP. The exact solution is $y(t) = t e^{-t}$.

starting values: (this is "cheating")

$$y(0) = 0, \quad y_j = t_j e^{-t_j}, \quad j = 1, \dots, k-1$$

$$\text{Adams - Moulton} \quad y_n = y_{n-1} + h \sum_{j=0}^k \beta_j F_{n-j}$$

$$\rightarrow y_n = y_{n-1} + h \beta_0 F_n + h \sum_{j=1}^k \beta_j F_{n-j}$$

$$F_n = e^{-t_n} - y_n$$

$$\rightarrow (1+h\beta_0) y_n = y_{n-1} + h \beta_0 e^{-t_n} + h \sum_{j=1}^k \beta_j F_{n-j}$$

$$\rightarrow y_n = \frac{1}{1+h\beta_0} \left[y_{n-1} + h \beta_0 e^{-t_n} + h \sum_{j=1}^k \beta_j F_{n-j} \right]$$

} another "cheat"
to avoid implicit
calculations - turn
into explicit using
exact solution

Observations

1) Adams-Moulton order p is generally more accurate than Adams-Basforth for the same h .

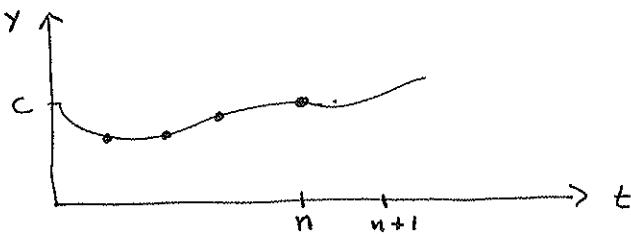
2) Compare back to RK methods.

AB4, $h = 0.00625$, error = 8.8×10^{-10}
 AM3, $h = 0.00625$, error = 6.6×10^{-11}
 RK4, $h = 0.00625$, error = 6.8×10^{-12}

(For this particular IVP)

Starting Values

$$\sum_{j=0}^k \alpha_j y_{n-j} = h \sum_{j=0}^k \beta_j F_{n-j}, \text{ where } n \geq k$$



Require values for y_i , $i = 1, \dots, k-1$ to get started

2 Basic Approaches

- 1) use a single step method of correct order
- 2) Bootstrap using multistep methods of the same family but with lower order. Need some form of error extrapolation.

Consider the Adams-Basforth Family

AB1 \Rightarrow Euler

$$y(h) = \underbrace{y(0)}_{Y_{AB1}(h)} + h F(0) + Ch^2 + \dots \quad (1)$$

$$y(h) = Y_{AB1}(h/2) + C(h/2)^2 + \dots \quad (2)$$

$$4(2) - (1) \Rightarrow 3y(h) = 4Y_{AB1}(h/2) - Y_{AB1}(h) + O(h^3)$$

$$y(h) = \underbrace{\frac{4Y_{AB1}(h/2) - Y_{AB1}(h)}{3}}_{O(h^2)} + O(h^3)$$

$O(h^2)$ accurate approximation

45

Order of Accuracy

$$\text{Define } N_n y_n = \frac{1}{n} \sum_{j=0}^k \alpha_j y_{n-j} - \sum_{j=0}^K \beta_j f_{n-j}$$

Local truncation error

$$d_n = N_h y(t_n) = \frac{1}{n} \sum_{j=0}^K \alpha_j y(t_{n-j}) - \underbrace{\sum_{j=0}^K \beta_j f(t_{n-j}) y(t_{n-j})}_{y'(t_{n-j})}$$

Expand using Taylor series.

$$y(t_{n-j}) = y(t_n - jh) = \sum_{m=0}^{\infty} \frac{(-jh)^m}{m!} y^{(m)}(t_n)$$

$$y'(t_{n-j}) = \sum_{m=0}^{\infty} \frac{(-j\hbar)^m}{m!} y^{(m+1)}(t_n)$$

$$\Rightarrow d_n = \frac{1}{n} \sum_{j=0}^K \alpha_j \sum_{m=0}^{\infty} \frac{(-jh)^m}{m!} y^{(m)}(t_n) - \sum_{j=0}^K \beta_j \sum_{m=0}^{\infty} \frac{(-jh)^m}{m!} y^{(m+1)}(t_n)$$

$$= (+)h^{-1} + (-)h^0 + (+)h^1 + \dots$$

"
O
for consistency

$\stackrel{0}{\rightarrow}$ For order of accuracy \rightarrow

$$\rightarrow d_n = \sum_{m=0}^{\infty} \left(\sum_{j=0}^k \alpha_j (-j)^m \right) \frac{h^{m+1}}{m!} Y^{(m)}(t_n) - \sum_{m=0}^{\infty} \left(\sum_{j=0}^k \beta_j (-j)^m \right) \frac{h^m}{m!} Y^{(m+1)}(t_n)$$

shift of m-index

$$d_n = \sum_{m=-1}^{\infty} \left(\sum_{j=0}^k \alpha_j (-j)^{m+1} \right) \frac{h^m}{(m+1)!} y^{(m+1)}(t_n) - \sum_{m=0}^{\infty} \left(\sum_{j=0}^k \beta_j (-j)^m \right) \frac{h^m}{m!} y^{(m+1)}(t_n)$$

$$h^{-1} : \sum_{i=0}^k d_i y(t_n) = 0$$

⇒

$\sum_{j=0}^k d_j = 0$, for consistency
 (don't forget to consider k^0 term)

$$h^m : \sum_{j=0}^K \alpha_j (-j)^{m+1} \frac{y^{(m+1)}(t_n)}{(m+1)!} - \sum_{j=0}^K \beta_j (-j)^m \frac{y^{(m+1)}(t_n)}{m!} = 0 \quad \begin{matrix} \text{For } p^{\text{th}} \\ \text{order method} \end{matrix}$$

for $m=0, \dots, p-1$

$$\Rightarrow \sum_{j=0}^k a_j (-j)^{m+1} = \sum_{j=0}^k b_j (-j)^m (m+1), \text{ for } m=0, 1, \dots, p-1$$

In summary:

For consistency:

$$\sum_{j=0}^k \alpha_j = 0 \quad \text{AND} \quad -\sum_{j=0}^k j \alpha_j = \sum_{j=0}^k \beta_j$$

Truncation Error:

$$e_n = \left| \frac{1}{(p+1)!} \left(\sum_{j=0}^k \alpha_j (-j)^{p+1} \right) - \frac{1}{p!} \sum_{j=0}^k \beta_j (-j)^p \right| h^p y^{(p+1)}(t_n) + \dots$$

$\underbrace{\phantom{\sum_{j=0}^k \alpha_j (-j)^{p+1}}}_{C_{p+1}}$

Example

use the order conditions to obtain the ~~order~~ obtain the α_j, β_j in ~~the~~ a linear multistep method.

$$y_n = y_{n-1} + h \sum_{j=0}^2 \beta_j f_{n-j} \quad (\text{AB12})$$

in this form, we've already chosen

$$\alpha_0 = 1, \alpha_1 = -1, \alpha_2 = 0 \quad (\text{Notice, } \sum \alpha_j = 0)$$

Problem is to choose $\beta_0, \beta_1, \beta_2$ st the method is 3rd order.

$$\rightarrow \sum_{j=0}^2 (-j)^{m+1} \alpha_j = (m+1) \sum_{j=0}^2 (-j)^m \beta_j, \quad m = 0, 1, 2$$

$$m=0: \quad 1 = \beta_0 + \beta_1 + \beta_2$$

$$m=1: \quad -1 = 2(-\beta_1 - 2\beta_2)$$

$$m=2: \quad 1 = 3(\beta_1 + 4\beta_2)$$

$$\Rightarrow \beta_0 = \frac{5}{12}, \quad \beta_1 = \frac{8}{12}, \quad \beta_2 = \frac{-1}{12}$$

Root Condition $\alpha_0 \xi^k + \alpha_1 \xi^{k-1} + \dots + \alpha_{k-1} \xi + \alpha_k = 0$

stability condition involving roots of this characteristic polynomial. What are the magnitudes of the roots?

O-Stability & The Root Condition
For Linear Multistep Methods

General k -step method gives

$$N_h y_n = \frac{1}{h} \sum_{j=0}^k \alpha_j y_{n-j} - \sum_{j=0}^k \beta_j f_{n-j}$$

O-stability requires that

$$|x_m - z_m| \leq K \left[\sum_{j=0}^{k-1} |x_j - z_j| + \max_{k \leq n \leq N} |N_h x_n - N_h z_n| \right]$$

for $0 \leq m \leq N$, $h \leq h_0$

$$\Rightarrow |x_m - z_m| \leq K \left[\sum_{j=0}^{k-1} |x_j - z_j| + \max_{k \leq n \leq N} \left| \frac{1}{h} \sum_{j=0}^k \alpha_j (x_{n-j} - z_{n-j}) - \sum_{j=0}^k \beta_j (f(x_{n-j}) - f(z_{n-j})) \right| \right]$$

dominant term
in this expression

Consider the behavior of solutions to the difference eqn

$$\sum_{j=0}^k \alpha_j v_{n-j} = 0$$

$$\Rightarrow \alpha_0 v_n + \alpha_1 v_{n-1} + \dots + \alpha_{k-1} v_{n-k+1} + \alpha_k v_{n-k} = 0$$

constant coeff linear difference equation.

set $v_n = \xi^n$, where $\xi = (\text{complex}) \text{ constant}$

$$\text{or equivalently, } v_n = e^{n \ln \xi} = e^{t_n \left(\frac{\ln \xi}{h} \right)}, t_n = nh$$

→ exponential solution for linear constant coeff difference eqn

$$\rightarrow \alpha_0 \xi^n + \alpha_1 \xi^{n-1} + \dots + \alpha_k \xi^{n-k} = 0$$

$$\xi^{n-k} (\alpha_0 \xi^k + \alpha_1 \xi^{k-1} + \dots + \alpha_k) = 0$$

ξ characteristic polynomial
of the difference eqn. There are k
roots of $P(\xi) = 0$, $\xi_1, \xi_2, \dots, \xi_k$.

general solution: $v_n = \sum_{m=1}^k c_m \xi_m^n$, c_m determined by initial conditions

Remarks

- 1) If we assume that multistep method is consistent
 Then $\sum_{j=0}^k \alpha_j = 0$ and $p(1) = \sum_{j=0}^k \beta_j = 0$
 $\Rightarrow \xi = 1$ is always a root of the characteristic polynomial for a consistent multistep method
 $\Rightarrow v_n = c_1 + \sum_{m=2}^k c_m \xi_m^n$
- 2) If the roots are repeated then the form changes slightly.
 e.g., suppose $\xi_2 = \xi_3 = \hat{\xi}$
 then $v_n = c_1 + (c_2 + nc_3) \hat{\xi}^n + \sum_{m=4}^k c_m \xi_m^n$
 triple root $\Rightarrow \hat{\xi}^n, n\hat{\xi}^n, n^2\hat{\xi}^n$, and so on...

3) terminology:

$\xi_1 = 1 \Rightarrow$ principal root

$\xi_m, m=2, \dots, k \Rightarrow$ extraneous roots

The extraneous roots determine stability

Root Condition

Require that $|\xi_m| \leq 1$ for all m for stability,
 (and if $|\xi_m| = 1$, ξ_m must be a simple root, i.e.
 non-repetable).

"strongly stable" \Rightarrow ~~weakly~~ $|\xi_m| < 1$ for $m \geq 2$

"weakly stable" \Rightarrow O-stable, but not strongly stable

Examples

① Adams Methods all have

$$\alpha_0 = 1, \alpha_1 = -1, \alpha_j = 0 \text{ for } j \geq 2$$

$$\Rightarrow p(s) = s^k - s^{k-1}$$

$$= s^{k-1}(s-1)$$

$$\Rightarrow \xi_1 = 0, \xi_m = 0, m \geq 2 \Rightarrow \boxed{\text{strongly stable}}$$

② "leap Frog" - 2 step method (2 roots to char poly)

$$y_n = y_{n-2} + 2h f_{n-1}$$

$$p(s) = s^2 - 1 \Rightarrow \xi_1 = 1, \xi_2 = -1 \quad \begin{matrix} \nearrow \\ \downarrow \\ \text{extraneous root sitting on boundary} \end{matrix}$$

$$\boxed{\text{weakly stable}}$$

③ $y' = f(t, y)$

so let take $y'(t_{n-1}) = f(t_{n-1}, y(t_{n-1}))$
and now approximate $y'(t_{n-1})$ using some finite difference

$$\Rightarrow \frac{1}{14h} (5y_n + 6y_{n-1} - 13y_{n-2} + 2y_{n-3}) = f_{n-1}$$

this is a 3-step method

$$5y_n + 6y_{n-1} - 13y_{n-2} + 2y_{n-3} = 14h f_{n-1}$$

$$\Rightarrow p(s) = 5s^3 + 6s^2 - 13s + 2$$

$$= (s-1) \underbrace{(5s^2 + 11s - 2)}_{\text{extraneous roots}}$$

$$\xi_1 = 1$$

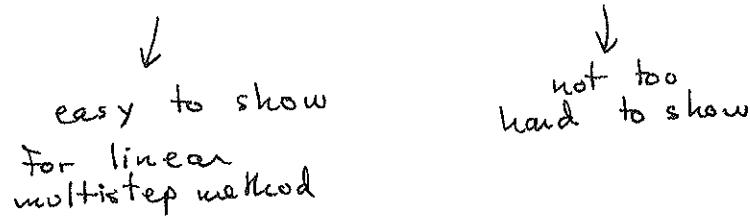
$$\xi_2, \xi_3 = \frac{-11 \pm \sqrt{161}}{10}$$

$$= 0.169, -2.37$$

↑
UNSTABLE

Test these using $y' = e^{-t} - y, y(0) = 0$
see plots on website

Thm : consistent method + O-stability \Rightarrow convergent



Absolute stability

Apply a general k-step method to test equation, $y' = \lambda y$

$$\text{Have } \sum_{j=0}^k \alpha_j y_{n-j} = h \lambda \sum_{j=0}^k \beta_j y_{n-j}$$

$$\text{let } z = h \lambda$$

$$\Rightarrow \sum_{j=0}^k \alpha_j y_{n-j} = z \sum_{j=0}^k \beta_j y_{n-j}$$

$$\text{set } y_n = s^n \Rightarrow \sum_{j=0}^k \alpha_j s^{n-j} = z \sum_{j=0}^k \beta_j s^{n-j}$$

$$\Rightarrow \underbrace{\sum_{j=0}^k \alpha_j s^{n-j}}_{p(s)} = z \underbrace{\sum_{j=0}^k \beta_j s^{n-j}}_{r(s)}$$

generally, $p(s)$ and $r(s)$ are k-degree polynomials

(if explicit, then $(k-1)$ degree polynomials)

roots s_1, \dots, s_k of the polynomial

$$p(s) - z r(s) = 0$$

for a given complex value for z

$$\text{general solution: } y_n = \sum_{m=1}^k c_m s_m^n$$

we want decay \Rightarrow

Region of absolute stability: $\{ z \text{ st } |s_m| \leq 1 \}$

To find region, set $\xi = e^{i\theta}$, $-N \leq \theta \leq N$

then compute $z = \frac{p(e^{i\theta})}{r(e^{i\theta})} = \text{boundary of region}$
of absolute stability

Implementation

(1) predictor/corrector methods

-used for non-stiff problems, Adams-Basforth and
Adams-Moulton methods are often used in a
predictor-corrector way.

predictor step: AB, explicit

example: 3-step

$$\tilde{y}_n = y_{n-1} + \frac{h}{12} [23F_{n-1} - 16F_{n-2} + 5F_{n-3}]$$

↑
predicted value

corrector step: AM, "implicit" (2-step)

$$y_n = y_{n-1} + \frac{h}{12} [5\tilde{f}_n + 8F_{n-1} - F_{n-2}]$$

↑
 $\tilde{f}_n = f(t_n, \tilde{y}_n)$

cost = 2 function evaluations per step

(2) Error estimates

For a predictor-corrector scheme, an error
estimate is natural and comes for free.

predictor step: $\tilde{y}_n - y(t_n) = \text{local error} = h d_n + \dots$
(AB)

$$\text{where } d_n = \tilde{c}_{p+1} h^p y^{(p+1)}(t_n) + \dots$$

$$\Rightarrow \tilde{y}_n - y(t_n) = \tilde{c}_{p+1} h^{p+1} y^{(p+1)}(t_n) + \dots$$

corrector step: $y_n - y(t_n) = \underbrace{c_{p+1} h^{p+1} y^{(p+1)}(t_n)}_{\text{assume this even though } y_n \text{ is}} + \dots$
(AM)
obtained using \tilde{f}_n

$$\text{subtract: } \tilde{y}_n - y_n = (\tilde{c}_{p+1} - c_{p+1}) h^{p+1} y^{(p+1)}(t_n) + \dots$$

$$\rightarrow h^{p+1} y^{(p+1)}(t_n) = \frac{\tilde{y} - y_n}{\tilde{c}_{p+1} - c_{p+1}} + \dots$$

substitute this error into local error of corrector step

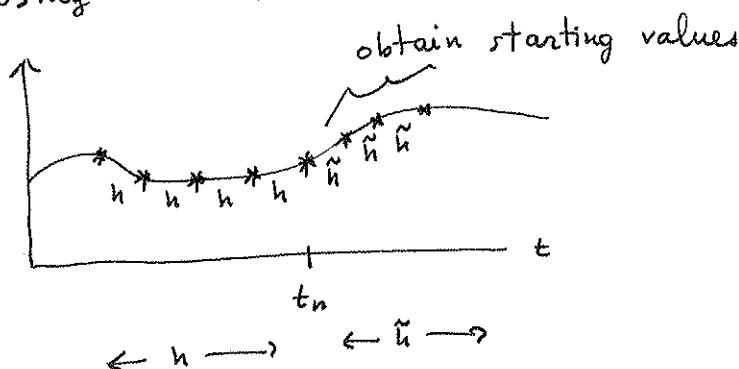
$$\rightarrow y_n - y(t_n) \approx \frac{c_{p+1}}{\tilde{c}_{p+1} - c_{p+1}} (\tilde{y}_n - y_n), \text{ Milne estimate}$$

(3) Variable step size

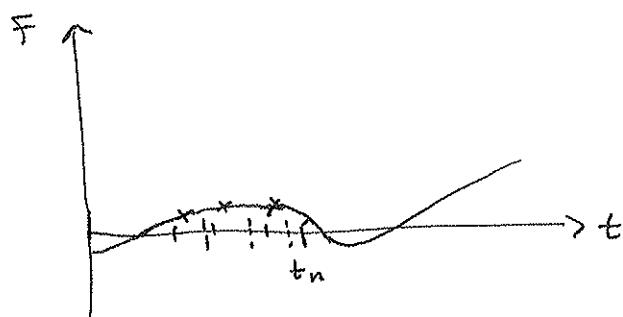
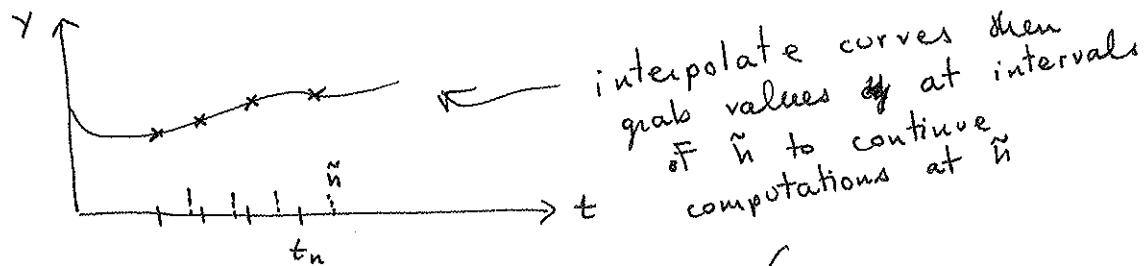
In an error control algorithm you may need to change h according to some error estimate.

