

Numerical Analysis Formulae

Page No.:

Date:

YOUNV
1

LAGRANGE'S INTERPOLATION TECHNIQUE

- Theorem (by Lagrange):

Take $n+1$ distinct points $x_0, x_1, \dots, x_n \in \mathbb{R}$

and take real nos. f_0, f_1, \dots, f_n .

Then \exists a UNIQUE Polynomial $p \in P_n$ such that $p(x_i) = f_i$ for $i=0, 1, \dots, n$.

- Fundamental Theorem of Algebra:

A polynomial of degree n can have at most n real zeros (or exactly n complex zeros)

Polynomial Interpolation — Recipe of Lagrange:

- Start with the $(n+1)$ pairs

$$(x_0, f_0), (x_1, f_1), (x_2, f_2), \dots, (x_n, f_n)$$

- Build the $n+1$ Lagrange polynomials, each of degree n :

$$L_k^n(x) = \prod_{\substack{j=0 \\ j \neq k}}^{n-1} \frac{(x - x_j)}{(x_k - x_j)}$$

- Construct the polynomial interpolant $p \in P_n$ as a weighted sum:

$$p(x) = \sum_{k=0}^n f_k L_k^n(x)$$

ERROR COMMITTED BY POLYNOMIAL APPROX.

- To measure closeness between two graphs f and g , we use $\max_{x \in [a, b]} |f(x) - g(x)|$

$\max_{x \in [a, b]} |f(x)|$ measures how far the fn is from zero.

- We define "Maximum Norm" $\| \cdot \|$ as:

$$\|f\| = \max_{x \in [a, b]} |f(x)|$$

- Weierstrass Approximation Theorem:

Take a fn. $f \in C[a, b]$.

Given a real no. $\epsilon > 0$, \exists a polynomial p such that $\|f - p\| < \epsilon$

(not a constructive proof)

In order to find p , we pick points on x -axis (using Chebyshev's idea; mostly preferred), then construct p using Lagrange's recipe.

- Theorem:

Take $f \in C^{n+1}[a, b]$.

Let x_0, x_1, \dots, x_n be $n+1$ distinct points in $[a, b]$. And let $p \in P_n$ be such that

$$p(x_i) = f(x_i) \text{ for } i=0, 1, \dots, n.$$

Then $\forall x \in [a, b]$, \exists a point $\xi \in [a, b]$ such that $f(x) - p(x) = \frac{1}{(n+1)!} f^{(n+1)}(\xi) \prod_{k=0}^n (x - x_k)$

(Thus error seems to depend on interpolation pts. x_k)

Using the theorem:

$$\max_{x \in [a,b]} |f(x) - p(x)| \leq \frac{1}{(n+1)!} \left\| f^{(n+1)} \right\| \prod_{k=0}^n |x - x_k|$$

∴ the error depends only on the no. of interpolation pts. and

$$\max_{x \in [a,b]} \prod_{k=0}^n |x - x_k|$$

NEWTON'S DIVIDED DIFFERENCES

Suppose $f \in C[a,b]$. Let $p \in P^n$ be the polynomial that agrees with f at x_0, x_1, \dots, x_n . Objective is to determine $n+1$ coefficients a_0, a_1, \dots, a_n of

$$p(x) = a_0 + a_1(x - x_0) + a_2(x - x_0)(x - x_1) + \dots$$

$$p(x_i) = f(x_i) \text{ for } i = 0, 1, \dots, n$$

$$\text{Clearly, } a_0 = p(x_0) = f(x_0)$$

$$a_1 = \frac{f(x_1) - f(x_0)}{x_1 - x_0}$$

Definition of divided difference:

Given distinct points x_0, x_1, \dots, x_n in $[a,b]$. Let $p \in P^n$ be the polynomial that interpolates f at those points. The coefficient of x^n in p is called the "divided difference" i.e. $f[x_0, x_1, \dots, x_n]$

$$f[x_0, x_1, \dots, x_n] = \sum_{k=0}^n f(x_k) \prod_{j=0, j \neq k}^n \frac{1}{(x_k - x_j)}$$

Recurrence Relation for divided difference :

$$f[x_0, x_1, \dots, x_{m+1}] = \underbrace{f[x_1, x_2, \dots, x_{m+1}] - f[x_0, x_1, \dots, x_m]}_{x_{m+1} - x_0}$$

Thus,

$$P_n(x) = f[x_0] + f[x_0, x_1](x - x_0) + \dots + f[x_0, x_1, \dots, x_n] \prod_{k=0}^{n-1} (x - x_k)$$

satisfying $p_n(x_i) = f(x_i)$ for $i = 0, 1, \dots, n$.

Divided difference table :-

$f[x_0]$	$f[x_0, x_1]$	\vdots	$f[x_0, x_1, \dots, x_n]$
$f[x_1]$	$f[x_1, x_2]$	\vdots	
\vdots	\vdots	\vdots	
$f[x_n]$	$f[x_{n-1}, x_n]$	\vdots	

OPERATION COUNT

Lagrange : $p(x) = f(x_0) \prod_{j=0}^{n-1} \frac{(x - x_j)}{(x_0 - x_j)} + \dots + f(x_n) \prod_{j=n}^{n-1} \frac{(x - x_j)}{(x_n - x_j)}$

Addition	Subtraction	Multiplication	Division
n	$2n(n+1)$	$2n^2 + n - 1$	$n+1$

Newton : $p(x) = f(x_0) + f[x_0, x_1](x - x_0) + \dots + f[x_0, \dots, x_n]$

Addition	Subtraction	Multiplication	Division
n	$\frac{3n(n+1)}{2}$	$\frac{n(n+1)}{2}$	$\frac{n(n+1)}{2}$

PIECEWISE / SPLINE INTERPOLATION

To interpolate a $f \in C[a, b]$.

Pick some points from $[a, b]$

$$a = x_0 < x_1 < x_2 < \dots < x_n = b$$

∴ Data points are $(x_0, f_0), (x_1, f_1), \dots, (x_n, f_n)$

Fix an integer $m \leq n$

On each subinterval $[x_{i-1}, x_i]$, build $\Psi \in C[a, b]$
such that $\Psi(x_{i-1}) = f_{i-1}$, $\Psi(x_i) = f_i$.

$$\Psi(x) = \begin{cases} a_0^1 + a_1^1 x + a_2^1 x^2 + \dots + a_m^1 x^m & \text{on } [x_0, x_1] \\ a_0^2 + a_1^2 x + a_2^2 x^2 + \dots + a_m^2 x^m & \text{on } [x_1, x_2] \\ \vdots & \vdots \\ a_0^n + a_1^n x + a_2^n x^2 + \dots + a_m^n x^m & \text{on } [x_{n-1}, x_n] \end{cases}$$

considering the sub-intervals $[x_{i-1}, x_i]$ for $i = 1, 2, \dots, n$

Hence, there are $n \cdot (m+1)$ coeffs to be determined:

- $n+1$ interpolation conditions:

$$\Psi(x_i) = f(x_i) = f_i \text{ for } i = 0, 1, \dots, n$$

- continuity of derivatives at interior pts:

$$\lim_{h \rightarrow 0^+} \Psi(x_i - h) = \lim_{h \rightarrow 0^+} \Psi(x_i + h) \quad i = 1, \dots, n-1$$

$$\lim_{h \rightarrow 0^+} \frac{d\Psi}{dx}(x_i - h) = \lim_{h \rightarrow 0^+} \frac{d\Psi}{dx}(x_i + h)$$

$$\lim_{h \rightarrow 0^+} \frac{d^{m-1}\Psi}{dx^{m-1}}(x_i - h) = \lim_{h \rightarrow 0^+} \frac{d^{m-1}\Psi}{dx^{m-1}}(x_i + h)$$

∴ $n(m+1) = (m+1)(n-1)m$ conditions.

Hence, total of $(m+1)n + (1-m)$ conditions
In need of $(m-1)$ more conditions.

ERROR BOUND IN LINEAR SPLINES

Theorem:

Take $f \in C^2[a, b]$.

