Dr. M.H. Shahnas

# Computational Physics
# Useful Books : Numerical Analysis
Richard L. Burden, J. Douglas Faires
# Numerical Recipes
William H. Press, Saul A. Teukolsky, William T. Vetterling, Brian P. Flannery
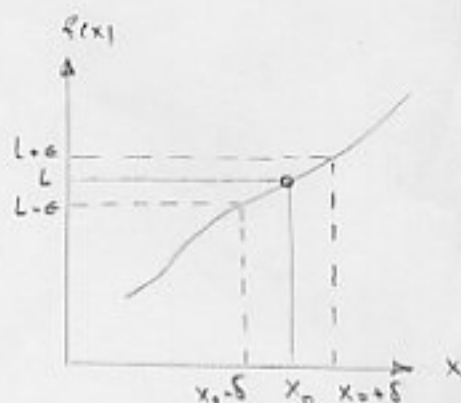
# Mathematical Preliminaries

1-1) Review of Calculus

Def. 1.1)

Let ; $f$ be a func. defined on a set $X$ of real numbers;

$$\text{Then} \quad \lim_{x \to x_0} f(x) = L$$

if $\forall \; \epsilon > 0 \quad \exists \; \delta > 0 \quad$ such that $|f(x) - L| < \epsilon$

whenever $x \in X$ and $0 < |x - x_0| < \delta$



Def. 1.2)

Let ; $f$ be a func. defined on a set $X$ of real numbers, and $x_0 \in X$ ;

$f$ is said to be <u>continuous</u> at $x_0$ if $\lim_{x \to x_0} f(x) = f(x_0)$

The func. $f$ is said to be continuous on $X$ if it is continuous at each number in $X$.

Def 1.3)

Let; $\{X_n\}_{n=1}^{\infty}$ be an infinite sequence of real or complex numbers. The sequence is said to Converge to a number $X$ (called the limit); if $\forall \epsilon > 0$ $\exists$ $N(\epsilon)$ (positive integer) such that $n > N(\epsilon)$ implies $|X_n - X| < \epsilon$.

The notation; $\lim_{n \to \infty} X_n = X$

or $X_n \to X$ as $n \to \infty$

means that the sequence $\{X_n\}_{n=1}^{\infty}$ converges to $X$

Theo. 1.4)

If $f$ is a func. defined on a set $X$ of real numbers and $X_0 \in X$, then the following are equivalent:

a) $f$ is continuous at $X_0$;

b) if $\{X_n\}_{n=1}^{\infty}$ is any sequence in $X$ converging to $X_0$,
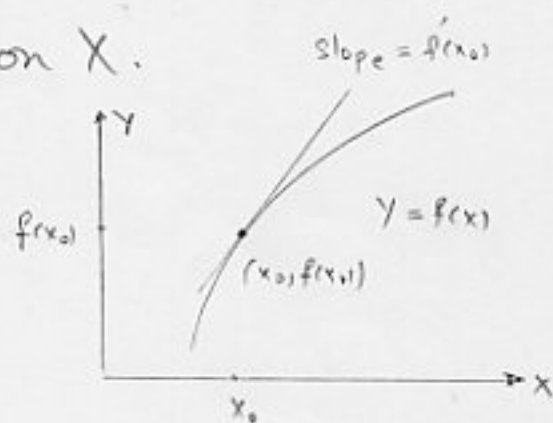
then;
$$\lim_{n \to \infty} f(X_n) = f(X_0)$$

Def. 1.5)

If $f$ is a func. defined in an open interval containing $x_0$, $f$ is said to be __differentiable__ at $x_0$, if

$$\lim \frac{f(x) - f(x_0)}{x - x_0} \qquad \text{exisits.}$$

$$\lim_{x \to x_0} \frac{f(x) - f(x_0)}{x - x_0} = f'(x_0)$$

A func. that has a derivative at each number in a set $X$ is said to be differentiable on $X$.

Theo. 1.6)

If the func. $f$ is __differentiable__ at $x_0$, then $f$ is __continuous__ at $x$.

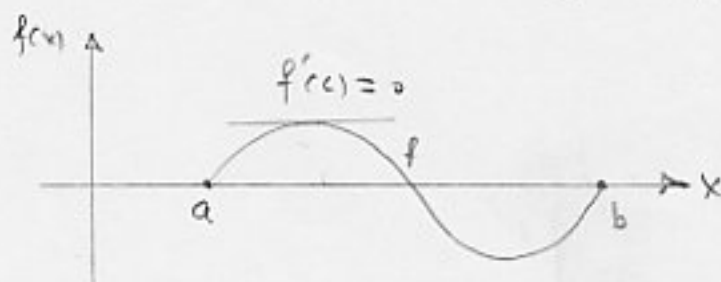Remark: $C(X)$ denotes the set of all funcs. continuous on $X$.

$C^n(X)$ denotes the set of all funcs. that have $n$ continuous derivatives on $X$

$C^\infty(X)$ denotes the set of all funcs. that have derivatives of all orders on $X$.

Theo. 1.7)    (Rolle's Theorem)

Suppose $f \in \overset{\text{continuous}}{C}[a,b]$ and $f$ is differentiable on $(a,b)$.

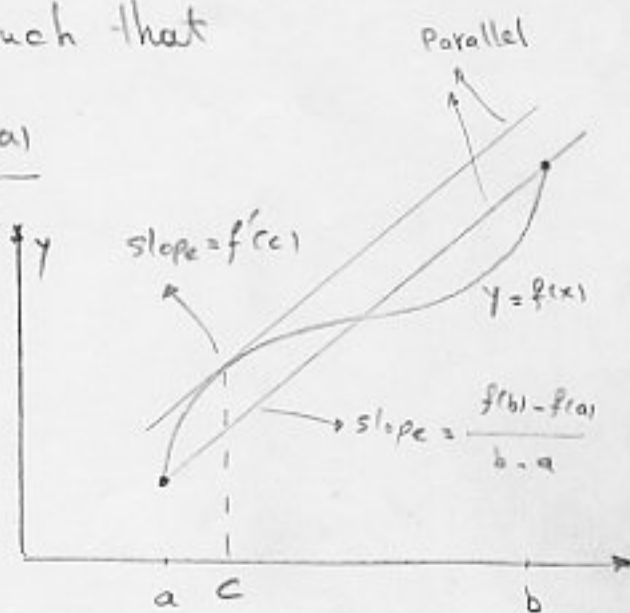If $f(a) = f(b) = 0 \xrightarrow{\text{then}} \exists\, c, \quad a < c < b$ with

$f'(c) = 0$



Theo. 1.8)    (Mean Value Theorem)

If $f \in C[a,b]$ and $f$ is differentiable on $(a,b)$;

$\xrightarrow{\text{then}} \exists\, c, \quad a < c < b$    such that
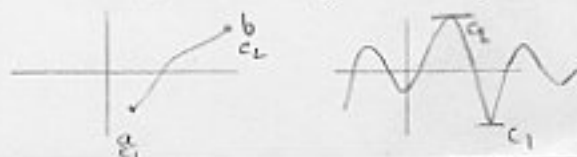
$$f'(c) = \frac{f(b) - f(a)}{b - a}$$



Theo. 1a)    (Extreme Value Theorem)

If $f \in C[a,b] \xrightarrow{\text{then}} \exists\, c_1, c_2 \in [a,b]$    with

$f(c_1) \leq f(x) \leq f(c_2), \quad \forall \quad x \in [a,b]$.

