

**A FIRST COURSE**  
**IN**  
**NUMERICAL ANALYSIS**



---

**A FIRST COURSE**  
**IN**  
**NUMERICAL ANALYSIS**  
**MAT4001 Notebook**

---

**Prof. Yutian Li**

*The Chinese University of Hong Kong, Shenzhen*



香港中文大學(深圳)

The Chinese University of Hong Kong, Shenzhen



# Contents

Acknowledgments	vii
Notations	ix
<b>1      Week1 .....</b>	<b>1</b>
1.1    Differentiation	1



# Acknowledgments

This book is from the MAT4001 in fall semester, 2018.

CUHK(SZ)





# Notations and Conventions

$\mathbb{R}^n$	$n$ -dimensional real space
$\mathbb{C}^n$	$n$ -dimensional complex space
$\mathbb{R}^{m \times n}$	set of all $m \times n$ real-valued matrices
$\mathbb{C}^{m \times n}$	set of all $m \times n$ complex-valued matrices
$x_i$	$i$ th entry of column vector $\mathbf{x}$
$a_{ij}$	$(i, j)$ th entry of matrix $\mathbf{A}$
$\mathbf{a}_i$	$i$ th column of matrix $\mathbf{A}$
$\mathbf{a}_i^T$	$i$ th row of matrix $\mathbf{A}$
$\mathbb{S}^n$	set of all $n \times n$ real symmetric matrices, i.e., $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $a_{ij} = a_{ji}$ for all $i, j$
$\mathbb{H}^n$	set of all $n \times n$ complex Hermitian matrices, i.e., $\mathbf{A} \in \mathbb{C}^{n \times n}$ and $\bar{a}_{ij} = a_{ji}$ for all $i, j$
$\mathbf{A}^T$	transpose of $\mathbf{A}$ , i.e., $\mathbf{B} = \mathbf{A}^T$ means $b_{ji} = a_{ij}$ for all $i, j$
$\mathbf{A}^H$	Hermitian transpose of $\mathbf{A}$ , i.e., $\mathbf{B} = \mathbf{A}^H$ means $b_{ji} = \bar{a}_{ij}$ for all $i, j$
$\text{trace}(\mathbf{A})$	sum of diagonal entries of square matrix $\mathbf{A}$
$\mathbf{1}$	A vector with all 1 entries
$\mathbf{0}$	either a vector of all zeros, or a matrix of all zeros
$\mathbf{e}_i$	a unit vector with the nonzero element at the $i$ th entry
$\mathcal{C}(\mathbf{A})$	the column space of $\mathbf{A}$
$\mathcal{R}(\mathbf{A})$	the row space of $\mathbf{A}$
$\mathcal{N}(\mathbf{A})$	the null space of $\mathbf{A}$
$\text{Proj}_{\mathcal{M}}(\mathbf{A})$	the projection of $\mathbf{A}$ onto the set $\mathcal{M}$



# Chapter 1

## Week1

### 1.1. Differentiation

**Definition 1.1** [Forward and Backward Difference Formula]

$$f'(x_0) = \frac{f(x_0 + h) - f(x_0)}{h} - \frac{h}{2}f''(\xi)$$

Consider the Taylor expansion, we have

$$f''(x_0) = \frac{1}{h^2} [f(x_0 - h) - 2f(x_0) + f(x_0 + h)] - \frac{h^2}{12}f^{(4)}(\xi)$$

**Definition 1.2** [General Derivative Approximation] The interpolation formula gives

$$f(x) = \sum_{k=0}^n f(x_k)L_k(x) + \frac{(x - x_0) \cdots (x - x_n)}{(n + 1)!}f^{(n+1)}(\xi_x)$$

Differentiating both sides gives ( $x_j$  is one of the node points)

$$f'(x_j) = \sum_{k=0}^n f(x_k)L'_k(x_j) + \frac{1}{(n + 1)!}f^{(n+1)}(\xi_{x_j}) \prod_{k=0, k \neq j}^n (x_j - x_k)$$