Let $s_L(x)$ be the linear spline interpolating f at the $(n+1)$ knots.

$$a = x_0 < x_1 < \dots < x_n = b$$

Let $h > 0$ be defined as

$$h = \max_{1 \leq i \leq n} h_i = \max_{1 \leq i \leq n} (x_i - x_{i-1})$$

Then the following error bound holds true:

$$\|f - s_L\| \leq \frac{h^2}{8} \|f''\|_{\text{uniform}}$$

Note: "Natural" cubic splines are those in which the second derivative at endpoints of domain is zero.

$$\text{i.e. } s''_0(x_0) = s''_{n-1}(x_n) = 0$$

NATURAL CUBIC SPLINES

$$s_i(x) = a_i(x-x_i)^3 + b_i(x-x_i)^2 + c_i(x-x_i) + d_i$$

From condition (i), $d_i = f_i$, $i=0, 1, \dots, n-1$

Define new variables

$$\sigma_i = s''(x_i) \text{ for } i=0, 1, \dots, n$$

For the case of Equally spaced interpolation points,

$$h = x_{i+1} - x_i$$

Since the spline is natural,

$$\sigma_0 = \sigma_n = 0$$

Note that $\sigma_i = s''(x_i) = 2b_i$,

$$\text{i.e. } b_i = \frac{\sigma_i}{2}$$

Using Condition (v),

$$s_i''(x_{i+1}) = s_{i+1}''(x_{i+1}) \\ \therefore \sigma_{i+1} = 6a_i h + 2b_i$$

$$\therefore a_i = \frac{\sigma_{i+1} - \sigma_i}{6h}, \quad i=0, 1, \dots, n-1$$

$$\therefore \text{As of now, } f_{i+1} = \left(\frac{\sigma_{i+1} - \sigma_i}{6h} \right) h^3 + \frac{\sigma_i}{2} h^2 + c_i h + f_i$$

This yields

$$\frac{f_{i+1} - f_i}{h} = \left(\frac{\sigma_{i+1} + 2\sigma_i}{6} \right) h + c_i$$

$$\Rightarrow c_i = \frac{f_{i+1} - f_i}{h} - \frac{h}{6} (2\sigma_i + \sigma_{i+1})$$

Using Condition (iv),

$$s'_i(x_{i+1}) = s'_{i+1}(x_{i+1})$$

$$3a_i h^2 + 2b_i h + c_i = c_{i+1}$$

Finally we get :

$$\sigma_{i-1} + 4\sigma_i + \sigma_{i+1} = \frac{6}{h^2} (f_{i-1} - 2f_i + f_{i+1})$$

$$\text{for } i = 1, \dots, n-1$$

$$\begin{pmatrix} 4 & 1 & 0 & \cdots & 0 & 0 \\ 1 & 4 & 1 & & & \\ 0 & 1 & 4 & 1 & & \\ \vdots & \vdots & \vdots & \ddots & 0 & 0 \\ 0 & \vdots & \vdots & \ddots & 1 & 0 \\ 0 & 0 & \cdots & \ddots & 1 & 4 & 1 \\ 0 & 0 & \cdots & \cdots & 1 & 4 & 1 \end{pmatrix} \begin{pmatrix} \sigma_1 \\ \sigma_2 \\ \vdots \\ \vdots \\ \vdots \\ \sigma_{n-1} \\ \sigma_n \end{pmatrix} = \frac{6}{h^2} \begin{pmatrix} f_0 - 2f_1 + f_2 \\ f_1 - 2f_2 + f_3 \\ \vdots \\ \vdots \\ f_{n-2} - 2f_{n-1} + f_n \end{pmatrix}$$

↓
Diagonally Dominant

Invertible

$$d_i = f_i$$

$$b_i = \frac{\sigma_i}{2}$$

$$a_i = \frac{\sigma_{i+1} - \sigma_i}{6h}$$

$$c_i = \frac{f_{i+1} - f_i}{6h} - \frac{h}{6}(2\sigma_i + \sigma_{i+1})$$

ERROR ESTIMATE IN NATURAL CUBIC SPLINE

Let $f \in C^4[a, b]$

Let $s \in C^2[a, b]$ be the natural cubic spline that interpolates f at $(n+1)$ equally-spaced knots.

$$a = x_0 < x_1 < \cdots < x_n = b$$

$$h = |x_{i+1} - x_i|, i = 0, 1, \dots, n-1$$

Then the following error bound holds true:

$$\|f - s\| \leq \|f^{(iv)}\| h^4$$

NUMERICAL INTEGRATION

(9)

Given a real-valued fn. f , compute

$$\int_a^b f(x) dx$$

Best method is to find antiderivative F then do $F(b) - F(a)$: But that may not always be possible.

NUMERICAL QUADRATURE

Partition $[a, b]$ into n parts

$$a = x_0 < x_1 < x_2 < \dots < x_n = b$$

Interpolate f at these $n+1$ points by some $p \in P^n$.

$$\int_a^b f(x) dx \approx \int_a^b p(x) dx$$

NEWTON-COTES INTEGRATION FORMULAE

$n+1$ equally spaced partitioning points.

$$p(x) = \sum_{i=0}^n f(x_i) L_i(x), \quad L_i(x) = \prod_{\substack{k=0 \\ k \neq i}}^n \frac{(x - x_k)}{(x_i - x_k)}$$

$$x_i = a + ih$$

Introduce variable t such that $x = a + th$

$$t = \frac{x-a}{h} = n \frac{x-a}{b-a} \in [0, n]$$

$$\therefore \frac{x-x_k}{x_i-x_k} = \frac{t-k}{i-k}$$

$$\therefore L_i(x=t) = \prod_{\substack{k=0 \\ k \neq i}}^n \frac{t-k}{i-k} = \varphi_i(t)$$

(10)

$$\int L_i(x) dx = h \int \varphi_i(t) dt$$

Quadrature Weights:

$$w_i = h \int \varphi_i(t) dt$$

$$\therefore \int f(x) dx \approx h \sum_{i=0}^n w_i f(x_i)$$

Notice that $w_i \in \mathbb{Q}_{\text{Rational}}$ and depend only on n .

$w_i \leftrightarrow f(x_i), x_i \in A$

$w_i \leftrightarrow a, b$ i.e. boundaries of integral.

$$\sum w_i = n$$

$$\sum_{i=0}^n w_i = n$$

Trapezium Rule ($n=1$)

$$w_0 = \int_a^b \frac{t-1}{1-0} dt = \left[-\frac{t^2}{2} + t \right]_0^1 = \frac{1}{2}$$

$$w_1 = \int_0^1 \frac{t-0}{1-0} dt = \frac{1}{2} (1-x) + \frac{1}{2} = (x)$$

$$\therefore \int_a^b f(x) dx \approx \frac{b-a}{2} (f(a) + f(b))$$

Simpson's Rule ($n=2$)

$$w_0 = \int_a^b \frac{(t-1)(t-2)}{(0-1)(0-2)} dt = \frac{1}{3}$$

$$w_1 = \frac{4}{3}$$

$$w_2 = \frac{1}{3}$$

$$\int_a^b f(x) dx \approx \frac{h}{3} \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right)$$

Quadrature Weights

$$n \quad w_0 \quad w_1 \quad w_2 \quad w_3 \quad w_4$$

$$1 \quad \frac{1}{2} \quad \frac{1}{2}$$

$$2 \quad \frac{1}{3} \quad \frac{4}{3} \quad \frac{1}{3}$$

$$3 \quad \frac{5}{8} \quad \frac{9}{8} \quad \frac{9}{8} \quad \frac{3}{8}$$

$$4 \quad \frac{1}{45} \quad \frac{64}{45} \quad \frac{24}{45} \quad \frac{64}{45} \quad \frac{1}{45}$$

$$I_{P_1} = \frac{h}{2} (f(a) + f(b)) \quad \{ \text{Trapezoid Rule} \}$$

$$I_{P_2} = \frac{h}{3} (f(a) + 4f(a+h) + f(b)) \quad \{ \text{Simpson's Rule} \}$$

$$I_{P_3} = \frac{3h}{8} (f(a) + 3f(a+h) + 3f(a+2h) + f(b))$$

$$I_{P_4} = \frac{h}{45} (14f(a) + 64f(a+h) + 24f(a+2h) + 64f(a+3h) + 14f(b))$$

{Dividing [a,b] in 4 parts}

Error Estimate in Newton-Cotes

Take $f \in C^{n+1}[a, b]$

$$|I_f - I_{P_n}| \leq \frac{1}{(n+1)!} \|f^{(n+1)}\| \int_a^b \prod_{i=0}^n |x-x_i| dx$$

- Hence, error in trapezoidal rule $\leq \frac{1}{2} \|f''\| \frac{(b-a)^3}{6}$

- Newton-Cotes formula need not give convergence as $n \rightarrow \infty$, \because weights may be negative.