If, in addition, $f$ is differentiable on $(a,b)$, $\xrightarrow{\text{then}}$ the

numbers $c_1$ and $c_2$ occur either at endpoints of $[a,b]$

or where $f' = 0$.

Theo 1.10)    (Weighted Mean value Theorem for Integrals)

If $f \in C[a,b]$, and, $g$ is integrable on $[a,b]$ and $g(x)$ does not change sign on $[a,b]$ $\xrightarrow{\text{then}}$
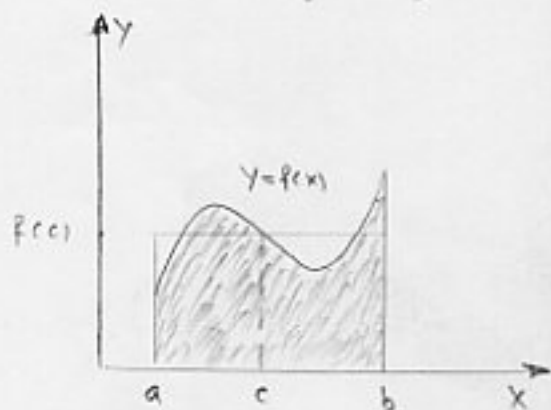
$\exists\, c$, $a < c < b$, such that;

$$\int_a^b f(x)\, g(x)\, dx = f(c) \int_a^b g(x)\, dx$$

when $g(x) \equiv 1$ $\longrightarrow$ $f(c) = \dfrac{1}{b-a} \int_a^b f(x)\, dx$ $\begin{pmatrix}\text{average value}\\ \text{of } f \text{ over}\\ [a,b]\end{pmatrix}$

Theo.1.11)  (Generalized Rolle's Theorem)

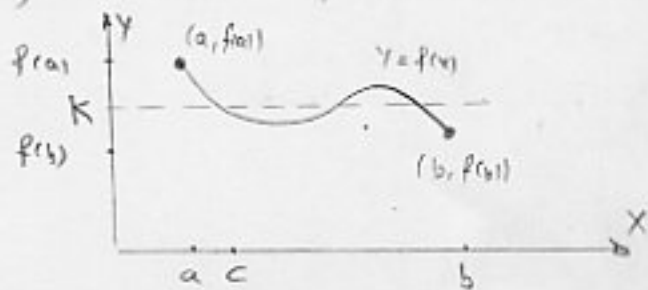Let; $f \in C[a,b]$ be $n$-times differentiable on $(a,b)$.

If $f = 0$ at $n+1$ distinct numbers $x_0 \cdots x_n$ in $[a,b]$.

$\xrightarrow{\text{then}}$ $\exists\, c \in (a,b)$ with $f^{(n)}(c) = 0$



Theo. 1.12)     (Intermediate Value Theorem)

If $f \in C[a,b]$ and $K$ is any number between $f(a)$ and $f(b)$, $\xrightarrow{\text{then}}$ $\exists\, c \in (a,b)$ for which $f(c) = K$

Ex.— Show that $x^5 - 2x^3 + 3x^2 - 1 = 0$ has a sol. in th interval $[0,1]$

Sol.

$$f(x) = x^5 - 2x^3 + 3x^2 - 1 \qquad \text{a polynomial continuous on } [0,1]$$

Since $f(0) = -1 < 0 < +1 = f(1)$

$\underset{\text{Intermediate Value Theo.}}{\underline{\hspace{3cm}}} \Rightarrow \exists x$ with $0 < x < 1$

for which $x^5 - 2x^3 + 3x^2 - 1 = 0$

Theo.1.13)   (Taylor's Theorem)

Suppose $f \in C^n[a,b]$ and $f^{(n+1)}$ exists on $[a,b]$.

$\overset{\text{continuous derivatives}}{\uparrow}$    $\underset{\text{derivative}}{\underset{\downarrow}{}}$

Let; $x_0 \in [a,b]$.

$\forall x \in [a,b]$, $\exists \xi(x)$ between $x_0$ and $x$ with

$$f(x) = P_n(x) + R_n(x)$$

Where;

$$P_n(x) = \sum_{k=0}^{n} \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k \quad \left( \begin{array}{l} \text{n-th Taylor Polynomial} \\ \text{for } f \text{ about } x_0 \end{array} \right)$$

and

$$R_n(x) = \frac{f^{(n+1)}(\xi(x))}{(n+1)!} (x - x_0)^{n+1} \quad \left( \begin{array}{l} \text{remainder term} \\ \text{or truncation error} \end{array} \right)$$

Remark: $n \to \infty$    ( Taylor series )

$x_0 = 0$    ( Maclaurin polynomial or series )

Ex.2 - Determine a) the second and b) the third Taylor polynomial for $f(x) = \cos x$ about $x_0 = 0$, and use these polynomials to approximate $\cos(0.01)$

Sol.

Since ; $f \in \overset{\infty}{C}(\mathbb{R})$, the previous theorem can be used for any $n > 0$

a) $n = 2$ and $x_1 = 0$ ; $\cos x = 1 - \frac{1}{2}x^2 + \frac{1}{6}x^3 \sin\xi(x)$

where $\xi(x)$ is a number between $0$ and $x$.



with $x = 0.01$

$\cos(0.01) = 1 - \frac{1}{2}(0.01)^2 + \frac{1}{6}(0.01)^3 \sin(\xi(x))$

$= 0.99995 + (0.1\overline{66}) \cdot 10^{-6} \sin(\xi(x))$

where $0 < \xi(x) < 0.01$      (rad.)

Since $|\sin(\xi(x))| < 1$ ——⟶ We can use $0.99995$ as an approximation to $\cos(0.01)$ with assurance of at least six-decimal place accuracy.