Choices to accommodate the change in h .

- [a] restart multi-step method with new h using a single step method.



- [b] Interpolate y or f .



(C) Variable step multistep method

$$\text{had } \sum_{j=0}^k \alpha_j y_{n-j} = h \sum_{j=0}^k \beta_j f_{n-j}$$

where α_j, β_j are found assuming uniform h

Instead find α_j, β_j assuming arbitrary set of nodes t_{n-j} . The formulas are more complicated and involve h_n .

General system of m ODEs.

$$\underline{y}' = \underline{F}(t, \underline{y}), \quad 0 \leq t \leq b$$

with m boundary conditions, $\underline{g}(\underline{y}(0), \underline{y}(b)) = \underline{0}$.

Solution of the BVP is not local, ie for IVP we can write solution locally using a Taylor series but this is not possible for a BVP in general.

This makes the question of existence and uniqueness of solutions more difficult.

Consider a corresponding IVP

$$\underline{y}' = \underline{F}(t, \underline{y}), \quad 0 \leq t \leq b, \quad \underline{y}(0) = \underline{\epsilon}$$

Suppose the IVP is well-posed. The solution exists uniquely and it depends continuously on $\underline{F}, \underline{c}$.

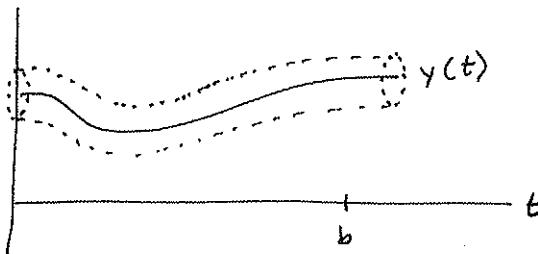
Solution of the IVP is $\underline{y} = \underline{y}(t, \underline{\epsilon})$. For the BVP, $\underline{\epsilon}$ is chosen st $\underline{g}(\underline{\epsilon}, \underline{y}(b, \underline{\epsilon})) = \underline{0}$ \leftarrow m algebraic equations for the m -components of $\underline{\epsilon}$. This eqn may have none, one, or many solutions.

For certain cases, you can show that a solution exists.

For example, $\underline{v}'' = \underline{F}(t, \underline{v}, \underline{v}')$, $0 \leq t \leq b$, $\underline{v}(0) = \underline{\alpha}$, $\underline{v}(b) = \underline{\beta}$ *

If $\underline{F}(t, \underline{v}, \underline{v}')$ is continuous and Lipschitz in the variables $\underline{v}, \underline{v}'$ for a domain $0 \leq t \leq b$ and $\underline{v}, \underline{v}'$ bounded, then * has a unique solution for b sufficiently small (see Ascher, Mattheij, Russel).

For the general case, suppose a solution exists. Consider now the local uniqueness of the solution \Rightarrow whether the solution is isolated.



A solution is isolated if there are no other solutions in a "tube" about it.

Variational problem

Linearize the general boundary value problem about a

solution $y(t)$: $\underline{z}'(t) = \underline{A}(t, \underline{y}(t)) \underline{z}(t), 0 \leq t \leq b \quad \left. \begin{array}{l} \text{linear,} \\ \text{homogeneous} \end{array} \right\}$ problem

$$\underline{B}_0 \underline{z}(0) + \underline{B}_1 \underline{z}(b) = 0$$

where $\underline{A} = \frac{\partial \underline{F}}{\partial \underline{y}}(t, \underline{y}(t))$, Jacobian of \underline{F} evaluated at t and $\underline{y}(t)$

$$\underline{B}_0 = \frac{\partial \underline{g}}{\partial \underline{u}}(\underline{y}(0), \underline{y}(b)) \quad \text{where } \underline{g}(\underline{y}(0), \underline{y}(b)) = \underline{g}(\underline{u}, \underline{u}) = \underline{0}$$

$$\underline{B}_1 = \frac{\partial \underline{g}}{\partial \underline{v}}(\underline{y}(0), \underline{y}(b))$$

The solution $y(t)$ is isolated if the variational problem has only the trivial solution, $z=0$.

Linear BVPs

Consider the BVP, $\underline{y}' = \underline{A}(t) \underline{y} + \underline{g}(t), 0 \leq t \leq b \quad \left. \begin{array}{l} \text{non-homogeneous} \\ \text{linear} \end{array} \right\}$ problem

with boundary conditions: ~~$\underline{y}(0) = \underline{c}$~~ $\underline{B}_0 \underline{y}(0) + \underline{B}_1 \underline{y}(b) = \underline{0}$

Corresponding IVP:

$$\underline{y}' = \underline{A}(t) \underline{y} + \underline{g} \quad (\text{same ODE}), \quad 0 \leq t \leq b$$

$$\text{with initial condition } \underline{y}(0) = \underline{c}$$

IF $\underline{A}(t)$ is continuous on $[0, b]$, then there is a fundamental solution $\underline{\Phi}(t)$ matrix st $\underline{\Phi}'(t)$

$$\underline{\underline{x}}'(t) = \underline{\underline{A}}(t) \underline{\underline{x}}, \quad 0 \leq t \leq b$$

$$\underline{\underline{x}}(0) = \underline{\underline{I}}$$

The solution of the IVP is

$$\underline{\underline{y}}(t) = \underline{\underline{x}}(t) \left[\underline{\underline{c}} + \int_0^t \underline{\underline{x}}^{-1}(z) \underline{\underline{g}}(z) dz \right]$$

The solution of the linear BVP occurs for some choice of $\underline{\underline{c}}$.

$$BC \Rightarrow \underline{\underline{B}}_0 \underline{\underline{y}}(0) + \underline{\underline{B}}_1 \underline{\underline{y}}(b) = \underline{\underline{x}}$$

$$\Rightarrow \cancel{\underline{\underline{B}}_0 \underline{\underline{c}}} + \cancel{\underline{\underline{B}}_1 \underline{\underline{x}}}$$

$$\Rightarrow \underline{\underline{B}}_0 \underline{\underline{c}} + \underline{\underline{B}}_1 \underline{\underline{x}}(b) \left[\underline{\underline{c}} + \int_0^b \underline{\underline{x}}^{-1}(z) \underline{\underline{g}}(z) dz \right] = \underline{\underline{x}}$$

$$\rightarrow \underbrace{\underline{\underline{B}}_0 + \underline{\underline{B}}_1 \underline{\underline{x}}(b)}_{\underline{\underline{Q}}} \underline{\underline{c}} = \underline{\underline{x}} - \underline{\underline{B}}_1 \underline{\underline{x}} \int_0^b \underline{\underline{x}}^{-1}(z) \underline{\underline{g}}(z) dz$$

A unique solution of the linear boundary value problem exists if $\underline{\underline{Q}}$ is nonsingular.

Example

$$v'' + \lambda v = 0, \quad 0 \leq t \leq b$$

$$v(0) = \alpha, \quad v'(b) = \beta$$

Dirichlet condition

Neumann condition

$$\text{Rewrite as a system: } \underline{\underline{y}} = \begin{pmatrix} v \\ v' \end{pmatrix}$$

$$\Rightarrow \underline{\underline{y}}' = \begin{pmatrix} 0 & 1 \\ -\lambda & 0 \end{pmatrix} \underline{\underline{y}}, \quad 0 \leq t \leq b$$

$$\text{boundary conditions to be written as } \underline{\underline{B}}_0 \begin{pmatrix} v(0) \\ v'(0) \end{pmatrix} + \underline{\underline{B}}_1 \begin{pmatrix} v(b) \\ v'(b) \end{pmatrix} = ?$$

$$\rightarrow \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}}_{\underline{\underline{B}}_0} \begin{pmatrix} v(0) \\ v'(0) \end{pmatrix} + \underbrace{\begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}}_{\underline{\underline{B}}_1} \begin{pmatrix} v(b) \\ v'(b) \end{pmatrix} = \begin{pmatrix} \alpha \\ \beta \end{pmatrix} \quad \leftarrow \text{separated BCs}$$

want to analyze $\underline{\underline{Q}}$, so calculate fundamental solution $\underline{\underline{x}}$

The general soln of $v'' + \lambda v = 0$ is $v(t) = c_1 \cos \sqrt{\lambda} t + c_2 \sin \sqrt{\lambda} t$

$$\rightarrow v'(t) = -\sqrt{\lambda} c_1 \sin \sqrt{\lambda} t + \sqrt{\lambda} c_2 \cos \sqrt{\lambda} t$$

$$\rightarrow \underline{v}(t) = \begin{pmatrix} v \\ v' \end{pmatrix} = c_1 \begin{pmatrix} \cos \sqrt{\lambda} t \\ -\sqrt{\lambda} \sin \sqrt{\lambda} t \end{pmatrix} + c_2 \begin{pmatrix} \sin \sqrt{\lambda} t \\ \sqrt{\lambda} \cos \sqrt{\lambda} t \end{pmatrix}$$

$$\underline{\underline{x}}(t) = \begin{bmatrix} 1 & 1 \\ x_1 & x_2 \end{bmatrix}, \quad \underline{x}_1(0) = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad \underline{x}_2(0) = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

Find $\underline{\underline{x}}(t) = \begin{bmatrix} \cos \sqrt{\lambda} t & \frac{1}{\sqrt{\lambda}} \sin \sqrt{\lambda} t \\ -\sqrt{\lambda} \sin \sqrt{\lambda} t & \cos \sqrt{\lambda} t \end{bmatrix}$

$$\underline{\underline{Q}} = \underline{\underline{B}}_0 + \underline{\underline{B}}_1 \underline{\underline{x}}(b) = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \cos \sqrt{\lambda} b & \frac{1}{\sqrt{\lambda}} \sin \sqrt{\lambda} b \\ -\sqrt{\lambda} \sin \sqrt{\lambda} b & \cos \sqrt{\lambda} b \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 \\ -\sqrt{\lambda} \sin \sqrt{\lambda} b & \cos \sqrt{\lambda} b \end{bmatrix}$$

$$\det(\underline{\underline{Q}}) = \cos \sqrt{\lambda} b = 0 \Rightarrow \underline{\underline{Q}} \text{ is singular}$$

$\Rightarrow \sqrt{\lambda} b \neq \pi(k - \frac{1}{2}), \quad k = 1, 2, \dots$. Then unique soln exists

Green's Functions

Had the following solution for the BVP

$$(\star) \quad y(t) = \underline{\underline{X}}(t) \left[c + \int_0^t \underline{\underline{X}}^{-1}(z) \underline{q}(z) dz \right]$$

where c solves $\underline{\underline{Q}}c = \underline{y} - \underline{\underline{B}}_1 \underline{\underline{X}}(b) \int_0^b \underline{\underline{X}}^{-1}(z) \underline{q}(z) dz$
Object is to rearrange this solution to the form

$$y(t) = \underline{\underline{X}}(t) \underline{y} + \int_0^t \underline{\underline{G}}(t, z) \underline{q}(z) dz$$

↑ ↑
 scaled kernel of integral
 fundamental - Greens Function
 solution matrix

substitute c into (\star) :

$$y(t) = \underline{\underline{X}}(t) \underline{\underline{Q}}^{-1} \underline{y} - \underline{\underline{X}}(t) \underline{\underline{Q}}^{-1} \underline{\underline{B}}_1 \underline{\underline{X}}(b) \int_0^b \underline{\underline{X}}^{-1}(z) \underline{q}(z) dz$$

$$+ \underline{\underline{X}}(t) \int_0^t \underline{\underline{X}}^{-1}(z) \underline{q}(z) dz$$

Let $\underline{\underline{E}}(t) = \underline{\underline{X}}(t) \underline{\underline{Q}}^{-1}$ $\rightarrow \underline{\underline{X}}(t) = \underline{\underline{E}}(t) \underline{\underline{Q}}$ and $\underline{\underline{E}}^{-1}(t) = \underline{\underline{Q}} \underline{\underline{X}}(t)^{-1}$

$$y(t) = \underline{\underline{E}}(t) \underline{y} - \underline{\underline{E}}(t) \underline{\underline{B}}_1 \underline{\underline{X}}(b) \underline{\underline{Q}}^{-1} \underline{\underline{Q}} \int_0^b \underline{\underline{X}}^{-1}(z) \underline{q}(z) dz + \underline{\underline{E}}(t) \underline{\underline{Q}}^{-1} \underline{\underline{Q}} \int_0^t \underline{\underline{X}}^{-1}(z) \underline{q}(z) dz$$

$$\rightarrow y(t) = \underline{\underline{E}}(t) \underline{y} + \underline{\underline{E}}(t) \underline{\underline{B}}_1 \underline{\underline{X}}(b) \int_0^b \underline{\underline{E}}^{-1}(z) \underline{q}(z) dz + \underline{\underline{E}}(t) \int_0^t \underline{\underline{E}}^{-1}(z) \underline{q}(z) dz$$

Note: $\underline{\underline{Q}} = \underline{\underline{B}}_0 + \underline{\underline{B}}_1 \underline{\underline{X}}(b)$
 $= \underline{\underline{B}}_0 \underline{\underline{X}}(0) + \underline{\underline{B}}_1 \underline{\underline{X}}(b)$, because $\underline{\underline{X}}(0) = \underline{\underline{I}}$
 $= \underline{\underline{B}}_0 \underline{\underline{E}}(0) \underline{\underline{Q}} + \underline{\underline{B}}_1 \underline{\underline{E}}(b) \underline{\underline{Q}}$

$$\Rightarrow \underline{\underline{I}} = \underline{\underline{B}}_0 \underline{\underline{E}}(0) + \underline{\underline{B}}_1 \underline{\underline{E}}(b)$$

then $\underline{\underline{B}}_1 \underline{\underline{X}}(b) \int_0^b \underline{\underline{E}}^{-1}(z) \underline{q}(z) dz = \underline{\underline{B}}_1 \underline{\underline{E}}(b) \int_0^t \underline{\underline{E}}^{-1}(z) \underline{q}(z) dz + \underline{\underline{B}}_1 \underline{\underline{E}}(b) \int_t^b \underline{\underline{E}}^{-1}(z) \underline{q}(z) dz$
 $\approx (\underline{\underline{I}} - \underline{\underline{B}}_0 \underline{\underline{E}}(0)) \int_0^t \underline{\underline{E}}^{-1}(z) \underline{q}(z) dz + \underline{\underline{B}}_1 \underline{\underline{E}}(b) \int_t^b \underline{\underline{E}}^{-1}(z) \underline{q}(z) dz$

$$\Rightarrow \underline{y}(t) = \underline{\Phi}(t) \underline{\delta} - \underline{\Phi}(t) \left((\underline{I} - \underline{B}_0 \underline{\Phi}(0)) \int_0^t \underline{\Phi}^{-1}(z) \underline{q}(z) dz + \underline{B}_1 \underline{\Phi}(b) \int_t^b \underline{\Phi}^{-1}(z) \underline{q}(z) dz \right)$$

$$+ \underline{\Phi}(t) \int_0^t \underline{\Phi}^{-1}(z) \underline{q}(z) dz$$

$$\rightarrow \underline{y}(t) = \underline{\Phi}(t) \underline{\delta} + \underline{\Phi}(t) \underline{B}_0 \underline{\Phi}(0) \int_0^t \underline{\Phi}^{-1}(z) \underline{q}(z) dz - \underline{\Phi}(t) \underline{B}_1 \underline{\Phi}(b) \int_t^b \underline{\Phi}^{-1}(z) \underline{q}(z) dz$$

$$\Rightarrow \underline{y}(t) = \underline{\Phi}(t) \underline{\delta} + \int_0^b \underline{G}(t, z) \underline{q}(z) dz$$

where $\underline{G}(t, z) = \begin{cases} \underline{\Phi}(t) \underline{B}_0 \underline{\Phi}(0) \underline{\Phi}^{-1}(z), & 0 \leq z \leq t \\ -\underline{\Phi}(t) \underline{B}_1 \underline{\Phi}(b) \underline{\Phi}^{-1}(z), & t \leq z \leq b \end{cases}$

GREENS
FUNCTION

Problem Stability

Consider the linear boundary value problem

$$\underline{y}' = \underline{A}(t) \underline{y} + \underline{q}(t), \quad 0 \leq t \leq b, \quad \underline{B}_0 \underline{y}(0) + \underline{B}_1 \underline{y}(b) = \underline{\delta}$$

solution is

$$\underline{y}(t) = \underline{\Phi}(t) \underline{\delta} + \int_0^b \underline{G}(t, z) \underline{q}(z) dz$$

The behavior of the solution depends on $\underline{\Phi}(t)$ and $\underline{G}(t, z)$.