Hence, we use Gaussian quadrature $G(f) =$

$$W_i = \int_a^b [\varphi_i(x)]^2 dx = \int_a^b \prod_{k=0}^{n-1} \left(\frac{x-x_k}{x_i-x_k} \right)^2 dx$$

Here quadrature points are not equally spaced.
 Gaussian quadrature weights are non-negative. (12)

Gauss quadrature - Convergence Theorem :

Let $f \in C[a, b]$

Let $G_n(f)$ be the Gaussian quadrature formula of order n . Then $\lim_{n \rightarrow \infty} |G_n(f) - I_f| = 0$

- Composite Newton-Cotes

alternative to building Newton-Cotes for large n .

- Apply Newton-Cotes to small set of quadr. points.
- Only requirement is Continuity.

- Composite Trapezoidal Rule :

$$\text{For } m=2 : h \left(\frac{1}{2}f(a) + \frac{1}{2}f(a+h) \right) + h \left(\frac{1}{2}f(a+h) + \frac{1}{2}f(b) \right)$$

$$= h \left(\frac{1}{2}f(a) + f(a+h) + \frac{1}{2}f(b) \right)$$

$$\text{For } m \geq 2 : h = \frac{b-a}{m}$$

$$\therefore C_p(f) = h \left(\frac{1}{2}f(a) + f(a+h) + f(a+2h) + \dots + f(a+(m-1)h) + \frac{1}{2}f(b) \right)$$

- Composite Simpson's $\frac{1}{3}$ Rule :

Divide $[a, b]$ into $2m$ subintervals with $m \geq 2$

$$\text{For } m=2 : \frac{h}{3} (f(a) + 4f(a+h) + f(a+2h)) + \frac{h}{3} (f(a+2h) + 4f(a+3h) + f(b))$$

$$h = \frac{b-a}{2m}$$

$$\text{For } m \geq 2 : C_{p_2}(f) = \frac{h}{3} (f(a) + 4f(a+h) + 2f(a+2h) + 4f(a+3h) + 2f(a+4h) + \dots + 2f(a+(2m-2)h) + h f(a+(2m-1)h) + f(b))$$

$$\text{- Error Estimate : } |C_{p_2}(f) - I_f| \leq \frac{m}{12} \|f''\| h^3 = \frac{b-a}{12} \|f''\| h^2 \quad \therefore h \rightarrow 0 \text{ is preferred.}$$

NUMERICAL SOLUTIONS TO ODES

- Peano's Theorem for existence of a solution for $y' = f(t, y)$,
 $y(0) = y_0$:
 Suppose $f: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$.
 f is continuous in $(-\delta, \delta) \times (y_0 - \gamma, y_0 + \gamma)$
 for some $\delta > 0, \gamma > 0$
 $\Rightarrow \exists \epsilon > 0$ and a function $y \in C^1(-\epsilon, \epsilon)$
 such that $y' = f(t, y), y(0) = y_0$.

Lipschitz Continuity

A function $g: [a, b] \rightarrow \mathbb{R}$ is said to be Lipschitz continuous if

$$|g(x) - g(y)| \leq L|x - y| \quad \forall x, y \in [a, b]$$

with some positive constant L .

Cauchy - Lipschitz Uniqueness Theorem

Suppose $f(t, y)$ be continuous in both variables
 and let it be Lipschitz Continuous in y variable.
 Then \exists an $\epsilon > 0$ and a unique function $y \in C^1(-\epsilon, \epsilon)$
 such that $y' = f(t, y), y(0) = y_0$.

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

$$y(t) = y_0 + \int_0^t f(s, y(s)) ds \quad \forall t \geq 0$$

Numerical Scheme

Define mesh points $t_n = nh$



$$y_n \approx y(t_n)$$

Euler Method (Explicit method)

$$y_{n+1} = y_n + h f(t_n, y_n)$$

$$f(s, y(s)) \approx f(t_n, y_n)$$

Trapezoidal Method (implicit method)

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

$$f(s, y(s)) \approx \frac{f(t_n, y_n) + f(t_{n+1}, y_{n+1})}{2}$$

$$\text{i.e. } y_{n+1} - \frac{h}{2} f(t_{n+1}, y_{n+1}) = y_n - \frac{h}{2} f(t_n, y_n)$$

• ORDER OF A NUMERICAL METHOD

Numerical Method given by the recurrence relation $y_{n+1} = F(t, f, y_0, y_1, \dots, y_n, y_{n+1})$.

The method is of order p if

$$y(t_{n+1}) - F(t, f, y(t_0), y(t_1), \dots, y(t_n)) = Ch^p$$

where $y(t)$ is the exact solution to the ODE $y' = f(t, y)$

• TAYLOR SERIES OF y around pt. a is given by

$$y(t) = y(a) + \frac{y'(a)(t-a)}{1!} + \frac{y''(a)(t-a)^2}{2!} + \dots$$

$$+ \frac{y^{(n)}(a)(t-a)^n + y^{(n+1)}(\xi)(t-a)^{n+1}}{n!}$$

for some $\xi \in [a, t]$

$\underbrace{R_n(t)}$

$R_n(t) \equiv$ remainder term

e.g) Euler's Method is of order 1.

e.g) Trapezoidal Method is of order 2.

Note: Method of order p recovers every polynomial solution of degree p or less.

ERRORS IN SOLUTIONS

• CONVERGENT Numerical Method $\Rightarrow \lim_{h \rightarrow 0^+} \max_{n=0,1,\dots,[T/h]} |e_{n,h}| = 0$

• In Euler's Method, local truncation error $\leq \frac{h^2}{2} y''(\xi)$

if $M = \max_{\xi \in [0, T]} |y''(\xi)| \Rightarrow$ local trunc. err. $\leq M \frac{h^2}{2}$

Also, $|e_{n,h}| \leq \frac{Ch}{L} ((1+Lh)^n - 1)$, $n=0,1,\dots,[T/h]$

ONE-STEP METHOD: approximation at n^{th} time step depends on approximation at $(n-1)^{\text{th}}$ time step.

eg: Euler, Trapezoidal

TWO-STEP METHOD: approximation at $(n+1)^{\text{th}}$ time step depends on approximation at n^{th} & $(n-1)^{\text{th}}$ time steps.

$$(y(t_{n+1}) - y(t_{n-1})) = \int_{t_{n-1}}^{t_{n+1}} f(s, y(s)) ds$$

eg: Simpson's $\frac{1}{3}$ rule

Eq) Adams-Basforth Method :

$$y_{n+1} = y_{n+3} + \frac{h}{24} (55f_{n+3} - 59f_{n+2} + 37f_{n+1} - 9f_n)$$

\rightarrow 4-step
 \rightarrow explicit

Eq) Adams-Moulton Method:

$$y_{n+3} = y_{n+2} + \frac{h}{24} (9f_{n+3} + 19f_{n+2} - 5f_{n+1} - 9f_n)$$

\rightarrow 3-step
 \rightarrow implicit.

- General s-STEP method: $\sum_{k=0}^s \alpha_k y_{n+k} = h \sum_{k=0}^s \beta_k f(t_{n+k}, y_{n+k})$

$\beta_s = 0 \Leftrightarrow$ Explicit method

$\beta_s \neq 0 \Leftrightarrow$ Implicit method

• Adam's recipe for s -Step Method

Suppose we have the following N approximations:

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

Interpolate $f(t_i, y_i)$ with $\Psi(t_i)$

$$\Psi(t) = \sum_{k=0}^{s-1} \Psi_k(t) f(t_{n+k}, y_{n+k})$$

$$\Psi_k(t) = \prod_{\substack{l=0 \\ l \neq k}}^{s-1} \left(\frac{t - t_{n+l}}{t_{n+k} - t_{n+l}} \right)$$

$$\text{So } y(t_{n+s}) = y(t_{n+s-1}) + \sum_{k=0}^{s-1} \int_{t_{n+s-1}}^{t_{n+s}} \Psi_k(\tau) f(t_{n+k}, y_{n+k}) d\tau$$

$$y_{n+s} = y_{n+s-1} + h \sum_{k=0}^{s-1} \beta_k f(t_{n+k}, y_{n+k})$$

$$\beta_k = \frac{1}{h} \int_{t_{n+s-1}}^{t_{n+s}} \Psi_k(\tau) d\tau = \frac{1}{h} \int_{t_{n+s-1}}^{t_{n+s}} \Psi_k(t_{n+s-1} + \tau) d\tau.$$

• Adam's recipe for 2-step method:

$$y_{n+2} = y_{n+1} + h \beta_0 f(t_n, y_n) + h \beta_1 f(t_{n+1}, y_{n+1}) \quad \text{for } n=0, 1, \dots$$

$$\beta_0 = \frac{1}{h} \int_{t_n}^{t_{n+1}} \left(\frac{2 - t_{n+1}}{t_{n+1} - t_n} \right) d\tau$$

$$\beta_1 = \frac{1}{h} \int_{t_n}^{t_{n+1}} \left(\frac{2 - t_n}{t_{n+1} - t_n} \right) d\tau$$