**R** Three-point formula:

$$f'(x_j) = f(x_0) \left[ \frac{2x_j - x_1 - x_2}{(x_0 - x_1)(x_0 - x_2)} \right] + f(x_1) \left[ \frac{2x_j - x_0 - x_2}{(x_1 - x_0)(x_1 - x_2)} \right] \\ + f(x_2) \left[ \frac{2x_j - x_0 - x_1}{(x_2 - x_0)(x_2 - x_1)} \right] + \frac{1}{6} f^{(3)}(\xi_j) \prod_{k=0, k \neq j}^2 (x_j - x_k)$$

Substituting  $x_j = x_0, x_1, x_2$ , we obtain:

$$f'(x_0) = \frac{1}{2h} [-3f(x_0) + 4f(x_0 + h) - f(x_0 + 2h)] + \frac{h^2}{3} f^{(3)}(\xi_0) \\ f'(x_0) = \frac{1}{2h} [-f(x_0 - h) + f(x_0 + h)] - \frac{h^2}{6} f^{(3)}(\xi_1)$$

**Definition 1.3** [Richardson Extrapolation] For the approximation with the form

$$M = N_1(h) + K_1 h^2 + K_2 h^4 + \dots$$

we have the approximation

$$N_j(h) = N_{j-1}\left(\frac{h}{2}\right) + \frac{N_{j-1}(h/2) - N_{j-1}(h)}{4^{j-1} - 1}$$

where  $N_j(h)$  has order  $O(h^{2j})$  ■

round-off error in  $N_1(h/2^k)$ ; we recommend comparing the final diagonal entries to ensure accuracy.

**Definition 1.4** [Quadrature formula]

$$\int_a^b f(x) dx \approx \sum_{i=0}^n a_i f(x_i)$$

where

$$a_i = \int_a^b L_i(x) dx$$

$$E(f) = \frac{1}{(n+1)!} \int_a^b \prod_{i=0}^n (x - x_i) f^{(n+1)}(\xi_x) dx$$

1. Trapezoidal Rule:

$$\int_a^b f(x) dx = \frac{h}{2} [f(x_0) + f(x_1)] - \frac{h^3}{12} f'(\xi)$$

Seperating into subintervals  $[x_{k-1}, x_k]$ , we have, with  $h = (b - a)/n$ ,

$$\int_a^b f(x) dx = \frac{h}{2} \left[ f(a) + 2 \sum_{j=1}^{n-1} f(x_j) + f(b) \right] - \frac{b-a}{12} h^2 f''(\mu)$$

2. Simpson's rule: for  $h = (x_2 - x_0)/2$ ,

$$\int_{x_0}^{x_2} f(x) dx = \frac{h}{3} [f(x_0) + 4f(x_1) + f(x_2)] - \frac{h^5}{90} f^{(4)}(\xi)$$

Error Bound: Taylor expansion of  $f(x)$ , and bound the integrand

$$\frac{1}{24} \int_{x_0}^{x_2} f^{(4)}(\xi_x) (x - x_1)^4 = \frac{f^{(4)}(\xi_1)}{120} (x - x_1)^5 \Big|_{x_0}^{x_2} = \frac{f^{(4)}(\xi_1)}{60} h^5$$

Seperating into subintervals  $[x_0, x_2], [x_2, x_4], \dots, [x_{n-2}, x_n]$ , we have, with  $h = (b - a)/n$ ,

$$\int_a^b f(x) dx = \sum_{j=1}^{n/2} \left\{ \frac{h}{3} [f(x_{2j-2}) + 4f(x_{2j-1}) + f(x_{2j})] \right\}$$

$$\text{Error} = -\frac{h^5}{90} \sum_{j=1}^{n/2} f^{(4)}(\xi_j) = -\frac{h^5}{90} \frac{n}{2} f^{(4)}(\mu) = -\frac{b-a}{180} h^4 f^{(4)}(\mu)$$

3. Composite Mid-point rule: for subintervals  $[x_{-1}, x_1], \dots, [x_{n-1}, x_{n+1}]$  with centers

$x_0, x_2, \dots, x_n$ , and  $h = (b - a)/(n + 2)$

$$\int_a^b f(x) dx = 2h \sum_{j=0}^{n/2} f(x_{2j}) + \frac{b-a}{6} h^2 f''(\mu)$$

**Definition 1.5** [Degree of Precision] The degree of precision of a quadrature formula is  $n$  if and only if the error is zero for all polynomials of degree  $k = 0, 1, \dots, n$ , but is not zero for some polynomial of degree  $n + 1$ . ■