Consider a perturbed problem:

$$\tilde{y}' = \underline{A}(t) \tilde{y} + \underline{q}(t) + \underline{\delta}(t), \quad \underline{B}_0 \tilde{y}(0) + \underline{B}_1 \tilde{y}(b) = \underline{\delta} + \underline{\epsilon}$$

Consider the difference $\underline{w}(t) = \tilde{y}(t) - \underline{y}(t)$

$$\underline{w}(t) \text{ solves } \underline{w}' = \underline{A}(t) \underline{w} + \underline{\delta}(t), \quad \underline{B}_0 \underline{w}(0) + \underline{B}_1 \underline{w}(b) = \underline{\epsilon}$$

$$\Rightarrow \underline{w}(t) = \underline{\Phi}(t) \underline{\epsilon} + \int_0^t \underline{G}(t, z) \underline{\delta}(z) dz$$

$$\text{If } K = \max_{0 \leq t, z \leq b} \left(\|\underline{G}(t, z)\|_\infty, \|\underline{\Phi}(t)\|_\infty \right)$$

$$\text{then } \|\underline{w}(t)\| = \max_{0 \leq t \leq b} \|\underline{w}(t)\| \leq K \left(\|\underline{\epsilon}\|_\infty + \int_0^b \|\underline{\delta}(z)\|_\infty dz \right)$$

$\hookrightarrow K$ is an absolute condition number

Problem Stability for Linear BVPs

BVP: $\underline{y}' = \underline{\underline{A}}(t) \underline{y} + \underline{g}(t), \quad 0 \leq t \leq b$

$$\underline{B}_0 \underline{y}(0) + \underline{B}_1 \underline{y}(b) = \underline{x}$$

Recall: Fundamental solution matrix for the IVP

$$\underline{\underline{X}}' = \underline{\underline{A}}(t) \underline{\underline{X}}, \quad \underline{\underline{X}}(0) = \underline{\underline{I}}$$

Define $\underline{\underline{Q}} = \underline{\underline{B}}_0 + \underline{\underline{B}}_1 \underline{\underline{X}}(b)$, then a solution of the BVP exists if $\underline{\underline{Q}}$ is nonsingular.

Suppose $\underline{\underline{Q}}$ is nonsingular. Define $\underline{\underline{\Phi}}(t) = \underline{\underline{X}}(t) \underline{\underline{Q}}^{-1}$, then the solution of BVP is:

$$\underline{y}(t) = \underline{\underline{\Phi}}(t) \underline{x} + \int_0^b \underline{\underline{\Phi}}(t, \tau) \underline{g}(\tau) d\tau,$$

$$\underline{\underline{\Phi}}(t, \tau) = \begin{cases} \underline{\underline{\Phi}}(t) \underline{\underline{B}}_0 \underline{\underline{\Phi}}(0) \underline{\underline{\Phi}}^{-1}(\tau) & t > \tau \\ -\underline{\underline{\Phi}}(t) \underline{\underline{B}}_1 \underline{\underline{\Phi}}(b) \underline{\underline{\Phi}}^{-1}(\tau) & t < \tau \end{cases}$$

Problem stability depends on the behavior of $\underline{\underline{\Phi}}(t)$, $\underline{\underline{\Phi}}(t, \tau)$.

- Assume the boundary conditions are separated:

$$\Rightarrow \underline{\underline{B}}_0 = \underbrace{\left[\begin{array}{cccc|c} x & \dots & x & \dots & x & \vdots \\ x & - & x & - & x & \vdots \\ \vdots & & \ddots & & \ddots & \\ 0 & & & & & \end{array} \right] \}_{m \times m}^{\{K\}}, \quad \underline{\underline{B}}_1 = \left[\begin{array}{c} 0 \\ \hline x & \dots & x & \dots & x \\ x & \dots & x & \dots & x \end{array} \right]_{m \times m}^{\{m-K\}}$$

" K modes of the solution are forced from the left whereas $m-K$ modes are forced from the right"

IF $\underline{\underline{\Phi}}(t)$ has first K columns that decay as t increases, and the last $m-K$ columns that decay as t decreases, then the BVP is stable. This is called an (exponential) dichotomy.

return to boundary conditions

$$\underline{B}_0 \underline{\Phi}(0) + \underline{B}_1 \underline{\Phi}(t) = \underline{I} \quad \text{, because } \underline{\Phi}(t) = \underline{x}(t) Q^{-1}$$

if $\underline{B}_0 = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix}$ and $\underline{B}_1 = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix}$

then $\underline{B}_0 \underline{\Phi}(0) = \begin{bmatrix} \underline{I} & 0 & 0 \\ 0 & \underline{I} & 0 \\ 0 & 0 & -\underline{I} \end{bmatrix} = \underline{P}$

The matrix \underline{P} is an orthogonal projector:

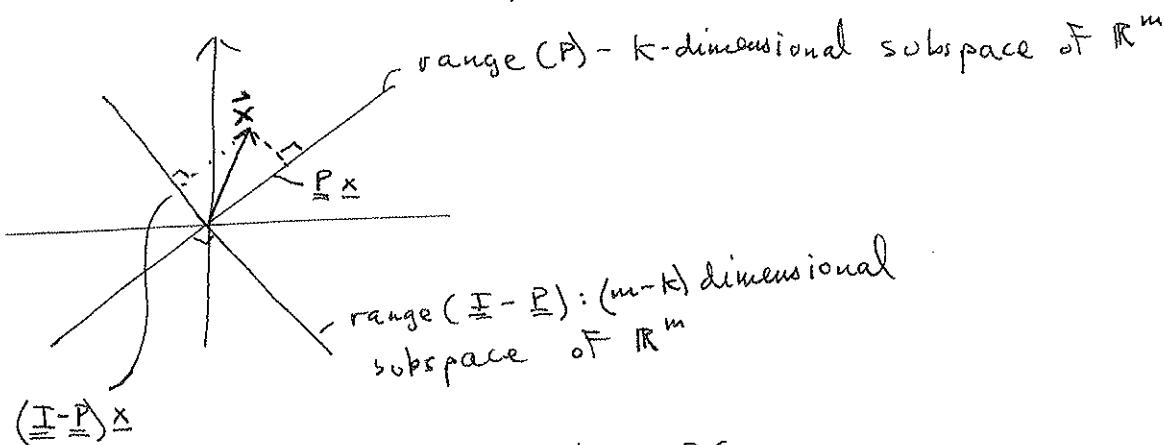
$$\Rightarrow \underline{P}^2 = \underline{P} \quad (\text{projection matrix})$$

and $\underline{B}_1 \underline{\Phi}(t) = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & \underline{I} \\ 0 & \underline{I} & 0 \end{bmatrix} = \underline{I} - \underline{P}$

$$\Rightarrow \underline{P}^T = \underline{P} \quad (\text{orthogonal})$$

The matrix \underline{P} projects \mathbb{R}^m into a k -dimensional subspace.

The matrix $\underline{I} - \underline{P}$ is also an orthogonal projector. It projects \mathbb{R}^m into a $(m-k)$ -dimensional subspace orthogonal to that of \underline{P} .



so for the case of separated BCs:

$$G(t, \tau) = \begin{cases} \underline{\Phi}(t) \underline{P} \underline{\Phi}^{-1}(\tau) & t > \tau \\ -\underline{\Phi}(t) (\underline{I} - \underline{P}) \underline{\Phi}^{-1}(\tau) & t < \tau \end{cases}$$

For dichotomy, we want \underline{P} to decay as t increases
and $\underline{I} - \underline{P}$ to decay as t decreases.

IF $\|\underline{\Phi}(t) \underline{P} \underline{\Phi}^{-1}(\tau)\| \leq K$ for $\tau < t$ and

if $\|\underline{\Phi}(t) (\underline{I} - \underline{P}) \underline{\Phi}^{-1}(\tau)\| < K$ for $\tau > t$ then

The boundary value problem possesses
a dichotomy solution (stable).

(BVP)

(Non ODE)

4/01/03

(54)

$$\text{IF } \|\underline{\underline{\underline{I}}}(t) \underline{\underline{\underline{P}}}^{-1}(x)\| \leq K e^{\alpha(x-t)}, \quad x < t$$

$$\text{and } \|\underline{\underline{\underline{I}}}(t)(\underline{\underline{\underline{I}}} - \underline{\underline{\underline{P}}}) \underline{\underline{\underline{P}}}^{-1}(x)\| \leq K e^{B(t-x)}, \quad t < x$$

for some $K > 0, \alpha, B > 0$ then the BVP
possesses an exponential dichotomy (asymptotic stability).

Example

$$v'' - v = 0, \quad 0 \leq t \leq b, \quad v(0) = \alpha, \quad v(b) = \beta$$

write as a first order system, $\underline{\underline{Y}}' = \underline{\underline{A}} \underline{\underline{Y}}$,

$$\text{where } \underline{\underline{Y}} = \begin{pmatrix} v \\ v' \end{pmatrix} \text{ and } \underline{\underline{A}} = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$$

Find fundamental solution matrix
general solution of original problem, $v(t) = c_1 \cosh t + c_2 \sinh t$

$$v'(t) = c_1 \sinh t + c_2 \cosh t$$

$$\Rightarrow \underline{\underline{Y}}(t) = c_1 \begin{pmatrix} \cosh t \\ \sinh t \end{pmatrix} + c_2 \begin{pmatrix} \sinh t \\ \cosh t \end{pmatrix}$$

pick c_1 and c_2 to solve for $\underline{\underline{X}}$: $\underline{\underline{X}} = \begin{pmatrix} \cosh t & \sinh t \\ \sinh t & \cosh t \end{pmatrix}$

satisfies $\underline{\underline{X}}(0) = \underline{\underline{I}}$

Note that $\underline{\underline{X}}(t)$ grows as t increases, which suggests the IVP is unstable, though it turns out the BVP is stable.

$$\text{From BCs: } v(0) = \alpha, \quad v(b) = \beta$$

$$\text{write as } \underline{\underline{B}}_0 \underline{\underline{Y}}(0) + \underline{\underline{B}}_1 \underline{\underline{Y}}(b) = \underline{\underline{Y}}$$

$$\underbrace{\begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}}_{\underline{\underline{B}}_0} \underbrace{\begin{pmatrix} v(0) \\ v'(0) \end{pmatrix}}_{\underline{\underline{Y}}(0)} + \underbrace{\begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}}_{\underline{\underline{B}}_1} \underbrace{\begin{pmatrix} v(b) \\ v'(b) \end{pmatrix}}_{\underline{\underline{Y}}(b)} = \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

$$\underline{\underline{Q}} = \underline{\underline{B}}_0 + \underline{\underline{B}}_1 + \underline{\underline{X}}(b) = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} + \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} \cosh b & \sinh b \\ \sinh b & \cosh b \end{pmatrix}$$

$$= \begin{pmatrix} 1 & 0 \\ \cosh b & \sinh b \end{pmatrix}$$

Notice $\det(\underline{\underline{Q}}) = \sinh(b) \Rightarrow \underline{\underline{Q}}$ is non-singular.

$$\underline{\underline{P}}(t) = \underline{\underline{X}}(t) \underline{\underline{Q}}^{-1} = \frac{1}{\sinh b} \begin{bmatrix} \sinh(b-t) & \sinh(t) \\ -\cosh(b-t) & \cosh(t) \end{bmatrix}$$

↑
this column
decays as
t increases

↑
this column decays
as t decreases from
b

Calculate $\underline{\underline{P}}(t) \underline{\underline{P}}^{-1}(z) = \underline{\underline{L}}$

$$\underline{\underline{P}} \underline{\underline{P}}^{-1}(z) = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \cosh z & -\sinh z \\ \cosh(b-z) & \sinh(b-z) \end{bmatrix} = \begin{bmatrix} \cosh z & -\sinh z \\ 0 & 0 \end{bmatrix}$$

$$\Rightarrow \underline{\underline{L}} = \frac{1}{\sinh b} \begin{bmatrix} \sinh(b-t) & \sinh(t) \\ \cosh(b-t) & \cosh(t) \end{bmatrix} \begin{bmatrix} \cosh z & -\sinh z \\ 0 & 0 \end{bmatrix}$$

$$\underline{\underline{L}} = \frac{1}{\sinh b} \begin{bmatrix} \cosh z \sinh(b-t) & -\sinh z \sinh(b-t) \\ -\cosh z \cosh(b-t) & \sinh z \cosh(b-t) \end{bmatrix}$$

Pick a norm: Frobenius norm: $\|\underline{\underline{L}}\|_F^2 = \text{trace}(\underline{\underline{L}} \underline{\underline{L}}^T)$

$$\Rightarrow \|\underline{\underline{L}}\|_F^2 = \frac{\cosh 2z \cdot \cosh 2(b-t)}{\sinh^2 b}$$

Notice $z \geq 0$, $b-t \geq 0$ then $\cosh 2t = \frac{1}{2}(e^{2t} + e^{-2t})$

$$\rightarrow \cosh 2xt = \frac{e^{2xt}}{2} \frac{e^{-2xt}}{2}$$

$$\rightarrow \cosh 2t = \frac{e^{2t}}{2} (1 + e^{-4t}),$$

$$\cosh 2(b-t) = \frac{1}{2} e^{2(b-t)} (1 + e^{-4(b-t)}),$$

$$\sinh 2b = \frac{1}{2} e^b (1 - e^{2b})$$

$$\Rightarrow \|\underline{\underline{L}}\|_F^2 = e^{2z+2(b-t)-2b} = e^{2(z-t)}.$$

$$\Rightarrow \|\underline{\underline{L}}\|_F = e^{z-t} \sqrt{\frac{C}{L}}$$

the
in theory

$d = L$
 $K = \frac{2}{(1-e^{-1})}$

You also $\|\underline{\underline{R}}\|_F \leq e^{t-z} K$

↑
where $K = \frac{2}{(1-e^{-1})}$

Stiff BVP

For an asymptotically stable BVP (ie one that possesses an exponential dichotomy) solution modes decay in from the boundaries depending on the number of BCs on each side.

If the eqn has constant coeffs then the rate of decay is determined by the eigenvalues of the coefficient matrix.

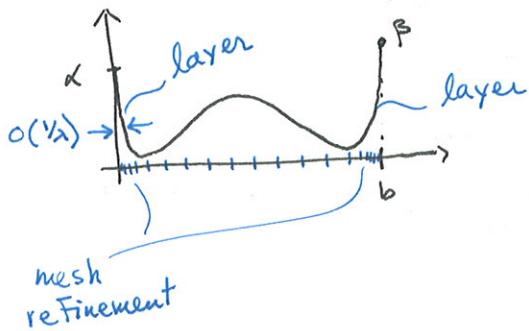
Suppose λ_j , $j = 1, \dots, m$ are these eigenvalues. The BVP would be stiff if the real part of λ_j times b is large for any j :

$$\boxed{b \cdot \operatorname{Re}(\lambda_j)} \gg 1 \text{ for any } j$$

Previously, for an IVP integrating for t increasing, the problem is stiff if $-(\operatorname{Re} \lambda_j)b \gg 1$. This differs from a BVP where both signs of $\operatorname{Re} \lambda_j$ must be considered.

For IVPs, the preferred methods had stiff decay (eg backward Euler, BDF, ...). There are no known methods for BVPs akin to this "stiff" decay because modes decay in both directions.