- Order of general S-step method:

$$\text{Method is } \sum_{k=0}^s \alpha_k y_{n+k} = h \sum_{k=0}^s \beta_k f(t_{n+k}, y_{n+k}) \text{ for } n=0, 1, \dots$$

Method is of order $p \geq 1$ iff:

$$\sum_{k=0}^s \alpha_k y(t_{n+k}) - h \sum_{k=0}^s \beta_k f(t_{n+k}, y(t_{n+k})) = O(h^{p+1})$$

- Theorem: Order of method is $p \geq 1$ only if: Then we should have

$$\sum_{k=0}^s \alpha_k = 0$$

$$\sum_{k=0}^s k^m \alpha_k = \sum_{k=0}^{m-1} k^{m-1} \beta_k \quad (\text{for } m = 1, \dots, p)$$

$$\sum_{k=0}^s k^{p+1} \alpha_k \neq (p+1) \sum_{k=0}^s k^p \beta_k$$

- Convergence of S-Step method:

Define 1st characteristic polynomial $g(z) = \sum_{k=0}^s \alpha_k z^k$

$$\text{"2nd"} \quad g(z) = \sum_{k=0}^s \beta_k z^k$$

- Dahlquist Equivalence Theorem:

An S-step method is convergent iff

* order of method p is ≥ 1

* roots of $g(z)$ lie in the closed unit disk $\{z \in \mathbb{C} \mid |z| \leq 1\}$

with those lying on $|z|=1$ being SIMPLE.

Runge-Kutta Methods

$$y' = f(t, y(t)), \quad y(0) = y_0$$

Integrating over $(0, h)$ yields

$$y(h) - y_0 = \int_0^h f(s, y(s)) ds$$

Change of variable $s = hT$ yields

$$y(h) - y_0 = h \int_0^1 f(hT, y(hT)) dT$$

We approximate: $\int f(hT, y(hT)) dT \approx \sum_{i=1}^{v-1} b_i f(c_i h, y(c_i h))$
via quadrature formula

We don't know value of $y(c_i h)$.

So approximate: $\xi_i \approx y(c_i h)$ for $i = 1, \dots, v$

For $i = 1$, we take: $c_1 = 0$ i.e. $\xi_1 = y_0$

For $i = 2$, we take: $\xi_2 = \xi_1 + a_{2,1} f(0, \xi_1)$

For $i = 3$, we take: $\xi_3 = y_0 + h a_{3,1} f(0, \xi_1) + h a_{3,2} f(c_2 h, \xi_2)$

$$\xi_v = y_0 + h \sum_{i=1}^{v-1} a_{v,i} f(c_i h, \xi_i)$$

∴ Problem reduces to "v-stage Runge-Kutta Method":

$$y_1 = y_0 + h \sum_{i=1}^{v-1} b_i f(c_i h, \xi_i)$$

$$\text{where } \xi_i = y_0 + h \sum_{j=1}^{i-1} a_{i,j} f(c_j h, \xi_j)$$

Here $A = (a_{i,j})$ is called RK matrix;

$b_i \rightarrow$ RK weights

$c_i \rightarrow$ RK nodes

$$\text{Generalizing: } y_{n+1} = y_n + h \sum_{i=1}^v b_i f(t_n + c_i h, \xi_i)$$

$$\xi_i = y_n + h \sum_{j=1}^{i-1} a_{i,j} f(t_n + c_j h, \xi_j)$$

$a_{i,j}, b_i$ & c_i are same for all subintervals $[t_n, t_{n+1}]$

* Comparing with Taylor Expansion of $y(t_{n+1})$

$$y(t_{n+1}) = y(t_n) + hf(t_n, y_n) + \frac{h^2}{2} \left(\frac{\partial f}{\partial t}(t_n, y_n) + f(t_n, y(t_n)) \right)$$

Example : 2-stage Runge-Kutta Method $\frac{\partial f}{\partial y}(t_n, y(t_n)) + O(h^3)$

$v=2$. We need to find $c_2, b_1, b_2, a_{2,1}$ ($c_1=0$)

We will find these constants imposing the order is ≥ 2 .

$$y_{n+1} = y_n + hb_1 f(t_n, y_n) + hb_2 f(t_n + c_2 h, \xi_2)$$

$$\text{where } \xi_2 = \xi_1 + ha_{2,1} f(t_n, y_n) = y_n + ha_{2,1} f(t_n, y_n)$$

Employing Taylor series in 2-variables :

$$\begin{aligned} f(t_n + c_2 h, \xi_2) &= f(t_n + c_2 h, y_n + ha_{2,1} f(t_n, y_n)) \\ &= f(t_n, y_n) + h \left(c_2 \frac{\partial f}{\partial t}(t_n, y_n) + a_{2,1} f(t_n, y_n) \frac{\partial f}{\partial y}(t_n, y_n) \right) \\ &\quad + O(h^2) \end{aligned}$$

Substituting this we get :

$$\begin{aligned} y_{n+1} &= y_n + hb_1 f(t_n, y_n) + hb_2 f(t_n + c_2 h, \xi_2) \\ &= y_n + hb_1 f(t_n, y_n) + hb^2 \left(f(t_n, y_n) + h \left(c_2 \frac{\partial f}{\partial t}(t_n, y_n) + a_{2,1} f(t_n, y_n) \frac{\partial f}{\partial y}(t_n, y_n) \right) \right) \\ &= y_n + h(b_1 + b_2) f(t_n, y_n) \\ &\quad + h^2 b_2 \left(c_2 \frac{\partial f}{\partial t}(t_n, y_n) + a_{2,1} f(t_n, y_n) \frac{\partial f}{\partial y}(t_n, y_n) \right) \\ &\quad + O(h^3) \end{aligned}$$

* Order 2 $\Rightarrow b_1 + b_2 = 1, a_{2,1} = c_2, b_2 c_2 = \frac{1}{2}$

We can take $b_1 = b_2 = \frac{1}{2}, a_{2,1} = c_2 = \frac{1}{2}$

\therefore Solution:
$$y_{n+1} = y_n + \frac{h}{2} (f(t_n, y_n) + f(t_n + h, y_n + hf(t_n, y_n)))$$

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

2-stage Runge Kutta :

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

4-stage RK :

$$\begin{array}{c|ccc} 0 & & & \\ \hline \frac{1}{2} & \frac{1}{2} & & \\ \frac{1}{2} & & 0 & \frac{1}{2} \\ \hline 1 & 0 & 0 & 1 \\ \hline \frac{1}{6} & \frac{1}{3} & \frac{1}{3} & \frac{1}{6} \end{array}$$

NUMERICAL SOLUTIONS TO SYSTEM OF LINEAR EQUATIONS

$Ax = b$ where $x, b \in \mathbb{R}^n$ and $A \in M_{m,n}(\mathbb{R})$

Undetermined System ($m < n$)

NO solution

Infinite solutions

Overdetermined System ($m > n$)

Solutions Exist (if no. of lin. indep.

equations \leq no. of unknowns)

Square System ($m = n$)

A is invertible: Unique soln.

A is not invertible: no soln.

infinite solns.

{Note: A is invertible iff null space $\text{null space}(A) = \{0\}$ }

- In square system, if A is invertible, unique solution given by $x = A^{-1}b$

- Cramer's Formula: $x_i = \frac{\det(A_i)}{\det(A)}$ for $i = 1, \dots, n$

where A_i is formed by replacing i^{th} column of A by b .

Special Case: A is upper Δ^{Lav} then solve for x by back-substitution.

$$x_i = b_i - \sum_{j=i+1}^n a_{ij} x_j$$

a_{ii}

$$\text{i.e. } x_n = \frac{b_n}{a_{nn}}, \quad x_{n-1} = \frac{b_{n-1} - a_{n-1,n} x_n}{a_{n-1,n}}$$

(Similarly for lower Δ^{Lav} matrices)

Gaussian Elimination is a direct method to solve the system.

→ Transform given system to an upper Δ^{tar} system

By elimination.

i.e. $Ax = b$, we want to find a matrix B such that BA is upper Δ^{tar} then solve the system $BAx = Bb$ by back-substitution.

- An Existence Theorem :

$A \rightarrow n \times n$ matrix

\exists at least one non-singular matrix B such that BA is upper Δ^{tar} .

- Diagonal Submatrix (definition) :

matrix $A = (A_{ij})$.