**Definition 1.6** [Romberg Method] First column: for  $k = 2, \dots, n$ ,  $h_k = (b - a)/(2^{k-1})$

$$R_{k,1} = \frac{1}{2} \left[ R_{k-1,1} + h_{k-1} \sum_{i=1}^{2^{k-2}} f(a + (2i-1)h_k) \right]$$

For  $k = j, j+1, \dots$ ,

$$R_{k,j} = R_{k,j-1} + \frac{1}{4^{j-1} - 1} (R_{k,j-1} - R_{k-1,j-1})$$

Stopping Criteria:  $|R_{n-1,n-1} - R_{n,n}| < \text{tol}$  and  $|R_{n-2,n-2} - R_{n-1,n-1}| < \text{tol}$ . Ensure two differently generated sets of approximations agree within the specified tolerance.

**Definition 1.7** [Gauss-Quadrature Rule] Generate Legendre polynomial  $P_0(x) = 1, P_1(x) = x$ ,

$$(n+1)P_{n+1}(x) = (2n+1)xP_n(x) - nP_{n-1}(x)$$

Let  $\{x_0, \dots, x_n\}$  be roots of  $P_{n+1}(x)$ , and

$$w_i = \int_{-1}^1 l_i(x) dx = \int_{-1}^1 \prod_{j=0, j \neq i}^n \frac{x - x_j}{x_i - x_j} dx, \quad i = 0, 1, \dots, n$$

which implies

$$\int_{-1}^1 f(x) dx = \sum_{j=0}^n w_j f(x_j) + \frac{f^{(2n+2)}(\xi)}{(2n+2)!} \int_{-1}^1 \prod_{i=0}^n (x - x_i)^2 dx$$

**Definition 1.8** [Pivoting] Why: small  $a_{kk}^{(k)}$  leads big error.

- Partial: select an element in the same column that is below the diagonal and has the largest absolute value.
- Scaled partial: first compute scale  $s_i$  for each row; then do the same as partial

**Proposition 1.1** •  $\|\mathbf{A}\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}|$

- $\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}$ . If  $\mathbf{A}$  symmetric, then  $\|\mathbf{A}\|_2 = \rho(\mathbf{A}) = \max |\lambda(\mathbf{A})|$
- $\|\mathbf{A}\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|$
- $\rho(\mathbf{A}) \leq \|\mathbf{A}\|$  for any natural norm.
- $\mathbf{A}$  is convergent iff  $\rho(\mathbf{A}) < 1$ .

*Proof.* For eigen-pair  $(\lambda, \mathbf{x})$  with  $\|\mathbf{x}\| = 1$ ,

$$|\lambda| = \|\lambda \mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{x}\| = \|\mathbf{A}\|$$

**Definition 1.9** [Jacobi's Method]

$$x_i^{(k)} = \frac{1}{a_{ii}} \left[ \sum_{j=1, j \neq i}^n (-a_{ij} x_j^{(k-1)}) + b_i \right]$$

Or for the matrix form  $\mathbf{x}^{(k)} = \mathbf{T}_j \mathbf{x}^{(k-1)} + \mathbf{c}$ :

$$\mathbf{x}^{(k)} = \mathbf{D}^{-1}(\mathbf{L} + \mathbf{U})\mathbf{x}^{(k-1)} + \mathbf{D}^{-1}\mathbf{b},$$

where  $\mathbf{D}$  are the diagonal of  $\mathbf{A}$ ;  $-\mathbf{L}, -\mathbf{U}$  are the strictly lower and upper part of  $\mathbf{A}$ . ■

**Definition 1.10** [Gauss-Seidel Method]

$$x_i^{(k)} = \frac{1}{a_{ii}} \left[ -\sum_{j=1}^{i-1} (a_{ij}x_j^{(k)}) - \sum_{j=i+1}^n (a_{ij}x_j^{(k-1)}) + b_i \right]$$

Or for the matrix form  $\mathbf{x}^{(k)} = \mathbf{T}_g \mathbf{x}^{(k-1)} + \mathbf{c}_g$ :

$$(\mathbf{D} - \mathbf{L})\mathbf{x}^{(k)} = \mathbf{U}\mathbf{x}^{(k-1)} + \mathbf{b} \implies \mathbf{x}^{(k)} = (\mathbf{D} - \mathbf{L})^{-1}\mathbf{U}\mathbf{x}^{(k-1)} + (\mathbf{D} - \mathbf{L})^{-1}\mathbf{b}$$