For BVPs, stiffness is handled numerically by resolving layers using some form of grid adaptation.



shooting methods

BVP: $\underline{y}' = \underline{f}(t, \underline{y})$, $0 \leq t \leq b$, $\underline{g}(\underline{y}(0), \underline{y}(b)) = 0$

construct solution by solving associated IVP
(the IVP has same diff eqn)

IVP: $\underline{y}' = \underline{f}(t, \underline{y})$
 $\underline{y}(0) = \underline{c}$ $\rightarrow \underline{y}(t; \underline{c})$

Need \underline{c} s.t. $\underline{g}(\underline{c}, \underline{y}(b; \underline{c})) = 0 \Rightarrow$ root-finding problem

use Newton's method

$$\underline{c}^{\text{NEW}} = \underline{c}^{\text{OLD}} - \frac{\underline{\Delta c}}{\text{correction}}$$

how to find correction?

Shooting Methods For BVPs

4/8/03

General 2-point BVP: $\underline{y}' = \underline{f}(t, \underline{y})$, $0 \leq t \leq b$

with boundary conditions: $\underline{g}(\underline{y}(0), \underline{y}(b)) = 0$

The idea behind shooting is to solve an associated initial value problem for some chosen initial condition and then iterate on the choice of the initial condition until BCs are satisfied.

Associated IVP: (same diff eq)

$$\underline{y}' = \underline{f}(t, \underline{y}), 0 \leq t \leq b, \underline{y}(0) = \underline{c}$$

The solution of the IVP is $\underline{y}(t, \underline{c})$. Want to

find \underline{c} s.t. $\underline{g}(\underline{c}, \underline{y}(b, \underline{c})) = 0$

⇒ Root-finding problem for \underline{c}

solve the rootfinding problem using

Newton:

$$\text{define } F(\underline{c}) \equiv \underline{g}(\underline{c}, \underline{y}(b, \underline{c})) = 0$$

We need Jacobian matrix of F

For $\underline{g} = \underline{g}(\underline{u}, \underline{v})$

shooting

Num DIFF

4/8/03

(56)

$$\frac{\partial F}{\partial \underline{c}} = \frac{\partial \underline{g}}{\partial \underline{u}} \frac{\partial \underline{u}}{\partial \underline{c}} + \frac{\partial \underline{g}}{\partial \underline{v}} \frac{\partial \underline{v}}{\partial \underline{c}}$$

$$\text{let } \underline{u} = \underline{c}, \underline{v} = \underline{y}(b, \underline{c})$$

$$= \frac{\partial \underline{g}}{\partial \underline{u}} \underline{I} + \frac{\partial \underline{g}}{\partial \underline{v}} \frac{\partial \underline{v}}{\partial \underline{c}}$$

need to compute this derivative,
which becomes a variational problem

to compute $\frac{\partial \underline{v}}{\partial \underline{c}}$ we have

$$\frac{\partial \underline{v}}{\partial t} = \underline{f}(t, \underline{y}(t, \underline{c}))$$

Differentiate wrt \underline{c}

$$\frac{\partial}{\partial \underline{c}} \left(\frac{\partial \underline{v}}{\partial t} \right) = \frac{\partial}{\partial \underline{c}} \left(\underline{f}(t, \underline{y}(t, \underline{c})) \right)$$

$$\rightarrow \frac{\partial}{\partial t} \left(\frac{\partial \underline{v}}{\partial \underline{c}} \right) = \frac{\partial \underline{f}}{\partial \underline{y}}(t, \underline{y}) \cdot \frac{\partial \underline{y}}{\partial \underline{c}}$$

$$\downarrow$$

$$\downarrow$$

$$\rightarrow \begin{cases} \underline{z}' = \underline{A}(t, \underline{y}) \underline{z} \\ \underline{A} = \frac{\partial \underline{f}}{\partial \underline{y}} \end{cases} \quad \text{with IC: } \underline{y}(0, \underline{c}) = \underline{c}$$

$$\downarrow$$

$$\frac{\partial \underline{y}}{\partial \underline{c}}(0, \underline{c}) = \underline{I}$$

$$\downarrow$$

$$\underline{z}(0) = \underline{I}$$

\Rightarrow Variational Problem

$$\underline{z}' = \underline{A} \underline{z}, \quad \underline{z}(0) = \underline{I}, \quad \underline{A} = \frac{\partial \underline{f}}{\partial \underline{y}},$$

Algorithm

- 1) Choose $\underline{c}^{(0)}$
- 2) for $k = 0, 1, \dots$
- 3) solve the IVPs

$$\begin{aligned} \dot{y} = f(t, y), \quad y(0) = \underline{c}^{(k)} \\ \dot{\underline{z}}_j' = \underline{A}(t, y) \underline{z}_j, \quad \underline{z}_j(0) = \underline{e}_j, \quad j = 1, \dots, m \end{aligned}$$

where $\underline{z} = \begin{bmatrix} \dot{z}_1 & \dot{z}_2 & \cdots & \dot{z}_m \end{bmatrix}^\top$
- 4) Compute $\underline{F}^{(k)} = g(\underline{c}^{(k)}, \underline{y}(b))$
- 5) compute

$$\underline{\Sigma}^{(k)} = \frac{\partial g}{\partial \underline{v}}(\underline{c}^{(k)}, \underline{y}(b)) + \frac{\partial g}{\partial \underline{v}}(\underline{c}^{(k)}, \underline{y}(b)) \underline{z}(b)$$
- 6) solve the $m \times m$ linear system

$$\underline{\Sigma}^{(k)} \underline{\Delta c}^{(k)} = \underline{F}^{(k)}$$
- 7) update $\underline{c}^{(k+1)} = \underline{c}^{(k)} - \underline{\Delta c}^{(k)}$
- 8) stop if $\|\underline{\Delta c}^{(k)}\| \leq tol$

Consider a specific case:

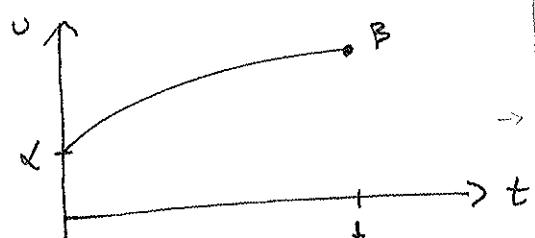
$$u'' - G(u) = 0, \quad 0 \leq t \leq 1, \quad u(0) = \alpha, \quad u(1) = \beta$$

Formulate and solve using shooting:

Associated IVP for a system of 1st order ODEs:

$$\underline{y} = \begin{pmatrix} u \\ u' \end{pmatrix} \rightarrow \dot{\underline{y}} = \begin{pmatrix} y_2 \\ G(y_1) \end{pmatrix}, \quad 0 \leq t \leq 1$$

pick initial condition: $\underline{y}(0) = \begin{pmatrix} \alpha \\ c \end{pmatrix}$



$$\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{pmatrix} u(0) \\ u'(0) \end{pmatrix} + \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix} \begin{pmatrix} u(1) \\ u'(1) \end{pmatrix} = \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

$$\Rightarrow g(\underline{y}(0), \underline{y}(b)) = \begin{pmatrix} u(0) - \alpha \\ u(1) - \beta \end{pmatrix}$$

\underline{A} is computed analytically

solve $m+1$ IVPs numerically

there's only one parameter, so it turns into a scalar problem:

$$F(c) = y_1(1, c) - \beta = 0$$

$$\text{Newton: } c^{(k+1)} = c^{(k)} - \frac{F(c^{(k)})}{F'(c^{(k)})}$$

Now calculate $F'(c)$

$$\text{we have } \underline{y}' = \begin{pmatrix} y_2 \\ G(y_1) \end{pmatrix}, \quad \underline{y}(0) = \begin{pmatrix} \alpha \\ c \end{pmatrix}$$

$$\underline{z} = \frac{\partial \underline{y}}{\partial c}, \quad \text{Differentiate wrt } c$$

$$\rightarrow \frac{\partial \underline{z}}{\partial t} = \begin{pmatrix} z_2 \\ G'(y_1) z_1 \end{pmatrix}, \quad \underline{z}(0) = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \quad \left. \right\} \text{ solve IVP numerically}$$

$$\Rightarrow c^{(k+1)} = c^{(k)} - \frac{y_1(1, c^{(k)}) - \beta}{z_1(1, c^{(k)})}$$

Example: $u'' - 9u = 0, \quad 0 \leq t \leq 1$

$$u(0) = 0, \quad u(1) = \sinh 3$$

this is a linear BVP whose solution is $u(t) = \sinh(3t)$

solve using shooting:

$$\text{set } \underline{w} = \begin{pmatrix} y_1 \\ y_2 \\ z_1 \\ z_2 \end{pmatrix} = \begin{pmatrix} u \\ u' \\ du/dc \\ d^2u/dc^2 \end{pmatrix}$$

$$\rightarrow \underline{w}' = \begin{pmatrix} w_2 \\ qw_1 \\ w_4 \\ qw_3 \end{pmatrix}, \quad \underline{w}(0) = \begin{pmatrix} \alpha \\ c \\ 0 \\ 1 \end{pmatrix}$$

$$c^{\text{new}} = c - \frac{w_4(1) - \sinh 3}{w_3(1)}$$

57
to calculate, integrate out to $t=1$ to obtain $y_1(1, c)$, then subtract Beta

Example

$$v''' - e^{t+1} = 0, \quad 0 \leq t \leq 1, \quad v(0) = v(1) = 0$$

"Bratu" problem

Here $v(t) = e^{w_1 t}$

$$\underline{w}' = \begin{pmatrix} w_2 \\ e^{w_1 + 1} \\ w_4 \\ w_3 e^{w_1 + 1} \end{pmatrix}, \quad \underline{w}(0) = \begin{pmatrix} 0 \\ c \\ 0 \\ 1 \end{pmatrix}$$

Difficulties

The main difficulty with shooting involves growing "modes" present in a stable BVP. If the modes grow too quickly then the linear system to determine the correction $\underline{\Delta c}$ become ill-conditioned.

BVP: $y' = f(t, y), \quad 0 \leq t \leq b, \quad g(y(0), y(b)) = 0$

For shooting, we solve some IVPs and a variational problem

$$\underline{y}' = \underline{f}(t, \underline{y}), \quad 0 \leq t \leq b, \quad \underline{y}(0) = \underline{c}$$

$$\underline{z}' = \underline{A} \underline{z}, \quad \underline{z}(0) = \underline{I}, \quad \underline{A} = \frac{\partial \underline{f}}{\partial \underline{y}}$$

Solve a linear system of the form $\underline{Q} \underline{\Delta c} = \underline{g}$

where $\underline{\Delta c}$ is the correction to \underline{c} and $\underline{Q} = \frac{\partial \underline{g}}{\partial \underline{y}} + \frac{\partial \underline{g}}{\partial \underline{z}} \cdot \underline{z}(b)$

If a solution mode (one of the columns of \underline{z}) grows too quickly as t increases, then \underline{Q} becomes ill-conditioned.

Example (pg 179 text)

$$v''' - 2\lambda v'' - \lambda^2 v' + 2\lambda^3 v = g(t), \quad 0 \leq t \leq 1$$

linear constant coeff ODE with nonhomogeneous term on right,
with BCs: $v(0) = \alpha, \quad v(1) = \beta, \quad v'(1) = \gamma$

First examine the homogeneous problem

substitute $v = e^{rt}$

$$\Rightarrow r^3 - 2\lambda r^2 - \lambda^2 r + 2\lambda^3 = 0 \quad (\text{char polynomial for homogeneous ODE})$$

$$\rightarrow (r+\lambda)(r-\lambda)(r-2\lambda) = 0$$

general solution: $v_h(t) = c_1 e^{-\lambda t} + c_2 e^{\lambda t} + c_3 e^{2\lambda t}$, homogeneous solution

$$\rightarrow v(t) = v_h(t) + v_p(t)$$

$$v(t) = c_1 e^{-\lambda t} + c_2 e^{\lambda t} + c_3 e^{2\lambda t} + v_p(t)$$

one mode that decays for increasing t
and two modes that grow as t increases.

\Rightarrow this suggests that IVP is unstable but the BVP would be nicely stable. If λ gets too big then shooting will become difficult.

The text chooses $g(t) = (\lambda^2 + \pi^2)(\pi \sin \pi t + 2\lambda \cos \pi t)$

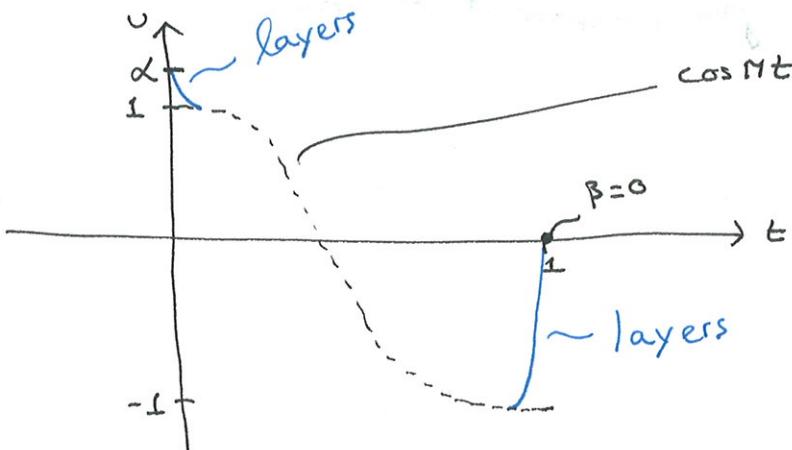
$$\alpha = \frac{3 + 2e^{-\lambda} + e^{-2\lambda}}{2 + e^{-\lambda}}$$

$$\beta = 0$$

$$\gamma = \lambda \left(\frac{3 - e^{-\lambda}}{2 + e^{-\lambda}} \right)$$

For any λ , with $g(t)$, β, α, γ as specified above, the exact solution is:

$$v(t) = \frac{1}{2 + e^{-\lambda}} \left(e^{-\lambda t} + e^{\lambda(t-1)} + e^{2\lambda(t-1)} \right) + \cos \pi t$$



the BCs force layers, depending on λ

Solve using shooting

IVP: $\underline{y}' = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -2\lambda^3 & \lambda^2 & 2\lambda \end{bmatrix} \underline{y} + \begin{pmatrix} 0 \\ 0 \\ g(t) \end{pmatrix}$, where $\underline{y} = \begin{pmatrix} u \\ u' \\ u'' \end{pmatrix}$

with $\underline{y}(0) = \begin{pmatrix} x \\ c_1 \\ c_2 \end{pmatrix}$ where c_1 and c_2 are shooting parameters
variational problems: $\underline{z}_1 = \frac{\partial \underline{y}_1}{\partial c_1}$, $\underline{z}_2 = \frac{\partial \underline{y}}{\partial c_2}$

$$\rightarrow \underline{z}'_j = \underline{A} \underline{z}_j, \quad j=1, 2$$

$$\rightarrow \underline{z}_j(0) = \begin{cases} \begin{pmatrix} 0 \\ 1 \end{pmatrix} & \text{if } j=1 \\ \begin{pmatrix} 0 \\ 1 \end{pmatrix} & \text{if } j=2 \end{cases}$$

so there are 3 ODEs
to solve, one ODE
for y , one for z_1 , and
for z_2 .