Δ^k = diagonal submatrix of order k .

$$= \begin{pmatrix} A_{11} & \cdots & A_{1k} \\ \vdots & \ddots & \vdots \\ A_{k1} & \cdots & A_{kk} \end{pmatrix}$$

Theorem :

Let A be a matrix s.t. all submatrices of order Δ^k $\forall k$ from $1, \dots, n$ are invertible, then Gaussian elimination does not need pivoting strategy.

Elimination Algorithm:

$$A_{ij}^{(k+1)} = A_{ij}^{(k)} - m_{ik} A_j^{(k)}$$

$$b_i^{(k+1)} = b_i^{(k)} - m_{ik} b_k^{(k)}$$

• LU factorization.

matrix multiplication

$$A = LU \quad (\text{if } A \text{ is square})$$

where $L_{ij} = \begin{cases} 1, & i=j \\ m_{ij}, & i > j \\ 0, & \text{otherwise} \end{cases}$

$U \equiv A^{(n)}$ (i.e., obtained at the end of Gaussian Elimination)

- When can we do LU factor.?

Ans: When we can do Gaussian Elim. without pivoting.

$A \Leftrightarrow LU \Leftrightarrow \text{Gauss Elim. without pivoting.}$

- LU factorization is unique

- Algorithm:

$$U_{1j} = A_{1j} \quad (j=1 \text{ to } n)$$

$$U_{jj} = A_{jj} - \sum_{k=1}^{j-1} L_{jk} U_{kj}$$

$$L_{j+1,j} = A_{j+1,j} - \sum_{k=1}^{j-1} L_{j+1,k} U_{kj}$$

$$L_{nj} = A_{nj} - \sum_{k=1}^{j-1} L_{nk} U_{kj}$$

$$U_{jj}$$

- Cholesky Factorization

$$A = BB^T \quad \text{where } B \text{ is lower } \Delta^{\text{low}}$$

Applicable for :

- Real entries
- Symmetric A
- Positive definite A . i.e. $\langle Ay, y \rangle > 0 \forall y \neq 0$ (i.e. eigenvalues > 0)

- Uniqueness Theorem (& Existence) :-

$A \rightarrow$ real, symmetric, pos. definite matrix

\exists a unique real lower Δ^{low} matrix B with positive diagonal entries s.t. $A = BB^T$

- Positive-Definiteness test of a matrix A :

Cholesky algorithm gives the required B
iff A is also pos. definite.

- Algorithm:

$$B_{jj} = \sqrt{A_{jj} - \sum_{k=1}^{j-1} B_{jk}^2}$$

$$B_{j+1,j} = A_{j+1,j} - \sum_{k=1}^{j-1} B_{j+1,k} B_{jk}$$

$$B_{n,j} = \frac{A_{nj} - \sum_{k=1}^{j-1} B_{nk} B_{jk}}{B_{jj}}$$

• QR Factorization

$A = QR$ where columns of A are $\{A_1, A_2, \dots, A_n\}$

where Q is orthogonal i.e. $Q^T Q = I$ &
 R is upper diag

$Ax = b$ reduces to $Rx = Q^T b$ which can
 be solved by back-substitution.

- Existence Theorem:

Any square matrix A can be QR factorized.
 if A is invertible then R is invertible.

- Uniqueness: If R is assumed to have positive diagonal entries then the pair (Q, R) is unique.

- Algorithm: (for invertible A)

Denote columns of A by A_1, A_2, \dots, A_n .

Column vectors are linearly indep.

Build an orthonormal basis using Gram-Schmidt.

$$q_1 = \frac{A_1}{\|A_1\|}$$

$$q_2 = \frac{A_2 - \langle q_1, A_2 \rangle q_1}{\| \cdot \|}$$

$$q_i = \frac{A_i - \sum_{k=1}^{i-1} \langle q_k, A_i \rangle q_k}{\| \cdot \|} \quad \text{for } 1 \leq i \leq n$$

Now make Q such that its i^{th} column is q_i .

To make R :

$$R_{ki} = \begin{cases} \langle q_k, A_i \rangle, & 1 \leq k \leq i-1 \\ \|A_i - \sum_{k=1}^{i-1} \langle q_k, A_i \rangle q_k\|, & k=i \\ 0, & k > i \end{cases}$$

NORMS

We use norms to estimate the error between true solution (x) and approximate solution (\tilde{x}) to $Ax = b$

$$e = \|x - \tilde{x}\|$$

is a measure of error
and bigger in \mathbb{R}^n

- Norms on \mathbb{R}^n :

Norm is a fn. $\|\cdot\| : \mathbb{R}^n \rightarrow [0, \infty)$ such that

$$\|x\| > 0 \quad \forall x \neq 0$$

$$\|\alpha x\| = |\alpha| \|x\| \quad \forall \alpha \in \mathbb{R}, x \in \mathbb{R}^n$$

$$\|x + y\| \leq \|x\| + \|y\|$$

- ℓ^p -norm for $p \in [1, \infty)$: (norm is on \mathbb{R}^n)

$$\|x\|_{\ell^p} = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}$$

- ℓ^∞ -norm:

$$\|x\|_{\ell^\infty} = \max_{i \in n} |x_i|$$

- Equivalent norms:

Two norms $\|\cdot\|$ and $\|\cdot\|'$ are equivalent if

$\exists m > 0$ and $\exists M > 0$ such that

$$m\|x\| \leq \|x\|' \leq M\|x\| \quad \forall x \in \mathbb{R}^n$$

- $\|\cdot\|_{\ell^\infty}$ and $\|\cdot\|_{\ell^p}$ norms are equivalent. $\therefore \|\cdot\|_{\ell^p}$

$$\|x\|_{\ell^\infty} \leq \|x\|_{\ell^p} \leq \sqrt[p]{n} \|x\|_{\ell^\infty}$$

$\|\cdot\|_{\ell^p}$ and $\|\cdot\|_{\ell^\infty}$ are equivalent

- $\|\cdot\|_{\ell^1}$ and $\|\cdot\|_{\ell^2}$ are equivalent:

$$\|x\|_{\ell^2} \leq \|x\|_{\ell^1} \leq \sqrt{n} \|x\|_{\ell^2}$$

- All norms on \mathbb{R}^n are equivalent!

• Norms on $M_{m,n}(\mathbb{R})$

- ℓ^p -norms:

Let $A \in M_{m,n}(\mathbb{R})$.

$$\|A\|_{\ell^p} = \left(\sum_{j=1}^n \left(\sum_{i=1}^m |A_{ij}|^p \right)^{1/p} \right) \text{ for } p \in [1, \infty)$$

$$\|A\|_{\ell^\infty} = \max_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n}} |A_{ij}| \text{ of max. of all}$$

- $\ell_{p,q}$ -norms: $\|A\|_{\ell_{p,q}}$

$$\|A\|_{\ell_{p,q}} = \left(\sum_{j=1}^n \left(\sum_{i=1}^m |A_{ij}|^p \right)^{q/p} \right)^{1/q} \text{ for } p, q \in [1, \infty)$$

→ Collect ℓ^p -norm of each column vector of size m .

→ Take ℓ^q -norm of this collection.

$$\|A\|_{\ell_{p,p}} = \|A\|_{\ell^p} \text{ for more standard}$$

$\|A\|_{\ell^2}$ is called Frobenius norm.

$$\|AB\|_{\ell^2} \leq \|A\|_{\ell^2} \|B\|_{\ell^2}$$

• Matrix Norms when $\|\cdot\|$ is

A norm $\|\cdot\|$ defined on $M_n(\mathbb{R})$ is Matrix norm

if $\forall A, B \in M_n(\mathbb{R})$

$$\|AB\| \leq \|A\| \|B\|$$

$$\|I - AB\| = \|BA\|$$

$$\|I_n\| \geq 1 \quad (\because \|A\| = \|AI_n\| \leq \|A\| \|I_n\|)$$

Frobenius norm is matrix norm.

$\|\cdot\|_\infty$ is not matrix norm. (Take A to have all ones $\therefore A^2$ has all n)

• Subordinate Matrix Norm:

Let $\|\cdot\|$ be a vector norm on \mathbb{R}^n . It induces a norm on M_n defined by

$$\|A\| = \sup_{x \in \mathbb{R}^n \setminus \{0\}} \frac{\|Ax\|}{\|x\|}$$

This is referred to as norm on M_n subordinate to the vector norm $\|\cdot\|$.