Using standard tables $\cos(0.01) = 0.99995000042$

b) For the third Taylor polynomial about $x_0 = 0$;

$$\cos x = 1 - \frac{1}{2}x^2 + 0 + \frac{1}{24}x^4 \cos(\xi(x))$$

where $\quad 0 < \xi(x) < 0.01 \quad$ and $\quad f'''(0) = 0$

the approximation $\longrightarrow \quad G(0.01) = 0.99995 \qquad$ as before

the accuracy $\longrightarrow \quad$ nine-decimal-place , since

$$\left|\frac{1}{24}x^4 \cos \xi(x)\right| \leq \frac{1}{24}(0.01)^4 (1) = 4.2 \times 10^{-10}$$

## 1.2 Round off Errors and Computer Arithmetic

Floating-point forms

$$a = \underbrace{\pm 0.d_1 d_2 \cdots d_k}_{\text{mantissa}} \underbrace{\times 10^n}_{\text{characteristic}} \qquad 1 \leq d_1 \leq 9 \;,\; 0 \leq d_i \leq 9$$

$$i = 2, \ldots k$$

$$-78 \leq n \leq 76$$

$$y = 0.d_1 d_2 \cdots d_k d_{k+1} d_{k+2} \cdots \times 10^n \qquad \text{(any positive real number)}$$

There are two ways to obtain the floating-point form;

i) chopping $\qquad fl(y) = 0.d_1 d_2 \cdots d_k \times 10^n$

The mantissa is terminated at $k$ decimal digits, just by chopping.

ii) Rounding $\qquad fl(y) = 0. \delta_1 \delta_2 \cdots \delta_k \times 10^n$

$$\begin{cases} \text{if } d_{k+1} \geq 5 \; ; & d_k \to d_{k+1} \quad \text{and chop off right most digits} \\ \text{. } d_{k+1} < 5 \; ; & d_k \to d_k \quad \text{" " " " " " " "} \end{cases}$$

### Ex.

$$\Pi = 3.14159265 \cdots$$

$$\Pi = 0.314159265 \cdots \times 10^1 \qquad \text{normalized decimal form}$$

$$fl(\Pi) = 0.31415 \times 10^1 = 3.1415 \qquad \text{(chopping)}$$

$$fl(\Pi) = (0.31415 + 0.00001) \times 10^1 = 3.1416 \qquad \text{(rounding)}$$

Round off error:

$$= |Y - fl(y)|$$

(regardless of whether the rounding or chopping method is used)

Def. 1.14) If $P^*$ is a approximation to $p$:

$$\text{absolute error} = |P - P^*|$$
$$\text{relative } \quad = \frac{|P - P^*|}{|P|} \qquad P \neq 0$$

EX.

a) $P = 0.3000 \times 10^1$      $P^* = 0.3100 \times 10^1$

     ab. error $= 0.1$ , rel. error $= 0.3\overline{333} \times 10^{-1}$

b) $P = 0.3000 \times 10^{-3}$      $P^* = 0.3100 \times 10^{-3}$

     ab. error $= 0.1 \times 10^{-4}$ , rel. error $= 0.3\overline{333} \times 10^{-1}$

c) $P = 0.3000 \times 10^4$      $P^* = 0.3100 \times 10^4$

     ab error $= 0.1 \times 10^3$ , rel. error $= 0.3\overline{333} \times 10^{-1}$

Relative error in computer;    (chopping)

$$= \left| \frac{y - fl(y)}{y} \right|$$

$$y = 0.d_1 d_2 \ldots - d_k d_{k+1} \cdots \times 10^n$$

$$\left| \frac{y - fl(y)}{y} \right| \underset{\substack{chopping}}{=} \left| \frac{0.d_1 d_2 \ldots - d_k d_{k+1} \ldots \times 10^n - 0.d_1 d_2 \ldots d_k \times 10^n}{0.d_1 d_2 \ldots \ldots \times 10^n} \right|$$

$$= \left| \frac{0.d_{k+1} d_{k+2} \ldots \times 10^{n-k}}{0.d_1 d_2 \ldots - \times 10^n} \right| = \left| \frac{0.d_{k+1} d_{k+2} \ldots}{0.d_1 d_2 \ldots} \right| \times 10^{-k}$$

Since   $d_1 \neq 0 \longrightarrow Min(denominator) = 0.1$

and     $Max(numerator) = 1$

$$\longrightarrow \left| \frac{y - fl(y)}{y} \right| \leq \frac{1}{0.1} \times 10^{-k} = 10^{-k+1}$$

In a similar manner, a bound for the relative error when using k-digit rounding arithmetic is:

$$\left| \frac{y - fl(y)}{y} \right|_{rounding} = \left| \frac{0.d_1 d_2 \ldots - d_k d_{k+1} \ldots \times 10^n - 0.\delta_1 \delta_2 \ldots \delta_k \times 10^n}{0.d_1 d_2 \ldots - \times 10^n} \right|$$

$$\leq 0.5 \times 10^{-k+1}$$

Def. 1.15)

The number $p^*$ is said to approximate $p$ to $t$ significant digits (or figures) if $t$ is the largest nonnegative integer for which

$$\left| \frac{p - p^*}{p} \right| < 5 \times 10^{-t}$$

Other Source of error:

In addition to inaccurate representation of numbers (floating point), the arithmatic performed in a computer is not exact.

Assume finite-digit arithmetic given by;

$$x \oplus y = fl(fl(x) + fl(y)) \qquad x \otimes y = fl(fl(x) \times fl(y))$$

$$x \ominus y = fl(fl(x) - fl(y)) \qquad x \oslash y = fl(fl(x) \div fl(y))$$

Ex.

Suppose; $x = \frac{1}{3}$, $y = \frac{5}{7}$ (use 5-digit chopping)

$$fl(x) = 0.33333 \times 10^0 \qquad fl(y) = 0.71428 \times 10^0$$

| Operation | Result | Actual value | Abs. error | Relative error |
|---|---|---|---|---|
| $x \oplus y$ | $0.10476 \times 10^1$ | $22/21$ | $0.190 \times 10^{-4}$ | $0.182 \times 10^{-4}$ |
| $y \ominus x$ | $0.38095 \times 10^0$ | $8/21$ | $0.238 \times 10^{-5}$ | $0.625 \times 10^{-5}$ |
| $x \otimes y$ | $0.23809 \times 10^0$ | $5/21$ | $0.524 \times 10^{-5}$ | $0.220 \times 10^{-4}$ |
| $y \oslash x$ | $0.21428 \times 10^1$ | $15/7$ | $0.571 \times 10^{-4}$ | $0.267 \times 10^{-4}$ |

Max. (Rel. error) $= 0.267 \times 10^{-4}$        satisfactory

Now consider;

$$u = 0.714251 , \quad v = 98765.9 , \quad w = 0.111111 \times 10^{-4}$$

$\longrightarrow fl(u) = 0.71425 \times 10^0$ , $fl(v) = 0.98765 \times 10^5$, $fl(w) = 0.11111 \times 10^{-4}$

| Operation | Result | Actual value | Abs. error | Rel. error |
|---|---|---|---|---|
| $Y \ominus u$ | $0.30000 \times 10^{-4}$ | $0.34714 \times 10^{-4}$ | $0.471 \times 10^{-5}$ | $0.136$ |
| $(Y \ominus u) \oplus w$ | $0.27000 \times 10^{1}$ | $0.31243 \times 10^{1}$ | $0.424$ | $0.136$ |
| $(Y \ominus u) \otimes V$ | $0.29629 \times 10^{1}$ | $0.34285 \times 10^{1}$ | $0.465$ | $0.136$ |
| $u \oplus V$ | $0.98765 \times 10^{5}$ | $0.98766 \times 10^{5}$ | $0.161 \times 10^{1}$ | $0.163 \times 10^{-4}$ |