To solve, create "super-system"

$$\underline{w} = \begin{pmatrix} \underline{y} \\ \underline{z}_1 \\ \underline{z}_2 \end{pmatrix} \Rightarrow \underline{w}' = \begin{bmatrix} \underline{A} & & \\ & \underline{A} & \\ & & \underline{A} \end{bmatrix} \underline{w}, \quad \underline{w}(0) = \begin{pmatrix} x \\ c_1 \\ c_2 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

(don't have to create "super-system" cause ODEs
are decoupled just a matter of style)

eqns to satisfy are

$$y_1(1) - \beta = 0, \quad y_2(1) - \gamma = 0 \quad (\text{these depend on } c_1, c_2)$$

Linear system

$$\begin{bmatrix} z_{11}(1) & z_{21}(1) \\ z_{12}(1) & z_{22}(1) \end{bmatrix} \begin{pmatrix} \Delta c_1 \\ \Delta c_2 \end{pmatrix} = \begin{pmatrix} y_1(1) - \beta \\ y_2(1) - \gamma \end{pmatrix}$$

where $z_{ij}(1)$ is j^{th} component of \underline{z}_i at $t=1$

BVP: $\underline{y}' = \underline{f}(t, \underline{y})$, $0 \leq t \leq b$, $\underline{g}(\underline{y}(0), \underline{y}(b)) = 0$

(59)

associated IVP (same diff eq, same interval)

$$\underline{y}' = \underline{f}(t, \underline{y}), 0 \leq t \leq b, \underline{y}(0) = \underline{c}$$

variational problem $\underline{z}' = \underline{A}(t)\underline{z}$, $\underline{z}(0) = \underline{I}$, $\underline{A}(t) = \frac{\partial \underline{f}}{\partial \underline{y}}(t, \underline{y}(t))$

solve $\underline{Q} \Delta \underline{c} = \underline{g}$, where $\underline{g}(\underline{c}, \underline{y}(b))$ $\underline{z}(t) = \frac{\partial \underline{y}}{\partial \underline{c}}(t)$

where $\underline{Q} = \underline{g}_0 + \underline{g}_Y \underline{z}(b)$

Main Difficulty:

then $\underline{c}^{new} = \underline{c}^{old} - \Delta \underline{c}$ occurs if this linear system is not solved accurately

Example

BVP $u''' - 2\lambda u'' - \lambda^2 u' + 2\lambda^3 u = g(t), 0 \leq t \leq 1$

$$u(0) = \alpha, u(1) = \beta, u'(1) = \gamma$$

Pick $\alpha, \beta, \gamma, g(t)$ such that exact solution is

$$u(t) = \frac{e^{-\lambda t} + e^{\lambda(t-1)} + e^{2\lambda(t-1)}}{2 + e^{-\lambda}} + \cos \pi t$$

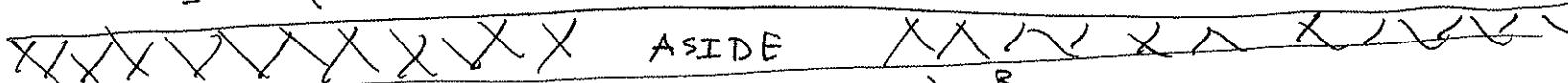
Notice, two growing modes in solution

associated IVP

$$\underline{y}' = \underline{A}\underline{y} + \underline{g}(t)$$

$$\underline{y}(0) = \begin{pmatrix} \alpha \\ c_1 \\ c_2 \end{pmatrix}$$

$$\text{where } \underline{y}(t) = \begin{pmatrix} u \\ u' \\ u'' \end{pmatrix}, \underline{A} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -2\lambda^3 & \lambda^2 & 2\lambda \end{bmatrix}, \underline{g}(t) = \begin{pmatrix} 0 \\ 0 \\ g(t) \end{pmatrix}$$



ASIDE

$$\text{Suppose } u'' - u = 0, u(0) = \alpha, u(1) = \beta$$

$$\text{then } \underline{y}' = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \underline{y}, \underline{g}(\underline{y}(0), \underline{y}(1)) = \begin{pmatrix} y_1(0) - \alpha \\ y_1(1) - \beta \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\frac{\partial \underline{g}}{\partial \underline{v}} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \frac{\partial \underline{g}}{\partial \underline{v}} = \begin{bmatrix} 0 & 0 \\ 1 & 0 \end{bmatrix}$$

Require c_1 and c_2 st

$$y_1(1) = B, \quad y_2(1) = \gamma$$

Iterate on (c_1, c_2) using Newton's method

The exact solution tells us there are two growing modes. In this case it is better to shoot in the reverse direction.

$$\underline{y}' = \underline{A} \underline{y} + \underline{q}(t), \quad 1 \geq t \geq 0$$

$$\underline{y}(1) = \begin{pmatrix} B \\ \gamma \end{pmatrix}$$

solve for some choice of c :

$$\text{let } \underline{z} \equiv \underline{z}(t) = \frac{\partial \underline{y}}{\partial c} \rightarrow \underline{z}' = \underline{A} \underline{z}, \quad \underline{z}(1) = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \quad 1 \geq t \geq 0$$

want c st $y_1(0) = \alpha$

scalar root finding problem: $F(c) = y_1(0, c) - \alpha = 0$

$$F'(c) = z_1(0, c)$$

$$\rightarrow \text{defn} \quad c^{\text{NEW}} = c^{\text{OLD}} - \frac{F(c^{\text{OLD}})}{F'(c^{\text{OLD}})}$$

Multiple Shooting

"Simple" shooting becomes ineffective if b is too big or

$\text{Re}(\text{eigenvalues of } \underline{I}_y)$ is too large. Rough estimate of problem is $b \|\underline{I}_y\|$ is too big. (Note, $b \|\underline{I}_y\|$ is dimensionless.) A fix is to use multiple shooting.

Consider a linear BVP:

$$\underline{y}' = \underline{A}(t) \underline{y} + \underline{q}(t), \quad 0 \leq t \leq b, \quad \underline{B}_0 \underline{y}(0) + \underline{B}_1 \underline{y}(b) = \underline{\gamma}$$

Introduce a partition of the interval:

$$0 = t_0 < t_1 < t_2 < \dots < t_N = b$$

let's assume the partition is known.

Idea: solve an IVP for each subinterval.

For $t_{n-1} \leq t \leq t_n$, define

$$\underline{Y}'_n(t) = \underline{\underline{A}}(t) \underline{Y}_n(t), \quad \underline{Y}_n(t_{n-1}) = \underline{\underline{I}} \quad \left. \right\} \text{homogeneous solution}$$

this is a fundamental solution matrix for IVP defined on subinterval with

$$\underline{v}'_n(t) = \underline{\underline{A}}(t) \underline{v}_n(t) + \underline{q}(t), \quad \underline{v}_n(t_{n-1}) = 0 \quad \left. \right\} \text{particular solution}$$

The general solution for the subinterval $t_{n-1} \leq t \leq t_n$ is

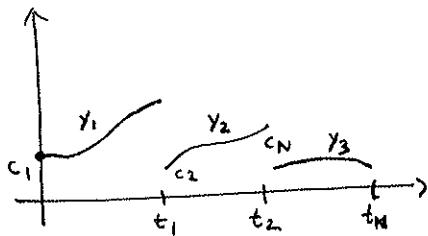
$$\underline{y}_n(t) = \underbrace{\underline{Y}_n(t)}_{\text{homogeneous}} \underline{\underline{c}}_n + \underbrace{\underline{v}_n(t)}_{\text{particular}}$$

Overall

$$\underline{y}(t) = \begin{cases} \underline{y}_1(t), & t_0 \leq t \leq t_1 \\ \underline{y}_2(t), & t_1 \leq t \leq t_2 \\ \vdots \\ \underline{y}_N(t), & t_{N-1} \leq t \leq t_N \end{cases}$$

this is a piecewise smooth solution, not necessarily continuous

For instance



choose $\underline{\underline{c}}$'s so that $\underline{y}(t)$ is continuous and satisfies the boundary conditions.

For continuity: $\underline{y}_n(t_n) = \underline{y}_{n+1}(t_n)$ for $n = 1, 2, \dots, N-1$ (interior nodes)

$$\begin{aligned} \rightarrow \underline{y}_n(t_n) &= \underline{Y}_n(t_n) \underline{\underline{c}}_n + \underline{v}_n(t_n) = \underline{\underline{Y}}_{n+1}(t_n) \\ &= \underline{y}_{n+1}(t_n) = \underline{\underline{c}}_{n+1} \end{aligned}$$

$$\Rightarrow \boxed{-\underline{Y}_n(t_n) \underline{\underline{c}}_n + \underline{\underline{c}}_{n+1} = \underline{v}_n(t_n)}$$

Boundary Conditions :

$$\underline{B}_0 \underline{Y}_*(0) + \underline{B}_1 \underline{Y}(b) = \underline{\delta}$$

\uparrow \uparrow

\underline{c}_1 $\underline{Y}_N(b) \underline{c}_N + \underline{Y}_N(b)$

$\underline{B}_1 \underline{Y}_N(b)$

→ $\underline{B}_0 \underline{c}_1 + \underline{B}_1 \underline{Y}_N(b) \underline{c}_N = \underline{\delta} - \underline{B}_1 \underline{Y}_N(b)$ \underline{w}

Unknowns: $N c$'s $\Rightarrow Nm$ unknowns

equations: $(N-1)m$ (From continuity)
 m (From BC's)

total: Nm equations, Nm unknowns

Collect the Nm linear eqns for c 's

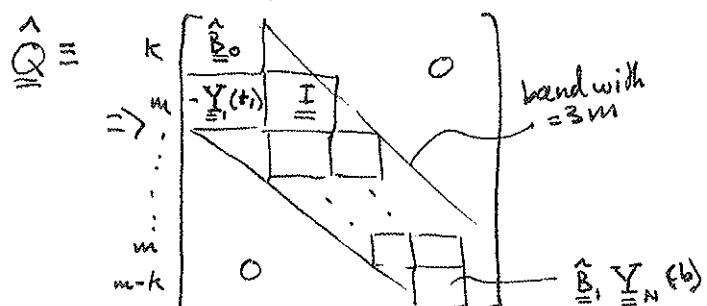
$$\underline{c} = \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_N \end{bmatrix} \in \mathbb{R}^{mN} \Rightarrow \begin{bmatrix} -\underline{Y}_1(t_1) & \underline{I} \\ -\underline{Y}_2(t_2) & \underline{I} \\ \ddots & \ddots \\ -\underline{Y}_{N-1}(t_{N-1}) & \underline{I} \\ \underline{B}_0 & \underline{B}_1 \underline{Y}_N(b) \end{bmatrix} \underline{c} = \begin{bmatrix} \underline{Y}_1(t_1) \\ \underline{Y}_2(t_2) \\ \vdots \\ \underline{Y}_{N-1}(t_{N-1}) \\ \underline{\delta} - \underline{B}_1 \underline{Y}_N(b) \end{bmatrix}$$

Solve the linear system for c , operation count = $O(m^3 N^3)$.
 The system is sparse so that the operation count can be reduced.

For example, if the BCs are separated then

$$\underline{B}_0 = \left[\begin{array}{c|c} \underline{\hat{B}}_0 & \\ \hline & 0 \end{array} \right] \stackrel{?}{=} \underbrace{\underline{\hat{B}}_0}_{m-k} \quad \underline{B}_1 = \left[\begin{array}{c|c} 0 & \\ \hline \underline{\hat{B}}_1 & \end{array} \right] \stackrel{?}{=} \underbrace{\underline{\hat{B}}_1}_{m-k}$$

Move the k eqns associated with k BCs at $t=0$ to the top of the linear system.



Now genuinely banded so solve using Gaussian elimination for a banded matrix with pivoting

operation count $\approx O(mN \cdot m^2)$ size of system times size band width squared
 (omit q due to order of mag $(bw)^2 = (3m)^2 = 9m^2$)

Remarks

- 1) For a given partition, each initial value problem is independent, so that these could be solved in parallel. \Rightarrow parallel shooting
- 2) May want to choose the nodes adaptively based on the growth of the solutions \Rightarrow IVP are no longer independent
- 3) The matrix $\hat{\underline{Q}}$ is better conditioned than \underline{Q} if the partition is fine enough.
- 4) IF the problem is nonlinear then a linearization leads to a system similar to $\hat{\underline{Q}}$.

Finite Difference Methods

4/14/03

BVP: $y' = f(t, y)$, $0 \leq t \leq b$, $g(y(0), y(b)) = 0$

Introduce a mesh or grid, ~~so that~~ $0 = t_0 < t_1 < \dots < t_n = b$

Set $h_n = t_n - t_{n-1} > 0$, $n=1, \dots, N$

- Choose some approximation

- Natural choice is a one-step method, symmetric, such as midpoint method or its higher order extensions (Gauss method); or trapezoidal method or its higher order extensions (Lobatto methods).

- There is no essential distinction between explicit versus implicit methods here because the plan is to solve for y of all grid points together.

Midpoint Method

Set $\underline{y}_n = \underline{y}(t_n)$, $n=0, \dots, N$

Let \underline{y}_n solve

$$\begin{cases} \underline{y}_n = \underline{y}_{n-1} + h_n \underline{f}(t_{n-1/2}, \frac{1}{2}(\underline{y}_n + \underline{y}_{n-1})) , n=1,2,\dots,N \\ \text{subject to } \underline{g}(\underline{y}_0, \underline{y}_N) = 0 \end{cases}$$

\Rightarrow set of $m(N+1)$ equations for $m(N+1)$ unknowns, $\{\underline{y}_n\}_{n=0}^N$

• Consider the linear case

$$\underline{f}(t, \underline{y}) = \underline{A}(t) \underline{y} + \underline{q}(t)$$

$$\underline{g}(\underline{u}, \underline{v}) = \underline{B}_0 \underline{u} + \underline{B}_1 \underline{v} - \underline{x}$$

$$\text{Have } \left\{ \begin{array}{l} \underline{y}_n = \underline{y}_{n-1} + h_n \left\{ \underline{A}(t_{n-1/2}) \frac{1}{2} (\underline{y}_{n-1} + \underline{y}_n) + \underline{q}(t_{n-1/2}) \right\} , n=1,2,\dots,N \\ \quad \downarrow \quad \downarrow \\ \quad \underline{A}_{n-1/2} \quad \underline{q}_{n-1/2} \\ \underline{B}_0 \underline{y}_0 + \underline{B}_1 \underline{y}_N = \underline{x} \end{array} \right.$$

Collect as a linear system, rewrite as:

$$\left(\underline{I} - \frac{h_n}{2} \underline{A}_{n-1/2} \right) \underline{y}_n - \left(\underline{I} + \frac{h_n}{2} \underline{A}_{n-1/2} \right) \underline{y}_{n-1} = h_n \underline{q}_{n-1/2}$$

$$\text{Set } \underline{D}_n = -\underline{I} - \frac{h_n}{2} \underline{A}_{n-1/2}, \quad \underline{U}_n = \underline{I} - \frac{h_n}{2} \underline{A}_{n-1/2}$$

$$\rightarrow \underline{D}_n \underline{y}_n + \underline{U}_n \underline{y}_{n-1} = h_n \underline{q}_{n-1/2}$$

$$\left[\begin{array}{cc} \underline{D}_1 & \underline{U}_1 \\ \underline{D}_2 & \underline{U}_2 \\ \vdots & \vdots \\ \underline{D}_N & \underline{U}_N \end{array} \right] \left[\begin{array}{c} \underline{y}_0 \\ \underline{y}_1 \\ \vdots \\ \vdots \\ \underline{y}_{N-1} \\ \underline{y}_N \end{array} \right] = \left[\begin{array}{c} h_1 \underline{q}_{1/2} \\ h_2 \underline{q}_{3/2} \\ \vdots \\ \vdots \\ h_N \underline{q}_{N-1/2} \\ \underline{x} \end{array} \right]$$

IF the boundary conditions are separated then

$$\underline{\underline{B}}_0 = \left[\begin{array}{c} \hat{\underline{\underline{B}}}_0 \\ -\hat{\underline{\underline{B}}}_0 \\ \vdots \\ 0 \end{array} \right]_{m-k}^k, \quad \underline{\underline{B}}_1 = \left[\begin{array}{c} 0 \\ \vdots \\ \hat{\underline{\underline{B}}}_1 \\ \vdots \\ 0 \end{array} \right]_{m-k}^k$$

perform manipulation to obtain banded matrix.

Now have

$$\left[\begin{array}{cccccc} \hat{\underline{\underline{B}}}_0 & & & & & \\ \hat{\underline{\underline{D}}}_1 & \hat{\underline{\underline{U}}}_1 & & & & \\ \vdots & \vdots & \ddots & & & \\ \hat{\underline{\underline{D}}}_2 & \hat{\underline{\underline{U}}}_2 & \ddots & & & \\ \vdots & \vdots & & \ddots & & \\ \text{BANDED} & & & & \ddots & \\ \text{MATRIX} & & & & & \end{array} \right] \cdot \left[\begin{array}{c} \underline{y}_0 \\ \underline{y}_1 \\ \underline{y}_2 \\ \vdots \\ \vdots \\ \underline{y}_N \end{array} \right] = \left[\begin{array}{c} \underline{x}_0 \\ h_1 g^{1/2} \\ h_2 g^{1/2} \\ \vdots \\ h_N g^{N-1/2} \\ \underline{x}_1 \end{array} \right]$$

(same structure as that for multishoot method, pg 60, back)

- Using direct banded solver, the operation count is $O(\text{size of system} \times \text{bandwidth}^2) \Rightarrow O(Nm^3)$
- tends to be a larger system than multishoot but don't have to integrate to get individual blocks

Example

$$u'' - 9u = 0, \quad 0 \leq t \leq 1, \quad u(0) = 0, \quad u(1) = \sinh 3, \quad u_{\text{exact}}(t) = \sinh(3t)$$

$$\text{let } \underline{y} = \begin{pmatrix} u \\ u' \end{pmatrix} \Rightarrow \underline{y}' = \underline{A} \underline{y}, \quad \underline{A} = \begin{bmatrix} 0 & 1 \\ 9 & 0 \end{bmatrix}$$

$$\hat{\underline{\underline{B}}}_0 = \hat{\underline{\underline{B}}}_1 = [1 \ 0], \quad \underline{y}_0 = 0, \quad \underline{y}_1 = \sinh(3)$$

→ see Matlab code, my linear BVP.m

$$\text{if have } \underline{\underline{B}}_0 = \begin{bmatrix} \hat{\underline{\underline{B}}}_0 & \\ -\hat{\underline{\underline{B}}}_0 & \cdot \\ \vdots & \cdot \\ 0 & \cdot \end{bmatrix} \text{ and } \underline{\underline{B}}_1 = \begin{bmatrix} 0 & \cdot \\ \vdots & \cdot \\ \hat{\underline{\underline{B}}}_1 & \cdot \\ \vdots & \cdot \\ 0 & \cdot \end{bmatrix}$$

$$\text{then define } \underline{\underline{B}} = \begin{bmatrix} \hat{\underline{\underline{B}}}_0 & \cdot \\ -\hat{\underline{\underline{B}}}_1 & \cdot \\ \vdots & \cdot \\ \hat{\underline{\underline{B}}}_1 & \cdot \\ \vdots & \cdot \\ 0 & \cdot \end{bmatrix} \text{ for storage}$$

use reshape function in Matlab

Example - stiffness

$$u''' - 2\lambda u'' - \lambda^2 u' + 2\lambda^3 u = g(t), \quad 0 \leq t \leq 1$$

$$u(0) = \alpha, \quad u(1) = \beta, \quad u'(1) = \gamma$$

$$\text{with } g(t) = (\lambda^2 + \pi^2)(\pi \sin \pi t + 2\lambda \cos \pi t)$$

$$\alpha = \frac{3 + 2e^{-\lambda} + e^{-2\lambda}}{2 + e^{-\lambda}}, \quad \beta = 0, \quad \gamma = \frac{\lambda(3 - e^{-\lambda})}{2 + e^{-\lambda}}$$

$$\text{The exact solution is } u_{\text{exact}}(t) = \frac{-\lambda t + e^{\lambda(t-1)} + e^{2\lambda(t-1)}}{2 + e^{-\lambda}} + \cos(\pi t)$$