**Proposition 1.2** If  $\rho(\mathbf{T}) < 1$ , then  $(\mathbf{I} - \mathbf{T})^{-1}$  exists, and

$$(\mathbf{I} - \mathbf{T})^{-1} = \mathbf{I} + \mathbf{T} + \mathbf{T}^2 + \dots$$

**Proposition 1.3** The iteration  $\mathbf{x}^{(k)} = \mathbf{T}\mathbf{x}^{(k-1)} + \mathbf{c}$  converges to the unique solution to

$$\mathbf{x} = \mathbf{T}\mathbf{x} + \mathbf{c}$$

iff  $\rho(\mathbf{T}) < 1$ . The error bound holds:

- $\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \|\mathbf{T}\|^k \|\mathbf{x}^{(0)} - \mathbf{x}\|$
- $\|\mathbf{x} - \mathbf{x}^{(k)}\| \leq \frac{\|\mathbf{T}\|^k}{1 - \|\mathbf{T}\|} \|\mathbf{x}^{(1)} - \mathbf{x}^{(0)}\|$
- $\|\mathbf{x} - \mathbf{x}^{(k)}\| \approx [\rho(\mathbf{T})]^k \|\mathbf{x}^{(0)} - \mathbf{x}\|$

*Proof.* Converse: Express  $\mathbf{x}^{(k)}$  as

$$\mathbf{x}^{(k)} = \mathbf{T}^k \mathbf{x}^{(0)} + (\mathbf{T}^{k-1} + \dots + \mathbf{T} + \mathbf{I})\mathbf{c}$$



Forward: Assume converge to  $\mathbf{x}$ , and therefore

$$\mathbf{x} - \mathbf{x}^{(k)} = \mathbf{T}(\mathbf{x} - \mathbf{x}^{(k-1)}) = \dots = \mathbf{T}^k(\mathbf{x} - \mathbf{x}^{(0)})$$

for arbitrary  $\mathbf{x}^{(0)}$ . Taking limit implies  $\lim_{k \rightarrow \infty} \mathbf{T}^k \mathbf{z} = \mathbf{0}$  for any  $\mathbf{z}$ . ■

**Proposition 1.4** Sufficient condition for convergence of Jacobi and Gauss-Seidel method:  $\mathbf{A}$  is strictly diagonally dominant, i.e.,  $|a_{ii}| > \sum_{j=1, j \neq i}^n |a_{ij}|$

**Proposition 1.5** For  $a_{ij} \leq 0, i \neq j$  and  $a_{ii} > 0$ , one and only one condition holds:

- $0 \leq \rho(\mathbf{T}_g) < \rho(\mathbf{T}_j) < 1$
- $1 < \rho(\mathbf{T}_j) < \rho(\mathbf{T}_g)$
- $\rho(\mathbf{T}_j) = \rho(\mathbf{T}_g) = 0$
- $\rho(\mathbf{T}_j) = \rho(\mathbf{T}_g) = 1$

**Definition 1.11** [Gauss-Seidel Method and Relaxation]

$$x_i^{(k)} = x_i^{(k-1)} + \frac{r_{ii}^{(k)}}{a_{ii}}$$

Relaxation:

$$x_i^{(k)} = x_i^{(k-1)} + \omega \frac{r_{ii}^{(k)}}{a_{ii}}$$

Or equivalently, ( $\omega < 1$  is under-relaxation methods,  $\omega > 1$  is over-relaxation methods.)

$$x_i^{(k)} = (1 - \omega)x_i^{(k-1)} + \frac{\omega}{a_{ii}} \left[ -\sum_{j=1}^{i-1} (a_{ij}x_j^{(k)}) - \sum_{j=i+1}^n (a_{ij}x_j^{(k-1)}) + b_i \right]$$

$$\mathbf{x}^{(k)} = (\mathbf{D} - \omega\mathbf{L})^{-1}[(1 - \omega)\mathbf{D} + \omega\mathbf{U}]\mathbf{x}^{(k-1)} + \omega(\mathbf{D} - \omega\mathbf{L})^{-1}\mathbf{b}$$

**Proposition 1.6** If  $a_{ii} \neq 0$ , then  $\rho(T_\omega) \geq |\omega - 1|$ , o.e., the SOR method converges only when  $0 < \omega < 2$ ; sufficient condition for convergence:  $\mathbf{A}$  is PD and  $0 < \omega < 2$ .