Other most common error source;

Subtraction of nearly equal numbers:

$$fl(x) = 0.d_1 d_2 \cdots d_p \alpha_{p+1} \alpha_{p+2} \cdots \alpha_k \times 10^{n}$$

$$fl(y) = 0.d_1 d_2 \cdots d_p \beta_{p+1} \beta_{p+2} \cdots \beta_k \times 10^{n}$$

$$fl(fl(x) - fl(y)) = 0.\underbrace{\alpha_{p+1} \alpha_{p+2} \cdots \alpha_k}_{k-p \text{ digits of significance}} \times 10^{n-p}$$

where

$$0.\alpha_{p+1} \alpha_{p+2} \cdots \alpha_k = 0.\alpha_{p+1} \alpha_{p+2} \cdots \alpha_k - 0.\beta_{p+1} \beta_{p+2} \cdots \beta_k$$

However in most calculation devices;

$$fl(fl(x) - fl(y)) = 0.\alpha_{p+1} \alpha_{p+2} \cdots \alpha_k \underbrace{0 0 0 0 0}_{P \text{ digits}} \times 10^{n-p}$$

Errors in chain;

Suppose;

$$fl(\text{Some arithmetic calculation}) = Z + \delta$$

(having finite digit)

$\delta$: error of arithmetic cal.

$Z$: actual floating point form of arithmetic cal.

$$\frac{Z}{\epsilon} \approx fl\left[\frac{(Z+\delta)}{fl(\epsilon)}\right] \qquad (\epsilon \neq 0)$$

suppose; $\epsilon = 10^{-n}$ $\qquad n > 0$

then

$$\frac{Z}{\epsilon} = Z \times 10^{n}$$

and

$$fl\left[\frac{(Z+\delta)}{fl(\epsilon)}\right] = (Z+\delta) \times 10^{n}$$

$$\text{abs. error} = |\delta| \times 10^{n}$$

Thus; if a finite-digit representation or calculation introduces an error, further enlargement of the error occurs when dividing by a number with small magnitude (or equivaletely, when multiplying by a number with large magnitude.).

## Ex.

$$ax^2 + bx + c = 0 \qquad a \neq 0$$

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

Consider; $\quad x^2 + 62.10x + 1 = 0$

$$x_1 = -0.01610723 \qquad x_2 = -62.08390 \qquad \text{(approx.)}$$

Since $\quad b^2 \gg 4ac \quad \longrightarrow \quad b \approx \sqrt{b^2 - 4ac}$

Now suppose 4-digit rounding arithmetic;

$$\sqrt{b^2 - 4ac} = \sqrt{(62.10)^2 - 4.000} = \sqrt{3856. - 4.000} = \sqrt{3852.} = 62.06$$

So;

$$fl(x_1) = \frac{-b + \sqrt{b^2 - 4ac}}{2a} = \frac{-62.10 + 62.06}{2.000} = \frac{-0.04000}{2.000} = -0.02000$$

is an approximation to $x_1 = -0.01610723$;

$$\text{Rel. error} = \left| \frac{-0.01611 - 0.02000}{-0.01611} \right| \approx 2.4 \times 10^{-1}$$

On the other hand;

$$fl(x_2) = \frac{-b - \sqrt{b^2 - 4ac}}{2a} = \frac{-62.10 - 62.06}{2.000} = \frac{-124.2}{2.000} = -62.10$$

is an approximation to $x_2 = -62.08$

$$\text{Rel. error} = \left| \frac{-62.08 + 62.10}{-62.08} \right| \approx 3.2 \times 10^{-4}$$

Remedy:

$$X_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a} \left( \frac{-b - \sqrt{b^2 - 4ac}}{-b - \sqrt{b^2 - 4ac}} \right) = \frac{-2c}{b + \sqrt{b^2 - 4ac}}$$

$$\longrightarrow fl(x_1) = \frac{-2.000}{62.10 + 62.06} = \frac{-2.000}{124.2} = -0.01610$$

$$\longrightarrow \text{Rel. error} = 6.2 \times 10^{-4}$$

The rationalization tequnique can also be applied to give the alternate form for $x_2$:

$$X_2 = \frac{-2c}{b - \sqrt{b^2 - 4ac}}$$

This would be the form to use if $b$ were a negative number.

But notice:

$$fl(x_2) = \frac{-2c}{b - \sqrt{b^2 - 4ac}} = \frac{2.000}{62.10 - 62.06} = \frac{-2.000}{0.04000} = -50.00$$

here then are two sources of error:

$$\begin{cases} 1 - \text{subtraction of nearly equal numbers} \\ 2 - \text{division by the small number.} \end{cases}$$

Ex.

Evaluate $f(x) = x^3 - 6x^2 + 3x - 0.149$ at $x = 4.71$
using three-digit arithmetic.

| | $x$ | $x^2$ | $x^3$ | $6x^2$ | $3x$ |
|---|---|---|---|---|---|
| Exact | 4.71 | 22.1841 | 104.487111 | 133.1046 | 14.13 |
| Three digit (chopping) | 4.71 | 22.1 | 104. | 132. | 14.1 |
| " " (rounding) | 4.71 | 22.2 | 105. | 133. | 14.1 |

Exact :  $f(4.71) = 104.487111 - 133.1046 + 14.13 - 0.149 = -14.636489$

Three-digit
(chopping)  $f(4.71) = 104. - 132. + 14.1 - 0.149 = -14.0$

Three-digit
(rounding)  $f(4.71) = 105. - 133. + 14.1 - 0.149 = -14.0$

Rel. error $= \left| \dfrac{-14.636489 - 14.0}{-14.636489} \right| \approx 0.04$

(for both method)

Alternative approach ; (nesting method)

$f(x) = x^3 - 6x^2 + 3x - 0.149 = ((x-6)x + 3)x - 0.149$

$\longrightarrow$ $f(4.71) = ((4.71 - 6)4.71 + 3)4.71 - 0.149 = -14.5$     3-digit chopping

     $f(4.71) = -14.6$    3-digit rounding

Rel. error $= \left| \dfrac{-14.636489 + 14.5}{-14.636489} \right| \approx 0.0093$ ,   Rel. error $\approx \left| \dfrac{-14.636489 + 14.6}{-14.636489} \right| \approx 0.0025$
chopping

## 1.3 Convergence:

**Ex.** The Taylor polynomial $P_N(x)$ for $f(x) = \ln x$ expanded about $x_0 = 1$ is

$$P_N(x) = \sum_{i=1}^{N} \frac{(-1)^{i+1}}{i} (x-1)^i$$

$$\ln 1.5 = 0.40546511 \qquad \text{to the 8-decimal places.}$$

We want to find $\text{Min.}(N) = ?$ for;

$$| \ln 1.5 - P_N(1.5) | < 10^{-5}$$

without using the Taylor polynomial truncation formula.