The system becomes, $\underline{y}'(t) = \underline{A}\underline{y} + \underline{g}(t)$

$$\underline{y} = \begin{pmatrix} u \\ u' \\ u'' \end{pmatrix}, \quad \underline{A} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ -2\lambda^3 & \lambda^2 & 2\lambda \end{bmatrix}, \quad \underline{g}(t) = \begin{bmatrix} 0 \\ 0 \\ g(t) \end{bmatrix}, \quad \underline{B} = \begin{bmatrix} \hat{\underline{B}}_0 \\ \vdots \\ \hat{\underline{B}}_n \\ \vdots \\ \hat{\underline{B}}_1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ \cdots & \cdots & \cdots \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$

$$\underline{y} = \begin{bmatrix} \alpha \\ -\frac{\lambda}{\underline{B}} \\ \beta \end{bmatrix} \quad \text{Recall we encountered problems shooting in either direction}$$

can play the "mesh game"

Non linear BVPs

consider non linear BVP, but with linear separated BCs

$$\underline{y}' = \underline{f}(t, \underline{y}), \quad a \leq t \leq b$$

$$\begin{bmatrix} \hat{\underline{B}}_0 \\ \vdots \\ \hat{\underline{B}}_n \\ \vdots \\ \hat{\underline{B}}_1 \end{bmatrix} \underline{y}(a) + \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix} \underline{y}(b) = \begin{bmatrix} \underline{y}_0 \\ \vdots \\ \underline{y}_n \\ \vdots \\ \underline{y}_1 \end{bmatrix}$$

with mesh: $a = t_0 < t_1 < \dots < t_n = b$

using midpoint method

$$\underline{y}_n = \underline{y}_{n-1} + h_n \underline{f}\left(t_{n-1/2}, \frac{1}{2}(\underline{y}_n + \underline{y}_{n-1})\right), \quad n = 1, \dots, N$$

$$\hat{\underline{B}}_0 \underline{y}_0 = \underline{y}_0, \quad \hat{\underline{B}}_1 \underline{y}_n = \underline{y}_1$$

$$\text{Set } \underline{y}_{n-1} = \tilde{\underline{y}}_{n-1} + \underline{\Delta y}_{n-1}$$

$$\underline{y}_n = \tilde{\underline{y}}_n + \underline{\Delta y}_n$$

where $\tilde{\underline{y}}_j$ are approximations to discrete solution, $\underline{\Delta y}_j$ are corrections
substitute into equations

$$\Rightarrow \tilde{\underline{y}}_n + \underline{\Delta y}_n = \tilde{\underline{y}}_{n-1} + \underline{\Delta y}_{n-1} + h_n \tilde{F}\left(t_{n-1/2}, \frac{1}{2}(\tilde{\underline{y}}_{n-1} + \tilde{\underline{y}}_n) + \frac{1}{2}(\underline{\Delta y}_{n-1} + \underline{\Delta y}_n)\right)$$

Expand for small $\underline{\Delta y}$'s:

$$\begin{aligned} \tilde{\underline{y}}_n + \underline{\Delta y}_n &= \tilde{\underline{y}}_{n-1} + \underline{\Delta y}_{n-1} + h_n \tilde{F}\left(t_{n-1/2}, \frac{1}{2}(\tilde{\underline{y}}_{n-1} + \tilde{\underline{y}}_n)\right) \\ &\quad + h_n \tilde{F}'\left(t_{n-1/2}, \frac{1}{2}(\tilde{\underline{y}}_{n-1} + \tilde{\underline{y}}_n)\right) \frac{\underline{\Delta y}_{n-1} + \underline{\Delta y}_n}{2} + \dots \end{aligned}$$

Notation:

$$\text{let } \tilde{F}(t_{n-1/2}, \frac{1}{2}(\tilde{\underline{y}}_{n-1} + \tilde{\underline{y}}_n)) = \tilde{F}_{n-1/2}$$

$$\text{and } \tilde{F}'(t_{n-1/2}, \frac{1}{2}(\tilde{\underline{y}}_{n-1} + \tilde{\underline{y}}_n)) = \tilde{A}_{n-1/2}$$

$$\Rightarrow \tilde{\underline{y}}_n + \underline{\Delta y}_n = \tilde{\underline{y}}_{n-1} + \underline{\Delta y}_{n-1} + h_n \tilde{F}_{n-1/2} + \frac{h_n}{2} \tilde{A}_{n-1/2} (\underline{\Delta y}_{n-1} + \underline{\Delta y}_n) + \dots$$

$$\Rightarrow \left(\tilde{F} - \frac{h_n}{2} \tilde{A}_{n-1/2}\right) \underline{\Delta y}_n - \left(\tilde{F} + \frac{h_n}{2} \tilde{A}_{n-1/2}\right) \underline{\Delta y}_{n-1} = +\tilde{\underline{y}}_{n-1} - \tilde{\underline{y}}_n + h_n \tilde{F}_{n-1/2}$$

$$n = 1, \dots, N$$

$$\text{BCs: } \hat{\underline{B}}_0 (\tilde{\underline{y}}_0 + \underline{\Delta y}_0) = \underline{y}_0 \rightarrow \hat{\underline{B}}_0 \underline{\Delta y}_0 = \underline{y}_0 - \hat{\underline{B}}_0 \tilde{\underline{y}}_0$$

$$\text{and similarly } \dots \hat{\underline{B}}_N \underline{\Delta y}_N = \underline{y}_N - \hat{\underline{B}}_N \tilde{\underline{y}}_N$$

$$\left[\begin{array}{c} \text{as before} \\ \vdots \\ \Delta y_0 \\ \vdots \\ \Delta y_N \end{array} \right] = \text{RHS} \rightarrow \text{drive to zero}$$

Finite Difference Methods

for Nonlinear BVPs

BVP: $\underline{y}' = \underline{f}(t, \underline{y})$, $0 \leq t \leq b$

with separated linear BCs: $\underline{\underline{B}}_0 \underline{y}(0) = \underline{y}_0$, $\underline{\underline{B}}_1 \underline{y}(b) = \underline{y}_1$

(k BCs on left & $m-k$ BCs on right)

Notice - $\underline{\underline{B}}_0, \underline{\underline{B}}_1$, not square matrices

Consider the midpoint method:

mesh: $0 = t_0 < t_1 < \dots < t_N = b$, $h_n = t_n - t_{n-1} > 0$

Let $\underline{y}_n \approx \underline{y}(t_n)$ and set $\underline{y}_n = \underline{y}_{n-1} + h_n \underline{f}(t_{n-1/2}, \frac{1}{2}(\underline{y}_{n-1} + \underline{y}_n))$, $n=1, \dots, N$

$$\cancel{\underline{\underline{B}}_0 \underline{y}(0)} \approx \underline{y}_0, \cancel{\underline{\underline{B}}_1 \underline{y}(b)} \quad \underbrace{\underline{\underline{B}}_0 \underline{y}_0 = \underline{y}_0}_{k \text{ eqns}}, \underbrace{\underline{\underline{B}}_1 \underline{y}_N = \underline{y}_1}_{m-k \text{ eqns}}$$

number of eqns & unknowns - $m(N+1)$

- This is generally a set of nonlinear algebraic eqns, generally so employ iterative scheme, like Newton's method.

Set $\underline{y}_n = \tilde{\underline{y}}_n + \Delta \underline{y}_n$, $n=0, \dots, N$

substitute into diff eqns, linearize for small ~~small~~ $\Delta \underline{y}_n$

$$\Rightarrow (\underline{\underline{I}} + \frac{h_n}{2} \underline{\underline{A}}_n) \underline{\Delta y}_{n-1} + (-\underline{\underline{I}} + \frac{h_n}{2} \underline{\underline{A}}_n) \underline{\Delta y}_n = \tilde{\underline{y}}_n - \tilde{\underline{y}}_{n-1} - h_n \underline{f}_n$$

where $\underline{f}_n = \underline{f}(t_{n-1/2}, \frac{1}{2}(\tilde{\underline{y}}_{n-1} + \tilde{\underline{y}}_n))$

$$\underline{\underline{A}}_n = \frac{\partial \underline{f}}{\partial \underline{y}}(t_{n-1/2}, \frac{1}{2}(\tilde{\underline{y}}_{n-1} + \tilde{\underline{y}}_n)), m \times m \text{ matrix}$$

$$\text{and } \underline{\underline{B}}_0 \underline{\Delta y}_0 = \underline{y}_0 - \underline{\underline{B}}_0 \tilde{\underline{y}}_0, \quad \underline{\underline{B}}_1 \underline{\Delta y}_N = \underline{y}_1 - \underline{\underline{B}}_1 \tilde{\underline{y}}_N$$

Linear system for Δy 's:

$$\left[\begin{array}{ccc|c} \hat{\underline{B}}_0 & 0 & 0 & \Delta y_0 \\ \frac{I + h_1 A_1}{2} & -\frac{I + h_1 A_1}{2} & 0 & \Delta y_1 \\ 0 & \frac{I + h_2 A_2}{2} & -\frac{I + h_2 A_2}{2} & \vdots \\ & & & \Delta y_{N-1} \\ & & & 0 & \Delta y_N \end{array} \right] = \left[\begin{array}{c} \underline{y}_0 - \hat{\underline{B}}_0 \tilde{\underline{y}}_0 \\ \tilde{\underline{y}}_1 - \tilde{\underline{y}}_0 - h_1 F_1 \\ \tilde{\underline{y}}_2 - \tilde{\underline{y}}_1 - h_2 F_2 \\ \vdots \\ \tilde{\underline{y}}_N - \tilde{\underline{y}}_{N-1} - h_N F_N \\ \underline{y}_1 - \hat{\underline{B}}_0 \tilde{\underline{y}}_N \end{array} \right]$$

For a given $\tilde{\underline{y}}_n$'s, you solve the linear system for Δy_n 's
then update:

$$\underline{y}_n^{\text{new}} = \tilde{\underline{y}}_n + \Delta y_n$$

then repeat with $\tilde{\underline{y}}_n = \underline{y}_n^{\text{new}}$ until $\max_n \|\Delta y_n\| < \text{tol}$

Example:

$$v'' - G(v) = 0, \quad 0 \leq t \leq 1$$

$$v(0) = \alpha, \quad v(1) = \beta$$

$$\text{set } \underline{y} = \begin{pmatrix} v \\ v' \end{pmatrix} \rightarrow \underline{y}' = \begin{pmatrix} y_2 \\ G(y_1) \end{pmatrix}, \quad \hat{\underline{B}}_0 = \hat{\underline{B}}_1 = \begin{bmatrix} 1 & 0 \end{bmatrix}$$

$$y_0 = \alpha, \quad y_1 = \beta$$

Midpoint

$$y_n = y_{n-1} + h_n \begin{pmatrix} \frac{1}{2}(y_{n-1} + y_n) \\ G\left(\frac{1}{2}(y_{n-1} + y_n)\right) \end{pmatrix}, \quad y_{10} = \alpha, \quad y_{1N} = \beta$$

$$\text{so } F_n = \begin{pmatrix} \frac{1}{2}(\tilde{y}_{n-1} + \tilde{y}_n) \\ G\left(\frac{1}{2}(\tilde{y}_{n-1} + \tilde{y}_n)\right) \end{pmatrix}$$

$$\text{remember, } \underline{F} = \begin{pmatrix} y_2 \\ G(y_1) \end{pmatrix} \quad \underline{F}' = \begin{pmatrix} 0 & 1 \\ G'(y_1) & 0 \end{pmatrix} \rightarrow \underline{A}_n = \begin{bmatrix} 0 & 1 \\ G'\left(\frac{1}{2}(\tilde{y}_{n-1} + \tilde{y}_n)\right) & 0 \end{bmatrix}$$

Case ① : $G(0) = 90$, $\lambda = 0$, $B = \sinh 3$

exact soln: $v(t) = \sinh(3t)$

Case ② $G(0) = -e^{0+2}$, $\lambda = B = 0$, "Bratu" problem

Aside

given $\tilde{y}_n = \tilde{y}_{1,n}$ and need $\tilde{y}_{2,n}$, where $y_2 = v'$

then use $\tilde{y}_{2,n} = \frac{\tilde{y}_{1,n} - \tilde{y}_{1,n-1}}{h_n}$

Consistency, O-stability, and convergence
For finite difference methods.

For midpoint, for example, have

$$N_h y_n = \frac{y_n - y_{n-1}}{h_n} - F\left(t, \frac{1}{2}(y_n + y_{n-1})\right), \quad (Ny = y' - F)$$

$$n = 1, \dots, N$$

$N_h y_n = 0$, if y_n solution of diff. eqn.

Local truncation error: $d_n = N_h y(t_n)$, $n = 1, \dots, N$

For midpoint, can show^(using Taylor series) that $d_n = O(h^2)$ for $h = \max_n h_n$

since $d_n \rightarrow 0$ as $h \rightarrow 0 \Rightarrow$ midpoint method is consistent

would like to show that the finite difference approximation is convergent.

$$\text{global error } e_n = y_n - y(t_n)$$

$$\Rightarrow \max_n \|e_n\| \rightarrow 0 \text{ as } h \rightarrow 0$$

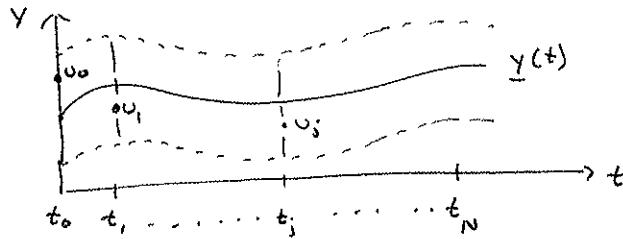
Specifically, $\max_n \|e_n\| = O(h^p)$, where p is order of accuracy
(for midpoint, $p = 2$)

The step from consistency to convergence requires
O-stability.

O-stability

define a "discrete tube" about a chosen trajectory $\underline{y}(t)$:

$$S_p(\underline{y}) = \left\{ \underline{u}_n \text{ st } \|\underline{u}_n - \underline{y}(t_n)\| < p \quad \forall n \right\}$$



The positive value p defines the radius of the tube.
A difference method is O-stable if p , h_0 and K (all positive) exist such that for any mesh with $h \leq h_0$ and any mesh functions \underline{x}_n and \underline{z}_n in $S_p(\underline{y})$, we have

$$(*) \quad \|\underline{x}_n - \underline{z}_n\| \leq K \left\{ \|\underline{g}(\underline{x}_0, \underline{x}_N) - \underline{g}(\underline{z}_0, \underline{z}_N)\| + \max_{1 \leq i \leq N} \|N_h \underline{x}_i - N_h \underline{z}_i\| \right\}$$

For all n .

IF $\underline{x}_n = \underline{y}_n$, $\underline{z}_n = \underline{y}(t_n)$ then $(*)$ tells us

$$\|\underline{e}_n\| \leq K \max_{1 \leq i \leq N} \|d_i\| \quad \forall n = 1, \dots, N$$

\uparrow
 $O(h^2)$

$$\Rightarrow \|\underline{e}_n\| = O(h^2)$$

(66)

O-stability for BVPsBVP $\underline{y}' = \underline{F}(t, \underline{y})$, $0 \leq t \leq b$

$$\underline{g}(\underline{y}(0), \underline{y}(b)) = 0$$

Finite difference method:

$$N_h \underline{y}_n = 0, n = 1, \dots, N$$

$$\underline{g}(\underline{y}_0, \underline{y}_N) = 0$$

$$\text{For Midpoint: } N_h \underline{y}_n = \frac{\underline{y}_n - \underline{y}_{n-1}}{h_n} - \underline{F}\left(t_{n-1/2}, \frac{\underline{y}_{n-1} + \underline{y}_n}{2}\right)$$

~~$N_h \underline{y}_n = 0$~~

Define: a discrete tube about a solution $\underline{y}(t)$:

$$S_p(\underline{y}) = \left\{ \underline{u}_n \text{ st } \|\underline{u}_n - \underline{y}(t_n)\| < p^{1/n} \right\}$$

A difference method is O-stable if p, k, h_0 exist
 st for any mesh with $h \leq h_0$ and any mesh
 functions $\underline{x}_n, \underline{z}_n$ in $S_p(\underline{y})$, we have

$$\|\underline{x}_n - \underline{z}_n\| \leq k \left\{ \|\underline{g}(\underline{x}_0, \underline{x}_N) - \underline{g}(\underline{z}_0, \underline{z}_N)\| + \max_{1 \leq j \leq N} \|N_h \underline{x}_j - N_h \underline{z}_j\| \right\} + n$$

In order to establish convergence we need to show consistency
 (easy) and O-stability (harder).

IF the BVP is non linear, then you would linearize about a
 solution $\underline{y}(t)$ and try to show O-stability for the linearized
 problem. This result would hold for p sufficiently small.