- Can also be written as $\|A\| = \sup_{y \in \mathbb{R}^n, \|y\|=1} \|Ay\|$
- $\exists z \in \mathbb{R}^n \setminus \{0\}$ such that $\|A\| = \frac{\|Az\|}{\|z\|}$
- Subordinate norm is a matrix norm, i.e. $\|AB\| \leq \|A\| \|B\|$
- Subordinate matrix norm of I_n is 1 ($\because \|I_n\| = \sup_{\|x\|} \frac{\|Ix\|}{\|x\|} = 1$)
- Frobenius norm of identity $= \sqrt{n} \therefore$ it is not subordinate to any vector norm.
- Let $\|\cdot\|_p$ denote the matrix norm subordinate to any vector norm $\|\cdot\|_p$. $\|A\|_p = \sup_{x \in \mathbb{R}^n \setminus \{0\}} \frac{\|Ax\|_p}{\|x\|_p}$
- For a unitary matrix A (*i.e.* $A^{-1} = A^T$), $\|A\|_2 = 1$
- For any unitary matrix A and $B \in M_n$, we have $\|AB\|_2 = \|BA\|_2 = \|B\|_2$
for $\|\cdot\|_2$

$\therefore \|\cdot\|_2$ -norm is invariant under multiplication by unitary matrices. theoretical

{ $\|\cdot\|_2$ -norm is useful for computations, because of inner product structure }

- A is a diagonal matrix: $A = \text{diag}(a_1, a_2, \dots, a_n)$
 $\|A\|_2 = \max_{1 \leq i \leq n} |a_i|$

NORMAL MATRICES

- Adjoint A^* of a matrix $A \in M(C)$ is defined as

$$A_{ij}^* = \overline{A_{ji}} = \overline{A_{ji}}$$

- A is normal if $A A^* = A^* A$

- Theorem (Characterization of Normal Matrices):

Let $S \subset M(C)$ denote set of all normal matrices.
Define set $T = \{A \in M(C) \text{ s.t. } A = U \text{diag}(\lambda_1, \dots, \lambda_n) U^*\}$
with U unitary (i.e. $U^{-1} = U^*$) and λ_i eigenvalues of A .
Then $S = T$.

- For a normal matrix A ,

$$\|A\|_2 = \|U \text{diag}(\lambda_1, \dots, \lambda_n) U^*\|_2 = \|\text{diag}(\lambda_1, \dots, \lambda_n)\|_2$$

$\|A\|_2 = \max_{1 \leq i \leq n} |\lambda_i|$

SPECTRAL RADIUS AND MATRIX NORMS

- Spectral Radius:

Maximum of moduli of eigen values of a matrix A is spectral radius of A : $\rho(A) = \max_{1 \leq i \leq n} |\lambda_i|$

$$\|A\|_2 = \rho(A) \quad \{A \text{ is normal}\}$$

ρ cannot be a norm $\because \rho = 0 \Rightarrow A = 0$

For any "matrix norm" $\|\cdot\|$ on $M_n(C)$, $\rho(A) \leq \|A\|$

For a given $A \in M_n$ and any $\epsilon > 0$, \exists a subordinate matrix norm $\|\cdot\|$ such that

$$\rho(A) < \|A\| \leq \rho(A) + \epsilon.$$

(The matrix norm depends on (A, ϵ))

CALCULATING MATRIX NORMS

- Practical recipe for computing $\|A\|_1$ of a given matrix :

$$\|A\|_1 = \max_{1 \leq j \leq n} \left(\sum_{i=1}^n |A_{ij}| \right)$$

(Hence $\|A\|_1$ is also called column-sum norm)

- $\|A\|_\infty = \max_{1 \leq i \leq n} \left(\sum_{j=1}^n |A_{ij}| \right)$ (Row-sum norm)

- Practical upper bound for spectral radius:

$$\rho(A) \leq \min(\text{column-sum norm}, \text{row-sum norm})$$

$$\|A\|_1 = \|U^* U(\lambda_1, \dots, \lambda_n) V\|_1 \leq \|U\| \cdot \|V\| \cdot \|\Lambda\|_1$$

- HERMITIAN MATRIX : $A = A^*$

\therefore Eigenvalue is always real.

- Normal matrix

- Proposition: $A \in M_n(\mathbb{C})$ is Hermitian iff \exists a unitary matrix U and n real eigenvalues of A : $\lambda_1, \lambda_2, \dots, \lambda_n$ such that

$$A = U \text{diag}(\lambda_1, \dots, \lambda_n) U^*$$

- For any $A \in M_n(\mathbb{C})$, $A^* A$ is Hermitian

& All eigenvalues of $A^* A$ are real and non-negative.

- SINGULAR VALUES: of a matrix $A \in M_n(\mathbb{C})$ are the non-negative square roots of the eigenvalues of $A^* A$ {by convention, we do not mention '0' as a sing. val even if it is}.

- Singular values of a normal matrix are the moduli of its eigenvalues (\because a normal matrix A can be written as $UDU^* = A$, where $D = \text{diag}(\lambda_1, \dots, \lambda_n)$)

- $\|A\|_2 = \max_{1 \leq i \leq n} \sqrt{\lambda_i} = \max(\text{eigenval. of } A^* A) = \text{largest singular val. of } A$

MATRIX SEQUENCES

Convergence:

A sequence $\{A^{(k)}\}$ of matrices converges to a limit A if $\lim_{k \rightarrow \infty} \|A^{(k)} - A\| = 0$ for some norm on $M_n(\mathbb{C})$.

It is also denoted by $\lim_{k \rightarrow \infty} A^{(k)} = A$

- Choice of norm is irrelevant (here) because all norms are equivalent in finite dimensional space.

• Sequence of iterated powers: $A^{(k)} = A^k$

The following are equivalent:-

i) $\lim_{k \rightarrow \infty} A^k = 0 \rightarrow$ zero matrix

ii) $\lim_{k \rightarrow \infty} A^k x = 0 \quad \forall x \in \mathbb{C}^n$

iii) $\rho(A) < 1$

iv) \exists at least one subordinate matrix norm

$\|\cdot\|$ such that $\|A\| < 1$

• Geometric matrix series:

$$S^{(0)} = I_n, \quad S^{(k)} = I_n + A + \dots + A^k$$

If $\{S^{(k)}\}$ converges, we get geometric series $\sum A^k = \lim_{k \rightarrow \infty} S^{(k)}$

- Geometric series converges if $\rho(A) < 1$.

If $\rho(A) < 1$, then $I_n - A$ is invertible

$$(I_n - A)^{-1} = \sum_{k=0}^{\infty} A^k$$

- Any invertible matrix has a neighbourhood in which there are invertible matrices.

Take an invertible matrix $A \in M_n(\mathbb{C})$. Take a subordinate matrix norm $\|\cdot\|$. Then the matrix B , s.t. $\|A - B\| < \frac{1}{\|A^{-1}\|}$, is invertible.

* - $\text{cond}(A) = \|A\| \|A^{-1}\| \geq \rho(A) \rho(A^{-1})$

The lower bound is attained if A is normal

$$\text{cond}_2(A) = \rho(A) \rho(A^{-1}) = |\lambda_{\max}| / |\lambda_{\min}|$$

Page No. _____ Date: _____

$$\|A\|_2 = \rho(A)$$

(32)

Relative Error: Let $y \in \mathbb{C}^n$ and $\tilde{y} \in \mathbb{C}^n$ be its perturbed value. Relative error wrt norm $\|\cdot\|$ on \mathbb{C}^n is defined as $\frac{\|\tilde{y} - y\|}{\|y\|}$

Condition Number: of an invertible matrix $A \in M_n(\mathbb{C})$ relative to a subordinate matrix norm $\|\cdot\|$ is

$$\text{cond}(A) = \|A\| \|A^{-1}\|$$

- $\text{cond}_p(A) = \|A\|_p \|A^{-1}\|_p$

- $\text{cond}(I_n) = 1$

- $\text{cond}(A) \geq 1 \quad (\because \|AA^{-1}\| \leq \|A\| \|A^{-1}\|)$

- $\text{cond}(A) = \text{cond}(A^{-1})$

- $\text{cond}(\alpha A) = \text{cond}(A) \quad \forall \alpha \neq 0$

- If A is invertible and $A + B$ is invertible, what is $(A + B)^{-1}$: $(A + B)^{-1} = (I_n + A^{-1}B)^{-1} A^{-1}$

- Condition Number and Amplification factor:

$$\frac{\|x_e - x\|}{\|x\|} \leq \|A\| \|A^{-1}\| \left\{ \frac{\|b_e - b\|}{\|b\|} + \frac{\|A_e - A\|}{\|A\|} \right\} + O(\epsilon^2)$$

→ Estimating relative errors in solution by relative errors in data.