From calculus we know;

$$\text{If} \quad A_n = \sum_{n=1}^{N} a_n \quad \text{and} \quad \lim_{N \to \infty} A_N = A \quad \underline{\text{(with decreasing terms)}}$$

$$\xrightarrow{\text{then}} \quad | A - A_N | \leq | a_{N+1} |$$

## Stability:

One criterion we will impose on an algorithm whenever possible is that small changes in the initial data produce correspondingly small changes in the final results. An algorithm that satisfies this property is called stable; it is unstable when this criterion is not fulfilled.

## Rounding Error Growth and its Connection to Algorithm Stability:

Suppose:    $E$ : an error introduced at some stage in the calculation

$E_n$ : the error after n subsequent operations

Def. 1.6 ) Suppose that $E_n$ represents the growth of an error after n subsequent operations.

If $|E_n| \approx c n E$    where c = const. indep. of n

→ the growth of error is said to be linear.

If $|E_n| \approx K^n E$    for some $k > 1$

→ the growth of error is called exponential.

Ex. —  The sequence $P_n = (\frac{1}{3})^n$, $n > 0$, can be generated recursively by letting $P_0 = 1$ and defining $P_n = (\frac{1}{3}) P_{n-1}$, whenever $n > 1$.

In 5-digit rounding arithmetic:

$P_0 = 0.10000 \times 10^1$      $P_1 = 0.33333 \times 10^0$      $P_2 = 0.11111 \times 10^0$

$P_3 = 0.37036 \times 10^{-1}$      $P_4 = 0.12345 \times 10^{-1}$ . —

Rounding error replacing $\frac{1}{3}$ by $0.33333$ produces an error of only $(0.33333)^n \times 10^{-5}$ in the $n$th term of the sequence.

$$\text{Rel. error} = \left| \frac{\frac{1}{3} - 0.33333}{\frac{1}{3}} \right| = 1 \times 10^{-5}$$

Abs. error = Rel. error × Data

Abs. error $= (0.33333) \times 10^{-5}$      for $n = 1$

Another way to generate the sequence:

Define $P_0 = 1$, $P_1 = \frac{1}{3}$

$$P_n = (\frac{10}{3}) P_{n-1} - P_{n-2} \qquad n \geq 2 \qquad (1)$$

This method is quite clearly <u>unstable</u>.

Note that formula (1) is satisfied whenever $P_n$ is of the form;

$$P_n = C_1 \left(\tfrac{1}{3}\right)^n + C_2 \, 3^n$$

| n | Computed $P_n$ | Correct value $P_n$ |
|---|---|---|
| 0 | $0.10000 \times 10^1$ | $0.10000 \times 10^1$ |
| 1 | $0.33333 \times 10^0$ | $0.33333 \times 10^0$ |
| 2 | $0.11110 \times 10^0$ | $0.11111 \times 10^0$ |
| 3 | $0.37000 \times 10^{-1}$ | $0.37037 \times 10^{-1}$ |
| 4 | $0.12230 \times 10^{-1}$ | $0.12346 \times 10^{-1}$ |
| 5 | $0.37660 \times 10^{-2}$ | $0.41152 \times 10^{-2}$ |
| 6 | $0.32300 \times 10^{-3}$ | $0.13717 \times 10^{-2}$ |
| 7 | $-0.26893 \times 10^{-2}$ | $0.45725 \times 10^{-3}$ |
| 8 | $-0.92872 \times 10^{-2}$ | $0.15242 \times 10^{-3}$ |

Verification:

$$\tfrac{10}{3} P_{n-1} - P_{n-2} = \tfrac{10}{3} \left[ C_1 \left(\tfrac{1}{3}\right)^{n-1} + C_2 3^{n-1} \right] - \left[ C_1 \left(\tfrac{1}{3}\right)^{n-2} + C_2 3^{n-2} \right]$$

$$= C_1 \left[ \tfrac{10}{3} \left(\tfrac{1}{3}\right)^{n-1} - \left(\tfrac{1}{3}\right)^{n-2} \right] + C_2 \left[ \tfrac{10}{3} 3^{n-1} - 3^{n-2} \right]$$

$$= C_1 \left(\tfrac{1}{3}\right)^n + C_2 3^n = P_n$$

To have $\begin{cases} P_0 = 1 \\ P_1 = \tfrac{1}{3} \end{cases}$ in equ (1) $\xrightarrow{\text{must}}$ $\begin{cases} C_1 = 1 \\ C_2 = 0 \end{cases}$

However, in the 5-digit approx. $\begin{cases} P_0 = 0.10000 \times 10^1 \\ P_1 = 0.33333 \times 10^0 \end{cases}$

$\xrightarrow{\text{thus}} \begin{cases} C_1 = 0.10000 \times 10^1 \\ C_2 = -0.12500 \times 10^{-5} \end{cases}$

This small change $C_1$ and $C_2$ results a rounding error;

$$0.10000 \times 10^1 \left(\tfrac{1}{3}\right)^n + (-0.12500 \times 10^{-5}) \times 3^n$$

↓ decreasing with n (exponential)

↓ increasing with n (exponential)

Remedy:

i) Using double - or multi-precision arithmetic.

ii) Using a more suitable algorithm.

Def. 1.17 -

Suppose $\{\alpha_n\}_{n=1}^{\infty}$ is a sequence that converges to a number $\alpha$. We say that $\{\alpha_n\}_{n=1}^{\infty}$ converges to $\alpha$ with <u>rate of convergence</u> $O(\beta_n)$, where $\{\beta_n\}_{n=1}^{\infty}$ is another sequence with $\beta_n \neq 0$ for each $n$, if

$$\frac{|\alpha_n - \alpha|}{|\beta_n|} \leq K \qquad \text{for sufficiently large } n$$

and $K = $ const. indep. of $n$. This is indicated by writing

$$\alpha_n = \alpha + O(\beta_n)$$

or

$$\alpha_n \longrightarrow \alpha \qquad \text{with rate of convergence } O(\beta_n)$$

Ex. -

Suppose that the sequences $\{\alpha_n\}$ and $\{\hat{\alpha}_n\}$ are described by

$$\alpha_n = \frac{n+1}{n^2}$$

$$\hat{\alpha}_n = \frac{n+3}{n^3} \qquad \forall n \geq 1 \qquad n : \text{int.}$$

Although;

$$\lim_{n \to \infty} \alpha_n = 0 \qquad \lim_{n \to \infty} \hat{\alpha}_n = 0$$

the sequence $\{\hat{\alpha}_n\}$ <u>converges</u> to this limit much <u>faster</u> than $\{\alpha_n\}$.