Consider the linear case: $\underline{F}(t, \underline{y}) = \underline{A}(t) \underline{y} + \underline{g}(t)$

$$\underline{g}(\underline{u}, \underline{v}) = \underline{B}_0 \underline{u} + \underline{B}_1 \underline{v} - \underline{v}$$

calculate:

$$g(\underline{x}_0, \underline{x}_N) - g(\underline{z}_0, \underline{z}_N) = \left(\underline{\underline{B}}_0 \underline{x}_0 + \underline{\underline{B}}_1 \underline{x}_N - \underline{\underline{d}} \right) - \left(\underline{\underline{B}}_0 \underline{z}_0 + \underline{\underline{B}}_1 \underline{z}_N - \underline{\underline{r}} \right)$$

$$\boxed{g(\underline{x}_0, \underline{x}_N) - g(\underline{z}_0, \underline{z}_N) = \underline{\underline{B}}_0 \underline{w}_0 + \underline{\underline{B}}_1 \underline{w}_N}$$

where $\underline{w}_0 = \underline{x}_0 - \underline{z}_0$, $\underline{w}_N = \underline{x}_N - \underline{z}_N$

$$N_h \underline{x}_j = \frac{\underline{x}_j - \underline{x}_{j-1}}{h_j} - \underline{\underline{A}}(t_{j-1/2}) \frac{1}{2} (\underline{x}_j + \underline{x}_{j-1}) \stackrel{\ominus}{\rightarrow} q(\underline{t}_{j-1/2})$$

$$N_h \underline{z}_j = \left(\frac{1}{h_j} \underline{\underline{I}} - \frac{1}{2} \underline{\underline{A}}_{j-1/2} \right) \underline{z}_j - \left(\frac{1}{h_j} \underline{\underline{I}} + \frac{1}{2} \underline{\underline{A}}_{j-1/2} \right) \underline{z}_{j-1} \stackrel{\ominus}{\rightarrow} q(\underline{t}_{j-1/2})$$

where $\underline{\underline{A}}_{j-1/2} = \underline{\underline{A}}(t_{j-1/2})$

and similarly,

$$N_h \underline{z}_j = \left(\frac{1}{h_j} \underline{\underline{I}} - \frac{1}{2} \underline{\underline{A}}_{j-1/2} \right) \underline{z}_j - \left(\frac{1}{h_j} \underline{\underline{I}} + \frac{1}{2} \underline{\underline{A}}_{j-1/2} \right) \underline{z}_{j+1} - g(t_{j+1/2})$$

$$\Rightarrow \boxed{N_h \underline{x}_j - N_h \underline{z}_j = \underline{\underline{R}}_j \underline{w}_j + \underline{\underline{S}}_j \underline{w}_{j+1}}$$

where $\underline{\underline{R}}_j = \frac{1}{h_j} \underline{\underline{I}} - \frac{1}{2} \underline{\underline{A}}_{j-1/2}$, $\underline{\underline{S}}_j = \left(\frac{1}{h_j} \underline{\underline{I}} + \frac{1}{2} \underline{\underline{A}}_{j-1/2} \right)$

Define: $\underline{s}_{j+1} = \underline{\underline{S}}_j \underline{w}_{j+1} + \underline{\underline{R}}_j \underline{w}_j$, $j = 1, \dots, N$

$$\underline{s}_N = \underline{\underline{B}}_0 \underline{w}_0 + \underline{\underline{B}}_1 \underline{w}_N$$

$$\Rightarrow \begin{bmatrix} \underline{\underline{S}}_1 & \underline{\underline{R}}_1 \\ \underline{\underline{S}}_2 & \underline{\underline{R}}_2 \\ \vdots & \ddots \\ \vdots & \ddots \\ \underline{\underline{S}}_N & \underline{\underline{R}}_N \\ \underline{\underline{B}}_0 & \underline{\underline{B}}_1 \end{bmatrix} \begin{pmatrix} \underline{w}_0 \\ \underline{w}_1 \\ \vdots \\ \vdots \\ \underline{w}_{N-1} \\ \underline{w}_N \end{pmatrix} = \begin{pmatrix} \underline{s}_0 \\ \vdots \\ \vdots \\ \vdots \\ \underline{s}_{N-1} \\ \underline{s}_N \end{pmatrix}$$

$\underbrace{\quad \quad \quad}_{\underline{\underline{I}}}$

$$\Rightarrow \underline{\underline{S}} \underline{w} = \underline{s}$$

(67)

$$\underline{\mathbb{I}} \text{ nonsingular} \Rightarrow \underline{w} = \underline{\mathbb{I}}^{-1} \underline{\delta}$$

For O-stability, must show that $\|\underline{\mathbb{I}}^{-1}\| \leq k$

Consider the inverse of $\underline{\mathbb{I}}$:

First look at calculation:

$$\underline{\mathbb{R}}_n^{-1} \underline{\Sigma}_n = -\left(\frac{1}{h_n} \underline{\mathbb{I}} - \frac{1}{2} \underline{\mathbb{A}}_{n-1/2}\right)^{-1} \left(\frac{1}{h_n} \underline{\mathbb{I}} + \frac{1}{2} \underline{\mathbb{A}}_{n-1/2}\right)$$

$$= -\left(\underline{\mathbb{I}} - \frac{h_n}{2} \underline{\mathbb{A}}_{n-1/2}\right)^{-1} \left(\underline{\mathbb{I}} + \frac{h_n}{2} \underline{\mathbb{A}}_{n-1/2}\right)$$

$$= -\left(\underline{\mathbb{I}} + \frac{h_n}{2} \underline{\mathbb{A}}_{n-1/2} + \dots\right) \left(\underline{\mathbb{I}} + \frac{h_n}{2} \underline{\mathbb{A}}_{n-1/2}\right)$$

$$= -\underbrace{\left(\underline{\mathbb{I}} + h_n \underline{\mathbb{A}}_{n-1/2} + O(h_n^2)\right)}$$

approximation of fundamental solution matrix $\underline{\Sigma}_n(t_n) + \dots$

$$\text{where } \underline{\Sigma}'_n = \underline{\mathbb{A}}(t) \underline{\Sigma}_n, \quad \underline{\Sigma}_n(t_{n+1}) = \underline{\mathbb{I}}$$

why? consider one step of midpoint/Euler

$$\underline{\Sigma}_n(t_n) = \underbrace{\underline{\Sigma}_n(t_{n-1})}_{\underline{\mathbb{I}}} + h_n \underline{\mathbb{A}}(t_{n-1/2}) \underbrace{\underline{\Sigma}_n(t_{n-1})}_{\underline{\mathbb{I}}} + \dots$$

Now set $\underline{\mathbb{D}} \underline{\mathbb{I}} = \underline{\mathbb{M}} + \underline{\mathbb{E}}$

where $\underline{\mathbb{D}} = \begin{bmatrix} \underline{\mathbb{R}}_1^{-1} & & & \\ & \underline{\mathbb{R}}_2^{-1} & & \\ & & \ddots & \\ & & & \underline{\mathbb{R}}_n^{-1} \\ & & & \underline{\mathbb{I}} \end{bmatrix}$, where $\underline{\mathbb{D}}$ is block diagonal

$$\Rightarrow \underline{\mathbb{M}} = \begin{bmatrix} -\underline{\Sigma}_1(t_1) & \underline{\mathbb{I}} & & & \\ & -\underline{\Sigma}_2(t_2) & \underline{\mathbb{I}} & & \\ & & -\underline{\Sigma}_3(t_3) & \underline{\mathbb{I}} & \\ & & & \ddots & \\ & & & & -\underline{\Sigma}_n(t_n) & \underline{\mathbb{I}} \\ & & & & & \underline{\mathbb{R}}_1 \end{bmatrix}$$

$\underline{\mathbb{E}}$ is a matrix of size $O(h^2)$ with the same zero structure as

$$\underline{\Sigma}^{-1} = \left(\underline{D}^{-1} (\underline{M} + \underline{E}) \right)^{-1} \rightarrow \underline{\Sigma}^{-1} = (\underline{M} + \underline{E})^{-1} \underline{D}$$

Let $(\underline{M} + \underline{E})^{-1} = \underline{M}^{-1}(\underline{I} - \underline{x})$, where \underline{x} is a small quantity

$$\begin{aligned} \rightarrow \underline{\Sigma} + \underline{E} &= (\underline{I} - \underline{x})^{-1} \underline{M} \\ &= (\underline{I} + \underline{x} + \dots) \underline{M} \\ &= \underline{M} + \underline{x} \underline{M} + \dots \rightarrow \underline{E} = \underline{x} \underline{M} \rightarrow \underline{x} = \underline{E} \underline{M}^{-1} \end{aligned}$$

$$\Rightarrow (\underline{M} + \underline{E})^{-1} = \underline{M}^{-1}(\underline{I} - E \underline{M}^{-1} + \dots)$$

Green's Functions

$$\rightarrow \underline{\Sigma}^{-1} = \underline{M}^{-1}(\underline{I} - E \underline{M}^{-1} + \dots) \underline{D}$$

$$\underline{M}^{-1} = \begin{bmatrix} \underline{G}(t_0, t_1) & \underline{G}(t_0, t_2) & \underline{G}(t_0, t_3) & \dots & \underline{G}(t_0, t_n) & \underline{E}(t_0) \\ \underline{G}(t_1, t_1) & \underline{G}(t_1, t_2) & - & - & \dots & \underline{G}(t_1, t_n) & \underline{E}(t_1) \\ \vdots & \vdots & & & & & \vdots \\ \underline{G}(t_n, t_1) & \underline{G}(t_n, t_2) & - & - & \dots & \underline{G}(t_n, t_n) & \underline{E}(t_n) \end{bmatrix}$$

IF $\|\underline{G}(t, z)\| \leq K$ and $\|\underline{E}\| \leq K$

$$\text{then } \|\underline{M}^{-1}\| \leq N K = \frac{K}{h}$$

and also, $\|\underline{D}\| = O(h)$, $\|\underline{E}\| = O(h^2)$, ~~for~~

$$\Rightarrow \|\underline{\Sigma}^{-1}\| \leq \tilde{K} (1 + O(h)) \Rightarrow \text{as } h \rightarrow 0, \|\underline{\Sigma}^{-1}\| \text{ is bounded}$$

Higher Order Methods

1) Collocation

The idea with collocation is to assume that the soln can be described by a chosen function (a polynomial, for instance) and then choose the coeffs st the ODEs are solved exactly.

⇒ Families: Radau, Gauss, Lebatta
 ↑ ↑ ↑
 B. Euler Midpoint Trapezoidal
 {
 symmetric methods - well suited for BVPs

General Form:

$$\underline{y}_n = \underline{y}_{n-1} + h_n \sum_{i=1}^s b_i \underline{\Sigma}(t_{n-1} + c_i h_n, \underline{Y}_i)$$

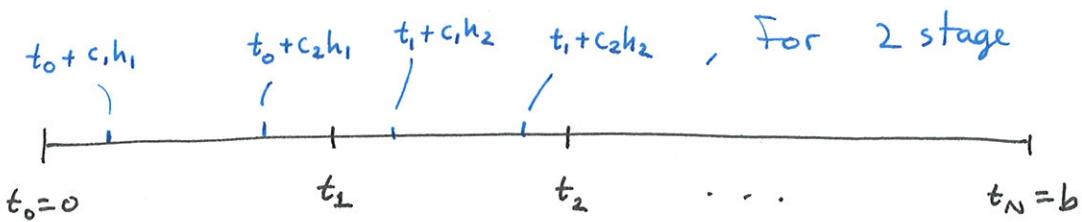
$$\underline{Y}_i = \underline{y}_{n-1} + h_n \sum_{j=1}^s a_{ij} \underline{\Sigma}(t_{n-1} + c_j h_n, \underline{Y}_j), \quad i = 1, \dots, s$$

e.g., Gauss, s=2 ($p=4$)

$$c_1 = \frac{3-\sqrt{3}}{6} \begin{pmatrix} \frac{1}{4} & \frac{3-2\sqrt{3}}{12} \\ \frac{3+2\sqrt{3}}{12} & \frac{1}{4} \end{pmatrix} \quad \left. \begin{array}{l} \{ a_{ij} \\ \} \\ \{ b_i \} \end{array} \right\}$$

$$c_2 = \frac{3+\sqrt{3}}{6} \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix}$$

Define a mesh: $0 = t_0 < t_1 < \dots < t_N = b$



so you have discretization points
and collocation points

Define \underline{y}_n 's on discretization points
and \underline{Y}_n 's on collocation points

Some Details for the Linear case

$$\underline{y}^{(t)} = \underline{\underline{A}} \underline{y}(t) + \underline{\underline{q}}(t)$$

discretize using Gauss with $s=2$, (4th order method)

$$\underline{y}_n = \underline{y}_{n-1} + h_n \left[b_1 \left(\underline{\underline{A}}(\underline{t}_{n-1} + c_1 h_n) \underline{Y}_{n1} + \underline{\underline{q}}(\underline{t}_{n-1} + c_1 h_n) \right) + b_2 \left(\underline{\underline{A}}(\underline{t}_{n-1} + c_2 h_n) \underline{Y}_{n2} + \underline{\underline{q}}(\underline{t}_{n-1} + c_2 h_n) \right) \right]$$

For notation, define $\underline{\underline{A}}_{n1} = \underline{\underline{A}}(\underline{t}_{n-1} + c_1 h_n)$

$$\underline{\underline{A}}_{n2} = \underline{\underline{A}}(\underline{t}_{n-1} + c_2 h_n)$$

$$\underline{\underline{q}}_{n1} = \underline{\underline{q}}(\underline{t}_{n-1} + c_1 h_n)$$

$$\underline{\underline{q}}_{n2} = \underline{\underline{q}}(\underline{t}_{n-1} + c_2 h_n)$$

$$\rightarrow \underline{y}_n = \underline{y}_{n-1} + h_n \left[b_1 (\underline{\underline{A}}_{n1} \underline{Y}_{n1} + \underline{\underline{q}}_{n1}) + b_2 (\underline{\underline{A}}_{n2} \underline{Y}_{n2} + \underline{\underline{q}}_{n2}) \right], n=1, \dots, N$$

Collocation PB:

$$\underline{Y}_n = \underline{X}$$

$$\underline{Y}_{n1} = \underline{y}_{n-1} + h_n (a_{11} (\underline{\underline{A}}_{n1} \underline{Y}_{n1} + \underline{\underline{q}}_{n1}) + a_{12} (\underline{\underline{A}}_{n2} \underline{Y}_{n2} + \underline{\underline{q}}_{n2}))$$

$$\underline{Y}_{n2} = \underline{y}_{n-1} + h_n (a_{21} (\underline{\underline{A}}_{n1} \underline{Y}_{n1} + \underline{\underline{q}}_{n1}) + a_{22} (\underline{\underline{A}}_{n2} \underline{Y}_{n2} + \underline{\underline{q}}_{n2}))$$

BCs

$$\hat{\underline{\underline{B}}}_0 \underline{y}_0 = \underline{y}_0, \quad \hat{\underline{\underline{B}}}_N \underline{y}_N = \underline{y}_1$$

Number of Eqns: $(3N+1)m = \underline{\text{Number of unknowns}}$

The basic structure of the resulting system of eqns is:

$$\left[\begin{array}{c|ccccc|c} BC & * & * & * & & & Y_0 \\ & * & * & * & & & Y_{1,1} \\ & * & * & * & & & Y_{1,2} \\ & * & * & * & * & & Y_2 \\ & & * & * & * & & Y_{2,1} \\ & & * & * & * & & Y_{2,2} \\ & & * & * & * & * & Y_3 \\ & & & \ddots & \ddots & \ddots & \vdots \\ & & & & & * & Y_N \end{array} \right] = q$$

4/

Higher Order Methods

BVP: $y' = f(t, y)$, $0 \leq t \leq b$, $g(y_0, y(b)) = 0$
 Various approaches to obtain higher order ($p > 2$) finite difference methods.

1) collocation methods

e.g. midpoint \rightarrow Gauss methods
 trapezoidal \rightarrow Dahlatto methods (sym RK schemes)

2) Richardson Extrapolation

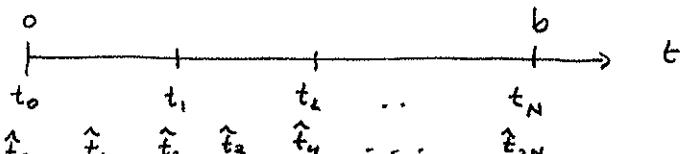
idea is to use linear combinations of approximate solutions of lower order methods to construct higher order approximations

e.g. approximate solution using midpoint on grid with uniform spacing h .

$$\textcircled{1} \quad y_n - y(t_n) = c(t_n) h^2 + O(h^4), \quad n=0, 1, \dots, N$$

Now consider the approximate solution using midpoint on grid with uniform spacing $\hat{h} = \frac{h}{2}$.

original grid
new grid



Have for the solution on the finer grid

$$\textcircled{2} \quad \hat{y}_n - \underline{y}(t_n) = \underline{\epsilon}(t_n) \hat{h}^2 + O(\hat{h}^4), \quad n=0, \dots, 2N$$

where $\underline{\epsilon}(t)$ is the same function for both grids - $\underline{\epsilon}(t)$ only depends on the method and is independent of grid size.