→ $\text{cond}(A) = \|A\| \|A^{-1}\|$ acts as amplifying factor.

- Well-conditioned Matrix: $\text{cond}(A) \approx 1$.

Ill-conditioned " : $\text{cond}(A) \gg 1$

Well-conditioned Matrices: If A is unitary ($A^{-1} = A^*$), $\|A\|_2 = \|A^*\|_2 = 1 \therefore \text{cond}_2(A) = 1$

largest sing. val. of A

$\text{cond}_2(\text{rotation matrix}) = 1$

* - $\|A\|_2 = \mu_{\max}$ (We know this). $\text{cond}_2(A) = \|A\|_2 \|A^{-1}\|_2 = \mu_{\max} / \mu_{\min}$

PAGE No.	
DATE	/ /

ITERATIVE METHODS FOR LINEAR SYSTEMS $Ax = b$

$x^{(0)} \rightarrow x^{(1)} \dots \rightarrow \infty$ such that sequence converges to desired soln. x . i.e. $\lim_{k \rightarrow \infty} \|x^{(k)} - x\| = 0$ for some norm on \mathbb{R}^n

- error @ k^{th} iteration : $e^{(k)} = x^{(k)} - x$ for $k=0,1,2,\dots$
- residual @ k^{th} iteration : $r^{(k)} = Ax^{(k)} - b$ for $k=0,1,\dots$

- Convergent iterative method: One that converges to exact solution for any choice of initial vector.

$$\lim_{k \rightarrow \infty} x^{(k)} = x \iff \lim_{k \rightarrow \infty} e^{(k)} = 0 \iff \lim_{k \rightarrow \infty} r^{(k)} = 0$$

- Splitting of A : $A = M_1 - N_1$, M_1 is invertible
 (M_1, N_1) is called splitting

- Method: Initialize with $x^{(0)} \in \mathbb{R}^n$.

Compute $x^{(k+1)} = M_1^{-1}N_1x^{(k)} + M_1^{-1}b$, $k=0,1,\dots$
 \Rightarrow yields $M_1x = N_1x + b$ i.e. $Ax = b$
 $\{M_1^{-1}N_1\}$ is called iteration matrix

- Convergence Theorem: The above method converges iff
 $\rho(M_1^{-1}N_1) < 1$

- Error: $e^{(k)} = (M_1^{-1}N_1)^k e^{(0)}$

Rate of Convergence Theorem: $\|e^{(k)}\| \sim \|e^{(0)}\| [\rho(M_1^{-1}N_1)]^k$
 (any arbitrary norm on \mathbb{R}^n) for $k \geq 1$

\therefore smaller $\rho \Rightarrow$ better convergence.

PAGE NO.	
DATE	/ /

Richardson's Iterative Method :

$$\text{Take } M = \frac{1}{\alpha} I_n, N = \frac{1}{\alpha} I_n - A, \alpha \neq 0$$

$$\text{Then } M^{-1} N = I_n - \alpha A$$

$$x^{(k+1)} = x^{(k)} + \alpha (b - Ax^{(k)})$$

- When does the method converge?

$M^{-1}N$ has eigenvalues $(1 - \alpha \lambda_i)$ where λ_i are eigen values of A .
 \therefore method converges iff $|1 - \alpha \lambda_i| < 1$ & eigenvals. of A λ_i .

- Real symmetric A with real eigenvalues :

$$\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$$

- if $\lambda_1, \lambda_n < 0 \Rightarrow$ method doesn't converge for any α
- if $\lambda_1 > 0 \neq \lambda_n$ then " converges iff $0 < \alpha < 2$
- if $\lambda_1 < 0 \neq \lambda_n$ $2 < \alpha < \lambda_n$

- Optimal α for A with $0 < \lambda_1 \leq \dots \leq \lambda_n$ (+ve. eigen vals.)

$$\alpha = \frac{2}{\lambda_1 + \lambda_n}$$

$$\rho(I_n - \tilde{\alpha} A) = \frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1}$$

$$\begin{aligned} \text{if further } A \text{ is normal, } \rho(I_n - \tilde{\alpha} A) \\ = \frac{\lambda_n/\lambda_1 - 1}{\lambda_n/\lambda_1 + 1} = \frac{\text{cond}_2(A) - 1}{\text{cond}_2(A) + 1} \rightarrow \text{increasing fn. of } \text{cond}_2. \end{aligned}$$

\therefore More well-conditioned the matrix, faster the convergence of Richardson method.

- Jacobi Method :

$$A = M - N$$

$M = D = \text{diag}(A_{11}, \dots, A_{nn})$, $A_{ii} \neq 0$, $A \neq 0$

$N = D - A$ $\therefore M$ is invertible

$$\therefore M^{-1}N = I_n - D^{-1}A$$

$$x^{(k+1)} = p(x^{(k)}) + D^{-1}(b - Ax^{(k)})$$

- For what matrices A does the method converge?

- if A is "strictly row-diagonally dominant" or "strictly column-diagonally dominant".

(In either of these cases, $\rho(M^{-1}N) < 1$)

- Defn: Strictly row diagonally dominant matrix A if $|A_{ii}| > \sum_{\substack{k=1 \\ k \neq i}} |A_{ik}|$

Strictly column diagonally dominant if $|A_{ii}| > \sum_{\substack{k=1 \\ k \neq i}} |A_{ki}|$

$$D = (D_{ij}) \quad D_{ii} = (-1)^i b_i$$

PAGE No.	
DATE	/ /

Gauss-Seidel Method :

$$A = M - N$$

$$M = D - E ; \quad N = F$$

$$D = \text{diag}(A_{11}, \dots, A_{nn}), \quad A_{ii} \neq 0$$

$$E_{ij} = \begin{cases} -A_{ij}, & i > j \\ 0, & \text{otherwise} \end{cases}; \quad F_{ij} = \begin{cases} -A_{ij}, & i < j \\ 0, & \text{otherwise} \end{cases}$$

strictly lower Δ^{low} strictly upper Δ^{high}

$$M^{-1}N = G = (I_n - L)^{-1}U$$

$$\text{where } L = D^{-1}E \quad U = D^{-1}F$$

For Jacobi method,

$J = \text{iteration matrix}$

$$J = L + U$$

- When does the method converge?

If A is "strictly row diagonally-dominant".

- Which method performs better and how?

Stein-Rosenburg Theorem: If J is non-negative, then one and only one of these hold:

- $\rho(G) = \rho(J) = 0$
- $0 < \rho(G) < \rho(J) < 1$
- $\rho(G) = \rho(J) = 1$
- $\rho(G) > \rho(J) > 1$

i.e. Either Jacobi & Gauss-Seidel converge or both diverge.

When both converge, Gauss-Seidel outperforms Jacobi.

• Relaxed Gauss-Seidel Method:

$$A = M - N$$

$$M = \frac{1}{\kappa} D - E \quad ; \quad N = \left(\frac{1}{\kappa} - 1 \right) D + F$$

{ E and F same as Gauss-Seidel }

$$M^{-1}N = \left(\frac{1}{\kappa} D - E \right)^{-1} \left(\left(\frac{1-\kappa}{\kappa} \right) D + F \right)$$

$$\text{i.e. } G_\kappa = \left(\frac{1}{\kappa} I_n - \frac{1}{\kappa} L \right)^{-1} \left(\left(\frac{1-\kappa}{\kappa} \right) I_n + U \right)$$

- Thm: For any A , we have $\rho(G_\kappa) \geq |\kappa - 1|$

So what conditions for convergence?

convergence $\Rightarrow 0 < \kappa < 2$

$\kappa \in \mathbb{R} \setminus (0, 2) \Rightarrow$ no convergence

- If A is symmetric & positive definite, $A \in M_n(\mathbb{R})$
then for any $\kappa \in (0, 2)$: $\rho(G_\kappa) < 1$

• Some other lemmas for positive definiteness & iteration matrices.

- A be Hermitian, positive definite

(M, N) be splitting of A .

Suppose $M + M^* - A$ is positive definite.

Eigenvalue \rightarrow i) Let λ be eigenvalue of $H \equiv A^{-1}(2M - A)$.