If we let $\qquad B_n = \frac{1}{n} \quad , \quad \hat{B}_n = \frac{1}{n^2}$

$$\left| \frac{\alpha_n - 0}{B_n} \right| = \left| \frac{(n+1)/n^2 - 0}{(1/n)} \right| = \frac{n+1}{n} \leq 2$$

$$\left| \frac{\hat{\alpha}_n - 0}{\hat{B}_n} \right| = \left| \frac{(n+3)/n^3 - 0}{(1/n^2)} \right| = \frac{n+3}{n} \leq 4$$

$$\longrightarrow \quad \alpha_n = 0 + O\left(\frac{1}{n}\right) \qquad \text{while} \quad \hat{\alpha}_n = 0 + O\left(\frac{1}{n^2}\right)$$

This implies that;

The rate of convergence of $\{\alpha_n\}$ is similar to the convergence of $\{\frac{1}{n}\}$ to zero

while " " " " " , $\{\hat{\alpha}_n\}$ " " " " $\{\frac{1}{n^2}\}$ " "

This concept generalizes to functions as follows:

Def 1.18 –

If $\lim\limits_{x \to 0} F(x) = L$,

the convergence is said to be $O(G(x))$

if $\exists \; k > 0$ indep. of $x$ for which

$$\frac{|F(x) - L|}{|G(x)|} \leq k \qquad \text{for sufficiently } \underline{small} \; |x| > 0$$

This situation is indicated by writing:

$$F(x) = L + O(G(x))$$

or $\quad F(x) \underset{x \to 0}{\longrightarrow} L \qquad$ with rate of convergence $O(G(x))$

Ex. — We found before using a third Taylor polynomial,

$$\cos x = 1 - \frac{1}{2} x^2 + \frac{1}{24} x^4 \cos \xi(x)$$

$$\longrightarrow \cos x + \frac{1}{2} x^2 = 1 + \frac{1}{24} x^4 \cos \xi(x) \qquad \begin{array}{ll} 0 \leq \xi(x) \leq x & \text{if } x \geq 0 \\ 0 \geq \xi(x) \geq x & \text{if } x \leq 0 \end{array}$$

$$\longrightarrow \cos x + \frac{1}{2} x^2 = 1 + O(x^4)$$

Since $\left| \dfrac{(\cos x + \frac{1}{2} x^2) - 1}{x^4} \right| = \left| \dfrac{1}{24} \cos \xi(x) \right| \leq \dfrac{1}{24}$

The implication is that $\cos x + \frac{1}{2} x^2$ converges to <u>its limit</u> , <u>1</u>, at approximately the <u>same rate</u> that $x^4$ converges to <u>zero</u>.

# Chapter 11

## Boundary-Value Probs. for Ordinary Differential Equs.

### 11.3 Finite Difference Methods for Linear Probs.

$$y'' = P(x)\, y' + q(x)\, y + r(x) \qquad a \le x \le b \;,\; y(a) = \alpha \;,\; y(b) = \beta$$

First, we select an integer $N > 0$ and divide the interval $[a, b]$ into $(N+1)$ equal subintervals

$$x_i = a + ih \quad \text{(mesh points)} \qquad h = \frac{b-a}{N+1} \qquad i = 0, 1, \ldots N+1$$

$$X_0 = a, \quad X_{N+1} = b$$

$$y''(x_i) = P(x_i)\, y'(x_i) + q(x_i)\, y(x_i) + r(x_i)$$

$$i = 1, 2, \ldots N \qquad \text{interior mesh points}$$

Expanding $y$ in third-deg. Taylor polynomial about $x_i$ evaluated at $x_{i+1}$ and $x_{i-1}$;

$$y(x_{i+1}) = y(x_i + h) = y(x_i) + h y'(x_i) + \frac{h^2}{2} y''(x_i) + \frac{h^3}{6} y'''(x_i) + \frac{h^4}{24} y^{(4)}(\xi_i^+) \qquad (1)$$

for some $\xi_i^+$, $\qquad x_i \le \xi_i^+ \le x_{i+1}$

$$y(x_{i-1}) = y(x_i - h) = y(x_i) - h y'(x_i) + \frac{h^2}{2} y''(x_i) - \frac{h^3}{6} y'''(x_i) + \frac{h^4}{24} y^{(4)}(\xi_i^-) \qquad (2)$$

for some $\xi_i^-$ $\qquad x_{i-1} \le \xi_i^- \le x_i$

assuming $y \in C^4[x_{i-1}, x_{i+1}]$

$(1)(2) \longrightarrow \ddot{y}(x_i) = \frac{1}{h^2}\left[y(x_{i+1}) - 2y(x_i) + y(x_{i-1})\right] - \frac{h^2}{24}\left[y^{(4)}(\xi_i^+) + y^{(4)}(\xi_i^-)\right]$

Using intermediate value Theo.;
$$\begin{cases} y^{(4)}(\xi_i^-) \equiv f(a), \; y^{(4)}(\xi_i^+) \equiv f(b) \\ \frac{1}{2}\left[y^{(4)}(\xi_i^-) + y^{(4)}(\xi_i^+)\right] \equiv k \;\; \text{(between } f(a) \text{ and)} \\ \qquad\qquad\qquad\qquad\qquad\qquad\qquad f(b) \\ \rightarrow \xi_i \equiv C \end{cases}$$

$$y''(x_i) = \frac{1}{h^2}\left[y(x_{i+1}) - 2y(x_i) - y(x_{i-1})\right] - \frac{h^2}{12}y^{(4)}(\xi_i)$$

$$(3)$$

$$x_{i-1} \leq \xi_i \leq x_{i+1}$$

Equ. (3) is called the Centered-difference formula for $y''(x_i)$

Similarly:

$$y'(x_i) = \frac{1}{2h}\left[y(x_{i+1}) - y(x_{i-1})\right] - \frac{h^2}{6}y''(\eta_i) \qquad (4)$$

$$x_{i-1} \leq \eta_i \leq x_{i+1} \qquad\qquad \text{centered-diff. formula for } y'(x_i)$$