**Proposition 1.7** If  $\mathbf{A}$  is PD and tridiagonal, then  $\rho(\mathbf{T}_g) = [\rho(T_j)]^2 < 1$ , the optimal choice is

$$\omega = \frac{2}{1 + \sqrt{1 - [\rho(T_j)]^2}}$$

under this choice, we have  $\rho(\mathbf{T}_\omega) = \omega - 1$

1. Direct method: Gaussian Elimination;

Advantage: Exact method, no truncation error;

Disadvantage: Computationally expansive, large round off error

Suitable: linear systems of small dimension

2. Iterative method:

Advantage: efficient in terms of both computer storage and computation

Disadvantage: not so efficient for small dimension

Suitable for: large linear sparse systems

**Definition 1.12** The condition number of nonsingular matrix  $\mathbf{A}$  is  $K(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$ .

If close to 1, then  $\mathbf{A}$  is well-conditioned, otherwise ill-conditioned. ■

**Theorem 1.1** For natural norm,

$$\|\tilde{\mathbf{x}} - \mathbf{x}\| \leq \|\mathbf{A}^{-1}\| \|\mathbf{r}\| = K(\mathbf{A}) \frac{\|\mathbf{r}\|}{\|\mathbf{A}\|}$$

$$\frac{\|\tilde{\mathbf{x}} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|} = K(\mathbf{A}) \frac{\|\mathbf{r}\|}{\|\mathbf{b}\|}, \quad \mathbf{x} \neq 0, \mathbf{b} \neq 0$$

*Proof.* Consider the first equality and  $\|\mathbf{b}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|$  ■

**Theorem 1.2** Suppose  $\mathbf{A}$  is nonsingular and  $\|\delta\mathbf{A}\| \leq \frac{1}{\|\mathbf{A}\|}$ , the solution  $\tilde{\mathbf{x}}$  to  $(\mathbf{A} + \delta\mathbf{A})\mathbf{x} = \mathbf{b} + \delta\mathbf{b}$  has the bound

$$\frac{\|\tilde{\mathbf{x}} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{K(\mathbf{A}) \|\mathbf{A}\|}{\|\mathbf{A}\| - K(\mathbf{A}) \|\delta\mathbf{A}\|} \left( \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|} + \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} \right)$$

**Proposition 1.8** For  $\mathbf{A} \in \mathbb{R}^{n \times n}$ ,  $\text{cond}(\mathbf{A}^T \mathbf{A}) = [\text{cond}(\mathbf{A})]^2$

**Theorem 1.3** For column full rank  $A \in \mathbb{R}^{m \times n}$ , we have  $A = QR$ , for  $Q \in \mathbb{R}^{m \times n}$  with orthogonal columns,  $R \in \mathbb{R}^{n \times n}$  is an upper triangular matrix.

As a result, the least square solution becomes

$$\mathbf{x}^* = R^{-1}Q^T \mathbf{b}$$

Gram-Schmidt Process: numerically unstable; ease of implementation

**Theorem 1.4 — Householder matrix.** For given vector  $\mathbf{x}$  and unit vector  $\mathbf{g}$ , define

$$H := I - 2\mathbf{u}\mathbf{u}^T,$$

$$H\mathbf{x} = \|\mathbf{x}\|\mathbf{g}, \quad \mathbf{u} = \frac{\mathbf{x} - \|\mathbf{x}\|\mathbf{g}}{\|\mathbf{x} - \|\mathbf{x}\|\mathbf{g}\|}$$

$$H_2 = \begin{pmatrix} 1 & 0 \\ 0 & \tilde{H}_2 \end{pmatrix}$$

Define  $Q = H_n \cdots H_1$ , and  $QA = R = [R_1; \mathbf{0}]$ , which implies

$$A = [Q_1, Q_2] \begin{pmatrix} R_1 \\ \mathbf{0} \end{pmatrix} = Q_1 R_1$$