Consider the finer solution at every other grid point, or where $\textcircled{1}$ and $\textcircled{2}$ have common points:

$$\hat{y}_{2n} - \underline{y}(t_{2n}) = \underline{\epsilon}(t_{2n}) \hat{h}^2 + O(\hat{h}^4), \quad n=0, \dots, N$$

$$\text{use } \hat{t}_{2n} = t_n, \quad \hat{h} = \frac{h}{2}$$

$$\textcircled{3} \rightarrow \hat{y}_{2n} - \underline{y}(t_n) = \underline{\epsilon}(t_n) \frac{h^2}{4} + O(h^4)$$

Eliminate 2nd order term in error from $\textcircled{1}$ and $\textcircled{3}$

$$\text{i.e. } 4 \times \textcircled{3} - \textcircled{1}$$

$$\Rightarrow 4 \hat{y}_{2n} - \underline{y}_n - 3 \underline{y}(t_n) = O(h^4)$$

$$\Rightarrow \underbrace{\frac{4 \hat{y}_{2n} - \underline{y}_n}{3} - \underline{y}(t_n)}_{\text{this is a 4th order accurate approximation of } \underline{y}(t_n)} = O(h^4)$$

Can also obtain an error estimate: eqn $\textcircled{1}$ - eqn $\textcircled{3}$

$$\textcircled{3} \quad \hat{y}_{2n} - \underline{y}(t_n) = \underline{\epsilon}(t_n) \frac{h^2}{4} + O(h^4)$$

$$\textcircled{1} \quad \underline{y}_n - \underline{y}(t_n) = \underline{\epsilon}(t_n) h^2 + O(h^4)$$

$$\Rightarrow \underline{y}_n - \hat{y}_{2n} = \frac{3}{4} \underline{\epsilon}(t_n) h^2 + O(h^4)$$

$$\rightarrow \underline{\epsilon}(t_n) h^2 = \underbrace{\frac{4}{3} (\underline{y}_n - \hat{y}_{2n})}_{\text{error estimate}} + O(h^4)$$

Continuation Methods

A discretization of a BVP (either by shooting), finite differences, or some other method) leads to a set of nonlinear algebraic eqns to solve for the discrete solution.

Let us refer to the collection of discrete unknowns as \underline{u} and the set of equations that define them as \underline{G} . We want to solve $\underline{G}(\underline{u}) = \underline{0}$. (In the case of shooting, \underline{u} is the shooting parameter(s). In the case of finite differences, \underline{u} is the solution on the grid...).

$\underline{G}(\underline{u}) = \underline{0}$ is a standard root-finding problem. The usual approach is some iteration, ie Newton's method or variation of it, and we need some initial guess to get started.

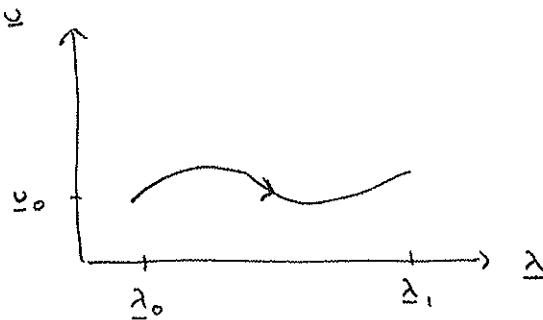
Continuation is an approach that can be used to obtain initial guesses in a rational way.

Typically, the equations (exact or discrete) involve a set of parameters, $\underline{\lambda}$ say.

$$\Rightarrow \underline{G}(\underline{u}, \underline{\lambda}) = \underline{0} \Rightarrow \underline{u} \text{ depends on } \underline{\lambda}$$

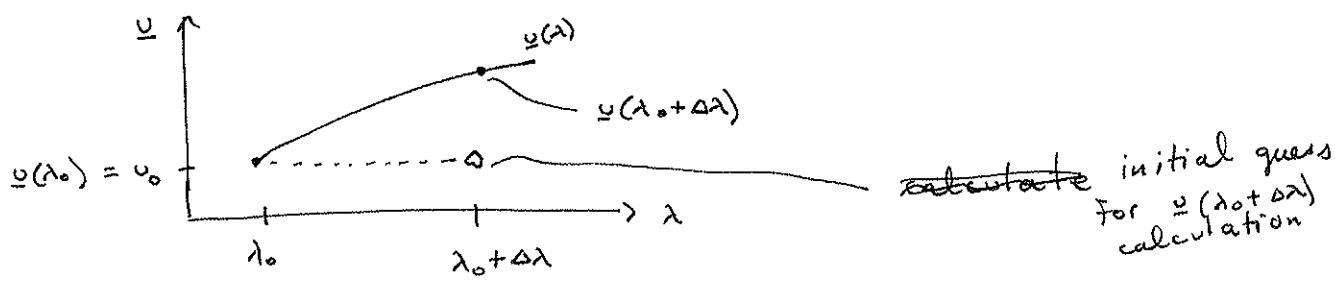
Exploit this dependence to obtain a good initial guess.

More generally, might want to determine how \underline{u} depends on $\underline{\lambda}$.

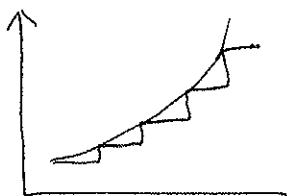


Suppose want to compute u at $\lambda = \lambda_+$ (hard) and we start at $\lambda = \lambda_0$ where the solution is easy (perhaps linear for instance). The question is, how can we move in parameter space in a systematic way \Rightarrow continuation.

Consider a step in parameter space: (with λ a scalar)



"easy", "trivial" continuation - use u_0 as initial guess for Newton at $\lambda = \lambda_0 + \Delta \lambda$ and then converge (hopefully) to $u(\lambda_0 + \Delta \lambda)$ using Newton, with λ fixed



this often works but there are more sophisticated ideas.

Instead of stair-casing along, why not find tangents...

"Euler-Newton" continuation: Regard u as a differentiable function of scalar parameter λ .

$$\Rightarrow G(u(\lambda), \lambda) = 0$$

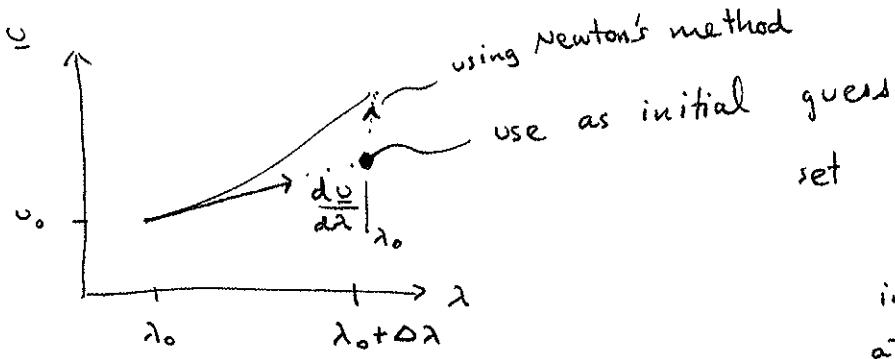
Take a derivative with respect to λ :

$$G_u(u, \lambda) \frac{du}{d\lambda} + G_\lambda(u, \lambda) = 0 \quad (\text{linear system})$$

for $\frac{du}{d\lambda}$

Apply at u_0, λ_0 . Define $G_u(u_0, \lambda_0) = G_u^0$, $G_\lambda(u_0, \lambda_0) = G_\lambda^0$

Solve $\left. G_u^0 \frac{du}{d\lambda} \right|_{\lambda=\lambda_0} = -G_\lambda^0$ for $\frac{du}{d\lambda}$



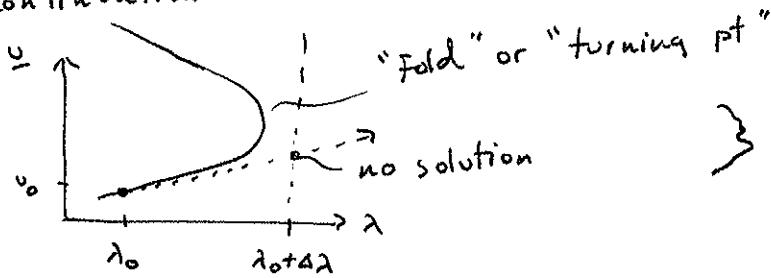
set $\tilde{u} = u_0 + \Delta \lambda \left. \frac{du}{d\lambda} \right|_{\lambda=\lambda_0}$

initial guess for Newton
at $\lambda_0 + \Delta \lambda$

Almost all of this information is available as is, in the code.
Certainly the Jacobian is known. The solution of the linear
system \bullet comes for free

Arc length Continuation

A difficulty occurs for either the "trivial" or "Euler-Newton"
continuation near "folds" or "turning points".



} multiple solutions exist

Regard u and λ depending on a parameter s :

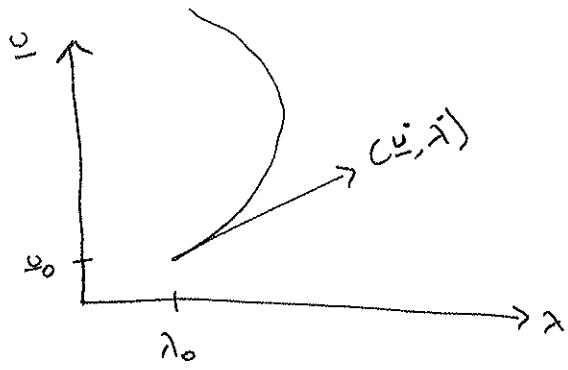
$$\Rightarrow G(u(s), \lambda(s)) = 0$$

Pick s : For example, if s is arc length then

$$\| \dot{\vec{u}}(s) \|^2 + \dot{\lambda}(s)^2 = 1, \text{ where } \cdot = \frac{d}{ds}$$

Find $\dot{u}, \dot{\lambda}$ at (λ_0, u_0) (ie the tangent at λ_0, u_0)

Stems from
 $ds^2 = dx^2 + dy^2$
 $\Rightarrow 1 = \dot{x}^2 + \dot{y}^2$



To find tangent,
differentiate wrt s
 $\underline{G}(\underline{v}(s), \lambda(s)) = 0$

$$\Rightarrow \underline{G}_v(v_0, \lambda_0) \dot{v} + \underline{G}_\lambda(v_0, \lambda_0) \dot{\lambda} = 0 \quad \star$$

matrix | vector |
 | |
 scalar

this represents m equations

set $\dot{v} = \alpha \underline{\phi}$, where α = scalar, $\underline{\phi}$ is a vector that solves

$$\underline{G}_v(v_0, \lambda_0) \underline{\phi} = \underline{G}_\lambda(v_0, \lambda_0)$$

Substitute $\dot{v} = \alpha \underline{\phi}$ into \star : $\underline{G}_v(v_0, \lambda_0) \alpha \underline{\phi} + \underline{G}_\lambda(v_0, \lambda_0) \dot{\lambda} = 0$

$$\rightarrow \cancel{\alpha} \underline{G}_\lambda(v_0, \lambda_0) + \underline{G}_\lambda(v_0, \lambda_0) \dot{\lambda} \rightarrow (\alpha + \dot{\lambda}) \underline{G}_\lambda = 0 \rightarrow \boxed{\alpha = -\dot{\lambda}}$$

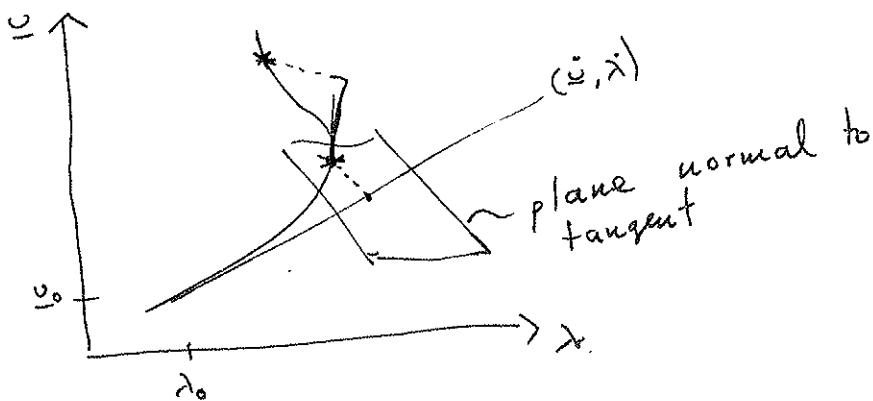
$$\Rightarrow \dot{v} = -\dot{\lambda} \underline{\phi}$$

therefore the equation $\|\dot{v}(s)\|^2 + \dot{\lambda}(s)^2 = 1$ becomes

$$\|\dot{\lambda} \underline{\phi}\|^2 + \dot{\lambda}^2 = 1 \rightarrow \dot{\lambda}^2 \|\underline{\phi}\|^2 + \dot{\lambda}^2 = 1 \rightarrow \boxed{\dot{\lambda} = \pm \frac{1}{\sqrt{1 + \|\underline{\phi}\|^2}}}$$

+ value for \nearrow direction, - value for \nwarrow

Now that we have calculated $(\dot{v}, \dot{\lambda})$ at (v_0, λ_0) . Now want to move back to curve but not at fixed λ , which gave us problems with this fold. Consider a plane normal to the tangent. Use Newton's method to move along the plane back to the curve.



Eqn for the curve: $\underline{G}(\underline{v}, \lambda) = 0$

Eqn for the plane: $N(\underline{v}, \lambda) = \dot{\underline{v}}^T (\underline{v} - \underline{v}_0) + \lambda (\lambda - \lambda_0) - \Delta s = 0$ (^{scalar} eqn)

Perform Newton on $\underline{G} = 0$, $N = 0$ to find \underline{v}, λ for a given Δs . Need to calculate a Jacobian:

$\underline{v} = \underline{\tilde{v}} + \underline{\Delta v}$, $\lambda = \tilde{\lambda} + \Delta \lambda$, where $\underline{\tilde{v}}, \tilde{\lambda}$ are initial guess

$$\underline{G}(\underline{\tilde{v}} + \underline{\Delta v}, \tilde{\lambda} + \Delta \lambda) = 0$$

$$\rightarrow \underline{G}(\underline{\tilde{v}}, \tilde{\lambda}) + \underline{G}_{\underline{v}}(\underline{\tilde{v}}, \tilde{\lambda}) \underline{\Delta v} + \underline{G}_{\lambda}(\underline{\tilde{v}}, \tilde{\lambda}) \Delta \lambda + \dots = 0$$

$$\textcircled{1} \rightarrow \underline{G}_{\underline{v}}(\underline{\tilde{v}}, \tilde{\lambda}) \underline{\Delta v} + \underline{G}_{\lambda}(\underline{\tilde{v}}, \tilde{\lambda}) \Delta \lambda = -\underline{G}(\underline{\tilde{v}}, \tilde{\lambda})$$

neglect higher order terms

- linear equation
for $\underline{\Delta v}, \Delta \lambda$

and

$$N(\underline{\tilde{v}} + \underline{\Delta v}, \tilde{\lambda} + \Delta \lambda) = 0 \quad \text{similarly:}$$

$$\rightarrow N(\underline{\tilde{v}}, \tilde{\lambda}) + \underbrace{N_{\underline{v}}(\underline{\tilde{v}}, \tilde{\lambda}) \underline{\Delta v}}_{\dot{\underline{v}}^T} + \underbrace{N_{\lambda}(\underline{\tilde{v}}, \tilde{\lambda}) \Delta \lambda}_{\tilde{\lambda}} = 0$$

$$\textcircled{2} \rightarrow \dot{\underline{v}}^T \underline{\Delta v} + \tilde{\lambda} \Delta \lambda = -N(\underline{\tilde{v}}, \tilde{\lambda})$$

Combine \textcircled{1} and \textcircled{2} into a linear system

$$\begin{bmatrix} \underline{G}_{\underline{v}} & \vdots & \underline{G}_{\lambda} \\ \hline \dot{\underline{v}}^T & \vdots & \tilde{\lambda} \end{bmatrix} \begin{pmatrix} \underline{\Delta v} \\ \Delta \lambda \end{pmatrix} = \begin{pmatrix} -\underline{G} \\ -N \end{pmatrix}$$

this becomes the problem to solve for the corrections $\underline{\Delta v}$ and $\Delta \lambda$

Must update both solutions

$$\underline{v}^{\text{NEW}} = \underline{\tilde{v}} + \underline{\Delta v}$$

$$\lambda^{\text{NEW}} = \tilde{\lambda} + \Delta \lambda$$

typically solve using efficient algorithms

Review - since Midterm

1) linear multi-step methods (IVPs)

- a) Adams Family For non-stiff problems
 - explicit ($B_0 = 0$) vs implicit ($B_0 \neq 0$)
- b) Backwards Difference Formula (BDF)
 - For stiff problems
- c) starting values
- d) order of accuracy \rightarrow order conditions
- e) O -stability, root-condition
- f) absolute stability
- g) implementation issues
 - predictor/corrector scheme
 - error estimates (Milne estimate)
 - variable mesh size

2) BVP - theory

- a) general problem, locally isolated solutions, linearization
- b) general solution of linear BVP \Rightarrow Green's Function
- c) problem stability \rightarrow dichotomy, exponential dichotomy

3) BVP - Shooting Methods

- a) general framework, Newton's method
- b) difficulties: if BVP is stable, IVP is unstable
- c) multiple shooting (for linear case)

4) BVP - finite difference Methods

- a) linear case, midpoint method \rightarrow linear system
- b) nonlinear case, midpoint method \rightarrow Newton \rightarrow linear system
- c) consistency, O -stability, convergence
- d) higher order methods
 - collocation and extrapolation
- e) continuation methods

Take-Home Final Exam Rules

- email to get the final
- 1st available time, saturday morning (May 3rd), to return monday morning
- 48 hrs to complete the exam
- All exams to be completed by 5pm wed May 7th
- NO COLLABORATION or TEXTBOOKS