The use of (3) and (4) in the main differential equ. yields:

$$\frac{y(x_{i-1}) - 2y(x_i) + y(x_{i+1})}{h^2} = P(x_i)\left[\frac{y(x_{i+1}) - y(x_{i-1})}{2h}\right] + q(x_i)y(x_i)$$

$$+ r(x_i) - \frac{h^2}{12}\left[2P(x_i)y''(\eta_i) - y^{(4)}(\xi_i)\right]$$

A Finite-Difference method with truncation error of order $O(h^2)$ results by using this equ. together with the boundary conds. $y(a) = \alpha$, and $y(b) = \beta$ to define

$$W_0 = \alpha \quad , \quad W_{N+1} = \beta$$

and

$$\left( \frac{2W_i - W_{i+1} - W_{i-1}}{h^2} \right) + P(x_i) \left( \frac{W_{i+1} - W_{i-1}}{2h} \right) + q(x_i) W_i$$

$$= -r(x_i)$$

$i = 1, --- N$

or;

$$-\left( 1 + \frac{h}{2} P(x_i) \right) W_{i-1} + \left( 2 + h^2 q(x_i) \right) W_i - \left( 1 - \frac{h}{2} P(x_i) \right) W_{i+1} = -h^2 r(x_i)$$

and the resulting system of equs. is expressed in the tridiagonal $N \times N$-matrix form;

$$AW = b$$

$$A = \begin{bmatrix} 2 + h^2 q(x_1) & -1 + \frac{h}{2} P(x_1) & 0 & - & - & - & - & - & - & - & - & 0 \\ -1 - \frac{h}{2} P(x_2) & 2 + h^2 q(x_2) & -1 + \frac{h}{2} P(x_2) & & & & & & & & & \vdots \\ 0 & & & & & & & & & & & 0 \\ \vdots & & & & & & & & & & -1 + \frac{h}{2} P(x_{N-1}) \\ 0 & - & - & - & - & - & - & 0 & -1 - \frac{h}{2} P(x_N) & 2 + h^2 q(x_N) \end{bmatrix}$$

$$W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ \\ w_N \end{bmatrix} \qquad b = \begin{bmatrix} -h^2 r(x_1) + (1 + \frac{h}{2} p(x_1)) w_0 \\ -h^2 r(x_2) \\ \vdots \\ -h^2 r(x_{N-1}) \\ -h^2 r(x_N) + (1 - \frac{h}{2} p(x_N)) w_{N+1} \end{bmatrix}$$

Linear Finite-Difference Algorithm 11.3

To approximate the sol. of the boundary-value problem:

$$y'' = p(x) y' + q(x) y + r(x) \qquad a \leq x \leq b, \quad y(a) = \alpha, \; y(b) = \beta$$

Input: $a, b, \alpha, \beta, N$

Output: approxs. $w_i$ to $y(x_i)$ for $i = 0, 1, \ldots N+1$

S1      $h = (b-a)/(N+1)$

        $x = a + h$

        $a_1 = 2 + h^2 q(x)$

        $b_1 = -1 + (h/2) p(x)$

        $d_1 = -h^2 r(x) + (1 + (h/2) p(x)) \alpha$

S2      Do $i = 2, N-1$

        $x = a + ih$

        $a_i = 2 + h^2 q(x)$

        $b_i = -1 + (h/2) p(x)$

$$c_i = -1 - (h/2) P(x)$$

$$d_i = -h^2 r(x)$$

Continue

S3  $x = b - h$

$$a_N = 2 + h^2 q(x)$$

$$c_N = -1 - (h/2) P(x)$$

$$d_N = -h^2 r(x) + (1 - (h/2) P(x)) \beta$$

To be continued by the algorithm for solving $\underline{\text{linear}}$ system of equs..

Remark: Non-linear differential equ.:

$$y'' = f(x, y, y') \qquad \text{general form}$$

Ex.:  $y'' = \frac{1}{8}(32 + 2x^3 - yy')$

$\underline{\text{Ex.}}$  $y'' + \frac{1}{x} y' - \frac{y}{x^2} = 3$  , $y(1) = 2$ , $y(2) = 3$

$\longrightarrow a = 1, \quad b = 2, \qquad \alpha = 2, \quad \beta = 3$

$P(x) = \frac{-1}{x}$ , $q(x) = +\frac{1}{x^2}$    $r(x) = 3$

$N = 5 \quad \longrightarrow h = 0.2, \quad x_1 = 1.2, \ x_2 = 1.4, \ x_3 = 1.6, \ x_4 = 1.8, \ x_5 = 2$

$$-2.0278\,Y_1 + 1.0833\,Y_2 \qquad\qquad = 0.12 - 1.8333$$

$$0.9286\,Y_1 - 2.0204\,Y_2 + 1.0714\,Y_3 \qquad\qquad = 0.12$$

$$0.9375\,Y_2 - 2.0156\,Y_3 + 1.0625\,Y_4 = 0.12$$

$$0.9444\,Y_3 - 2.0123\,Y_4 = 0.12 - 3.1667$$

| X | analytic Sol. | Numerical Sol. |
|---|---|---|
| 1 | 2 | 2 |
| 1.2 | 1.9067 | 1.9083 |
| 1.4 | 1.9886 | 1.9904 |
| 1.6 | 2.2100 | 2.2115 |
| 1.8 | 2.5511 | 2.5570 |
| 2 | 3 | 3 |

Approximating Eigenvalues:

Jacobi Method;

If A is symmetric, it can be shown there exists an <u>orthogonal</u> matrix $P$ such that:

$$P^T A P = D$$

where $D$ is diagonal matrix whose diagonal elements are the <u>eigenvalues</u> of A.

Jacobi's method transforms A into diagonal form by annihilating its off-diagonal elements one-by-one.

It makes use of <u>plane rotation</u> matrices $R(p,q)$ which are basically <u>unit</u> matrices except for elements;

$$r_{pp} = C\theta \qquad r_{pq} = -S\theta$$

$$r_{qp} = S\theta \qquad r_{qq} = C\theta$$

$R(p,q)$ orthogonal matrix;

then a tr. of the form $\bar{A} = R^T A R$ preserves the eigenvalues of A.

$\theta$ is chosen in such a way that $\bar{a}_{pq}$ is reduced to <u>zero</u>.

$$\bar{A} = R^T A R = \begin{bmatrix} 1 & \cdots & 0 & \cdots & 0 & \cdots & 0 \\ 0 & \cdots & c & \cdots & s & \cdots & 0 \\ \vdots & & & & & & \\ 0 & \cdots & -s & \cdots & c & \cdots & 0 \\ \vdots & & & & & & \\ 0 & \cdots & 0 & \cdots & 0 & \cdots & 1 \end{bmatrix} \begin{bmatrix} a_{11} & \cdots & a_{1p} & \cdots & a_{1q} & \cdots & a_{1n} \\ a_{p1} & \cdots & a_{pp} & \cdots & a_{pq} & \cdots & a_{pn} \\ \vdots & & & & & & \\ a_{q1} & \cdots & a_{qp} & \cdots & a_{qq} & \cdots & a_{qn} \\ a_{n1} & \cdots & a_{np} & \cdots & a_{nq} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} 1 & \cdots & 0 & & 0 & \cdots & 0 \\ 0 & \cdots & c & \cdots & -s & \cdots & 0 \\ \vdots & & & & & & \\ 0 & \cdots & s & \cdots & c & \cdots & 0 \\ \vdots & & & & & & \\ 0 & & 0 & \cdots & 0 & \cdots & 1 \end{bmatrix}$$

when $c = C\theta$, $s = S\theta$

$$(R^T A)_{pp} = c\, a_{pp} + s\, a_{qp}$$

$$(R^T A)_{pq} = c\, a_{pq} + s\, a_{qq}$$

and so

$$\bar{a}_{pq} = (c\, a_{pp} + s\, a_{qp})(-s) + (c\, a_{pq} + s\, a_{qq})\, c$$

$$= - cs\, a_{pp} - s^2 a_{qp} + c^2 a_{pq} + cs\, a_{qq}$$

Since $A$ is symmetric,

$$\bar{a}_{pq} = (c^2 - s^2) a_{pq} + cs\,(a_{qq} - a_{pp}) = \tfrac{1}{2} S 2\theta\,(a_{qq} - a_{pp}) + C 2\theta\, a_{pq}$$

Hence; if

$$\theta = \begin{cases} \tfrac{1}{2}\tan^{-1}\left(\dfrac{2 a_{pq}}{a_{pp} - a_{qq}}\right) & a_{pp} \neq a_{qq} \\[2mm] \pm\dfrac{\Pi}{4} & a_{pp} = a_{qq} \end{cases} \qquad (1)$$

then the similarity tr. reduces $\bar{a}_{pq}$ to <u>Zero</u>.