Then $\operatorname{Real}(\lambda) > 0$

Spectral radius \rightarrow ii) Take $H \equiv A^{-1}(2M - A)$. Then

$$(H - I_n)(H + I_n)^{-1} = M^{-1}N$$

And,

$$\rho(M^{-1}N) < 1$$

PAGE NO.	
DATE	/ /

EIGEN-VALUES OF A MATRIX

eigen pair (λ, α) where $\alpha \neq 0$
associated with A if $A\alpha = \lambda\alpha$

- How to find eigen vals.?

$$|A - \lambda I_n| = 0$$

- characteristic polynomial is of degree n .

can have at most n roots, along with repeated roots if any. i.e. $\lambda_i = \lambda_j$ for some $i \neq j$.

$$\det(A) = \lambda_1 \lambda_2 \dots \lambda_n$$

$$\text{Tr}(A) = \lambda_1 + \dots + \lambda_n$$

- Eigen values of A^{inv} matrix are its diagonal values.
- "real symm. matrix" are real.
- orthogonal matrix satisfy $|A| = 1$.

- " A and B are similar" means $B = P^{-1}AP$ for some invertible P .

$$\max_i |\lambda_i| = \|A\| \leq \sqrt{\text{any matrix norm}}$$

i.e. all eigen vals. lie in the ball centred around origin in C with radius $\|A\|$.

- Gershgorin Circle Theorem:

Any eigenvalue λ of matrix A belongs to the set $\bigcup_{i=1}^n D_i$.

$$\text{where } D_i = \{z \in C \mid |z - A_{ii}| \leq R_i\}$$

$$\text{and } R_i = \sum_{j \neq i} |A_{ij}|$$

ROOTS OF NONLINEAR EQUATIONS

Intermediate Value Theorem

f be continuous $[a, b] \rightarrow \mathbb{R}$.

If $f(a)f(b) \leq 0$ then $\exists \xi \in [a, b]$ s.t. $f(\xi) = 0$.

- Fixed pt: $y \in \mathbb{R}$ is a fixed pt. of $\Psi: \mathbb{R} \rightarrow \mathbb{R}$ if $\Psi(y) = y$

Fixed point theorem

g be continuous fn. $g: [a, b] \rightarrow [a, b]$.

Then $\exists \xi \in [a, b]$ s.t. $\xi = g(\xi)$

- Finding roots of $f(x) = 0$ is equivalent to

solving $x - g(x) = 0$ for some cont. fn. g .

$$\text{Then } f(\xi) = 0 \Leftrightarrow \xi = g(\xi)$$

- Iterative process to get fixed pt. of continuous fn. g :
($g: [a, b] \rightarrow [a, b]$)

Idea is to approximate fixed pt. ξ by sequence $x^{(k)}$.

$x^{(0)} \in [a, b]$ = initial guess.

Further iterates computed using $x^{(k+1)} = g(x^{(k)})$

Suppose the sequence converges, then to show

that it converges to the fixed pt. of g :

$$\xi = \lim_{k \rightarrow \infty} x^{(k+1)} = \lim_{k \rightarrow \infty} g(x^{(k)}) = g\left(\lim_{k \rightarrow \infty} x^{(k)}\right) = g(\xi)$$

$\therefore \xi$ must be fixed pt. of g .

For what functions g does the sequence of iterates converge?

- contraction mapping theorem:

If $g: [a, b] \rightarrow [a, b]$ is a contraction on $[a, b]$, then g has a "unique" fixed pt. $\xi \in [a, b]$.

Also, Sequence of iterates given by

$$x^{(k)} = g(x^{(k-1)}), k = 1, 2, \dots$$

converges to fixed pt. ξ for any initial guess $x^{(0)} \in [a, b]$

• Contraction: let g be cont. fn. $g: [a, b] \rightarrow \mathbb{R}$.

g is called a contraction if \exists a constant $L \in (0, 1)$ such that

$$|g(x) - g(y)| \leq L|x-y| \quad \forall x, y \in [a, b]$$

Can we lift the condition of contraction on entire $[a, b]$?

- Theorem:

continuous fn. $g: [a, b] \rightarrow [a, b]$

Let $\xi \in [a, b]$ be fixed pt. of g .

Assume g' is continuous in $[\xi-h, \xi+h]$ for some $h > 0$. Assume $|g'(\xi)| < 1$.

Then the sequence of iterates $x^{(n+1)} = g(x^{(n)})$

converges to ξ provided initial guess $x^{(0)}$ is sufficiently close to ξ .

On the other hand, in above thm. if we assume $|g'(\xi)| > 1$, then sequence of iterates given by $x^{(n+1)} = g(x^{(n)})$ does not converge to ξ for whatever initial guess $x^{(0)} \neq \xi$.

* Conclusion: continuous fn. $f: [a, b] \rightarrow \mathbb{R}$

Simple iteration $x^{(k+1)} = f(x^{(k)})$

Relaxation iteration $x^{(k+1)} = x^{(k)} - \lambda f(x^{(k)})$

Page Fixed
Date:

41
Rouva

How to approximate zeros of function f ?

• Relaxation iteration:

Continuous fn. $f: [a, b] \rightarrow \mathbb{R}$.

initial guess $= x^{(0)} \in [a, b]$

Iteration $x^{(k+1)} = x^{(k)} - \lambda f(x^{(k)})$, $k=0, 1, 2, \dots$

$\lambda \neq 0$ is some real parameter.

- How do we know the above sequence converges to root of f ?

Suppose the sequence converges to ξ .

$$\xi = \lim_{k \rightarrow \infty} x^{(k+1)} = \lim_{k \rightarrow \infty} x^{(k)} - \lambda \lim_{k \rightarrow \infty} f(x^{(k)}) = \xi - \lambda \lim_{k \rightarrow \infty} f(x^{(k)})$$

$$\Rightarrow f(\xi) = 0$$

Relaxation iteration to find zeros of $f \Leftrightarrow$

Simple iteration to find fixed pt. of $g(x) = x - \lambda f(x)$

- What should value of λ be to guarantee convergence?

Suppose f is differentiable about ξ .

$\therefore g$ is also diff. around ξ . $g' = 1 - \lambda f'$

$$|g'(\xi)| = |1 - \lambda f'(\xi)| < 1 \text{ if :}$$

• λ has same sign as $f'(\xi)$.

• λ is not too large.

Theorem: cont. fn. $f: [a, b] \rightarrow [a, b]$

Let $\xi \in [a, b]$ be zero of f , i.e. $f(\xi) = 0$.

Assume f' is continuous around ξ and let $f'(\xi) \neq 0$.

Then $\exists \lambda > 0, \delta > 0$ s.t. sequence of iterates given by $x^{(k+1)} = x^{(k)} - \lambda f(x^{(k)})$ converges to ξ for whatever initial guess $x^{(0)} \in [\xi - \delta, \xi + \delta]$

* • Bisection Method: Suppose $a_k & b_k$ s.t. $f(a_k)f(b_k) < 0$. Define 42

$$c_k = \frac{a_k + b_k}{2} \quad \text{if } f(c_k) \neq 0$$

• Secant Method: Replace $f'(x_k)$ in Newton's meth. by $f'(x_k) \approx \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}$
Gives $e^{(k+1)} \leq \frac{2}{3} e^{(k)}$

• Error sequence (for a sequence of iterates) is

for any sequence of iterates $x^{(k)}$ trying to approximate a point ξ , the error sequence is
 $e^{(k)} = x^{(k)} - \xi$ for $k = 0, 1, 2, \dots$

• Order of convergence of sequence of iterates is $p (\geq 1)$
if \exists constant $C \geq 0$ (with $C < 1$ if $p = 1$)
and integer N s.t. $\forall k \geq N$:
 $|e^{(k+1)}| \leq C |e^{(k)}|^p$

$p = 1 \Rightarrow$ linear convergence

$p = 2 \Rightarrow$ quadratic convergence

- What is the order of simple iteration $x^{(k+1)} = g(x^{(k)})$?

$$|e^{(k+1)}| = |x^{(k+1)} - \xi| = |g(x^{(k)}) - g(\xi)| \leq L |e^{(k)}|$$

linear convergence

• Newton-Raphson Method

It is a relaxed iteration with best convergence.

$$g(x) = x - \frac{f(x)}{f'(x)}$$

- When does this method converge?

Thm: let $f: [a, b] \rightarrow \mathbb{R}$ be C^2 on I_8 on $[\xi - \delta, \xi + \delta]$ for some $\delta > 0$. Suppose $f(\xi) = 0$, $f'(\xi) \neq 0$, $f''(\xi) \neq 0$.

Suppose further \exists constant M s.t. $|f''(x)| \leq M \forall x, y \in I_8$

If initial guess $x^{(0)}$ is s.t. $|x^{(0)} - \xi| \leq h$ with

$h = \min\left\{\delta, \frac{1}{M}\right\}$, then iteration sequence converges quadratically to ξ . $\lim_{k \rightarrow \infty} \frac{e^{(k+1)}}{(e^{(k)})^2} \leq \frac{M}{2}$