Unfortunately, as the calculation progresses, subsequent trs. will probably change the values of elements previously reduced to zero.

Thus the method becomes an iterative one in which we construct the sequence of matrices

$$A^{(k+1)} = R^T (P,q) \, A^{(k)} \, R(P,q) \qquad k = 0, 1, 2, \dots$$

with $A^{(0)} = A$.

We take $\theta$ as:

$$-\frac{\pi}{4} \leq \theta \leq \frac{\pi}{4} \qquad (\text{if } a_{pp} \neq a_{qq}) \qquad \text{satisfying eqn(1)}$$

$$\theta = (\text{sign of } a_{pq}^{(k)}) \frac{\pi}{4} \qquad (\text{if } a_{pp} = a_{qq})$$

For each value of $k$, we have to decide which off-diagonal element $a_{pq}$ is to be reduced to zero.
It seems reasonable to choose the element of maximum modulus and this gives the standard Jacobi method.

These additional restrictions ensure that the iteration tends to a fixed diagonal matrix.

Because of the symmetry, only the elements above the diagonal need be considered.

$$R = R_1 R_2 R_3 \cdots R_n$$

**Ex.**

$$A = \begin{pmatrix} 10 & 3 & 2 \\ 3 & 5 & 1 \\ 2 & 1 & 0 \end{pmatrix} \equiv A^{(0)}$$

The element of $A^{(0)}$ with max. modulus and above the diagonal is in position $(p,q) = (1,2)$

$$R(1,2) = \begin{pmatrix} c\theta & -s\theta & 0 \\ s\theta & c\theta & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

$$\theta = \frac{1}{2} \tan^{-1}\left( \frac{2a_{12}}{a_{11} - a_{22}} \right) = \frac{1}{2} \tan^{-1}\left( \frac{2 \times 3}{10 - 5} \right) = 0.43803$$

$$\longrightarrow \quad \begin{cases} s\theta = 0.42416 \\ c\theta = 0.90559 \end{cases}$$

$$A^{(1)} = R^T A R = \begin{pmatrix} c\theta & s\theta & 0 \\ -s\theta & c\theta & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 10 & 3 & 2 \\ 3 & 5 & 1 \\ 2 & 1 & 0 \end{pmatrix} \begin{pmatrix} c\theta & -s\theta & 0 \\ s\theta & c\theta & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

$$= \begin{pmatrix} 11.40518 & 0 & 0 \\ 0 & 3.59489 & 0.05727 \\ 2.23534 & 0.05727 & 0 \end{pmatrix}$$

Next : we rotate in the $(1,3)$ plane using

$$R(1,3) = \begin{pmatrix} c\theta & 0 & -s\theta \\ 0 & 1 & 0 \\ s\theta & 0 & c\theta \end{pmatrix}$$

$$\theta = \frac{1}{2} \tan^{-1}\left( \frac{2a_{13}}{a_{11} - a_{22}} \right) = \frac{1}{2} \tan^{-1}\left( \frac{2 \times 7.23534}{11.40518 - 0} \right) = 0.18679$$

$$A^{(2)} = \begin{pmatrix} 11.82777 & 0.01064 & 0 \\ 0.01064 & 3.59489 & 0.05627 \\ 0 & 0.05627 & -0.42246 \end{pmatrix}$$

Note that the second _iteration_ has destroyed the zero elements in position (1,2) and (2,1).

$$R(2,3) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & C\theta & -S\theta \\ 0 & S\theta & C\theta \end{pmatrix} \qquad \theta = \frac{1}{2} \tan^{-1} \left( \frac{2 \times 0.05627}{3.59481 - (-0.42246)} \right)$$

$$= 0.01401$$

giving;

$$A^{(3)} = \begin{pmatrix} 11.82777 & 0.01064 & 0 \\ 0.01064 & 3.59567 & 0 \\ 0 & 0 & -0.42325 \end{pmatrix}$$

Finally; $R(1,2)$, with $\theta = 0.00129$

$$A^{(4)} = \begin{pmatrix} 11.82780 & 0 & 0 \\ 0 & 3.59567 & 0 \\ 0 & 0 & -0.42325 \end{pmatrix} \qquad \text{(tollerance: 0.0001)}$$

$$\lambda_1 = 11.8278 \qquad \lambda_2 = 3.5957 \qquad \lambda_3 = -0.4233$$

Since $\quad D = R^T A R \qquad$ where $R = R_1 R_2 \cdots R_n$

$\rightarrow AR = RD$

$\longrightarrow$ Eigenvectors of $A$ are the columns of $R$

$$R = R(1,2) \, R(1,3) \, R(2,3) \, R(1,2)$$

$$= \begin{pmatrix} 0.88929 & -0.42762 & -0.16222 \\ 0.41795 & 0.90386 & -0.09145 \\ 0.18573 & 0.01226 & 0.98251 \end{pmatrix}$$

The elements of $\bar{A}$ are:

$$\bar{a}_{ip} = a_{ip} \, c + a_{iq} \, s = \bar{a}_{pi} \left.\begin{array}{c} \\ \\ \end{array}\right\} \; i \neq p, q$$

$$\bar{a}_{iq} = a_{ip}(-s) + a_{iq} \, c = \bar{a}_{qi}$$

$$\bar{a}_{pp} = (c \, a_{pp} + s \, a_{qp}) \, c + (c \, a_{pq} + s \, a_{qq}) \, s$$

$$= c^2 a_{pp} + 2 \, cs \, a_{pq} + s^2 a_{qq}$$

$$\bar{a}_{qq} = (-s \, a_{pp} + c \, a_{qp})(-s) + (-s \, a_{pq} + c \, a_{qq}) \, c$$

$$= s^2 a_{pp} - 2 \, cs \, a_{pq} + c^2 a_{qq}$$

$$\bar{a}_{pq} = \bar{a}_{qp} = 0$$

$$\bar{a}_{ij} = a_{ij} \qquad i, j \neq p, q$$