

NEW YORK CQF WORKSHOP

- Mathematical Methods: The Heat Equation; Laplace Transforms and Applications.
- Introduction to Numerical Methods: Numerical Integration; Root Finding; Numerical Linear Algebra; Finite Difference Method

Mathematical Methods

The Heat Equation

We have seen on numerous occasions the one dimensional heat/diffusion equation

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2}$$

for the unknown function of the form $u(x, t) = t^\alpha \phi\left(\frac{x}{t^\beta}\right)$. The corresponding solution derived using the similarity reduction technique is the *fundamental solution*

$$u(x, t) = \frac{1}{2\sqrt{\pi t}} \exp\left(-\frac{x^2}{4t}\right).$$

Some books refer to this as a *source solution*.

Let's consider the following integral

$$\lim_{t \rightarrow 0} \int_{-\infty}^{\infty} u(y, t) f(y) dy$$

which can be simplified by the substitution

$$s = \frac{y}{2\sqrt{t}} \implies 2\sqrt{t} ds = dy$$

to give

$$\lim_{t \rightarrow 0} \frac{1}{2\sqrt{\pi t}} \int_{-\infty}^{\infty} \exp(-s^2) f(2\sqrt{t}s) 2\sqrt{t} ds.$$

In the limiting process we get

$$\begin{aligned} f(0) \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} \exp(-s^2) ds &= f(0) \frac{1}{\sqrt{\pi}} \sqrt{\pi} \\ &= f(0). \end{aligned}$$

Hence

$$\lim_{t \rightarrow 0} \int_{-\infty}^{\infty} u(y, t) f(y) dy = f(0).$$

A slight extension of the above shows that

$$\lim_{t \rightarrow 0} \int_{-\infty}^{\infty} u(x-y, t) f(y) dy = f(x),$$

where

$$u(x-y, t) = \frac{1}{2\sqrt{\pi t}} \exp\left(-\frac{(x-y)^2}{4t}\right).$$

Let's derive the result above. As earlier we begin by writing $s = \frac{x-y}{2\sqrt{t}} \implies y = x - 2\sqrt{t}s$ and hence $dy = -2\sqrt{t}ds$. Under this transformation the limits are

$$y = \infty \longrightarrow s = -\infty$$

$$y = -\infty \longrightarrow s = \infty$$

$$\begin{aligned} & \frac{1}{2\sqrt{\pi t}} \int_{-\infty}^{\infty} \exp(-s^2) f(x - 2\sqrt{t}s) (-2\sqrt{t}ds) ds \\ &= \frac{1}{\sqrt{\pi}} \int_{-\infty}^{\infty} \exp(-s^2) f(x - 2\sqrt{t}s) ds \end{aligned}$$

and $\lim_{t \rightarrow 0} f(x - 2\sqrt{t}s) = f(x)$, so

$$\begin{aligned} &= f(x) \frac{1}{\sqrt{\pi}} \sqrt{\pi} \\ &= f(x) \end{aligned}$$

$$\lim_{t \rightarrow 0} \int_{-\infty}^{\infty} u(x - y, t) f(y) dy = f(x).$$

Since the heat equation is a constant coefficient PDE, if $u(x, t)$ satisfies it, then $u(x - y, t)$ is also a solution for any y .

Recall what it means for an equation to be linear:

Since the heat equation is linear,

1. if $u(x - y, t)$ is a solution, so is a multiple $f(y) u(x - y, t)$

2. we can add up solutions. Since $f(y) u(x - y, t)$ is a solution for any y , so too is the integral

$$\int_{-\infty}^{\infty} u(x - y, t) f(y) dy.$$

Recall, adding can be done in terms of an integral. So we we can summarize by specifying the following initial value problem

$$\begin{aligned} \frac{\partial u}{\partial t} &= \frac{\partial^2 u}{\partial x^2} \\ u(x, 0) &= f(x) \end{aligned}$$

which has a solution

$$u(x, t) = \frac{1}{2\sqrt{\pi t}} \int_{-\infty}^{\infty} \exp\left(-\frac{(x - y)^2}{4t}\right) f(y) dy.$$

This satisfies the initial condition at $t = 0$ because we have shown that at that point the value of this integral is $f(x)$. Putting $t < 0$ gives a non-existent solution, i.e. the integrand will blow up.

Example Consider the IVP

$$\begin{aligned}\frac{\partial u}{\partial t} &= \frac{\partial^2 u}{\partial x^2} \\ u(x, 0) &= \begin{cases} 0 & \text{if } x > 0 \\ 1 & \text{if } x < 0 \end{cases}\end{aligned}$$

We can write down the solution as

$$\begin{aligned}u(x, t) &= \frac{1}{2\sqrt{\pi t}} \int_{-\infty}^{\infty} \exp\left(-\frac{(x-y)^2}{4t}\right) u(y, 0) dy \\ &= \frac{1}{2\sqrt{\pi t}} \int_{-\infty}^0 \exp\left(-\frac{(x-y)^2}{4t}\right) \cdot 1 dy\end{aligned}$$

put

$$s = \frac{y-x}{\sqrt{4t}}$$

$\int_{-\infty}^0$ becomes $\int_{-\infty}^{\frac{-x}{2\sqrt{t}}}$

$$\begin{aligned} & \frac{1}{2\sqrt{\pi t}} \int_{-\infty}^{\frac{-x}{2\sqrt{t}}} \exp(-s^2/2) \sqrt{2t} ds \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\frac{-x}{2\sqrt{t}}} \exp(-s^2/2) ds \\ &= N\left(\frac{-x}{2\sqrt{t}}\right) \end{aligned}$$

So we have expressed the solution in terms of the CDF.

This can also be solved by using the substitution

$$\hat{s} = \frac{-(y - x)}{2\sqrt{t}} \longrightarrow -dy = 2\sqrt{t}d\hat{s}$$

$$\int_{-\infty}^0 \text{ becomes } \int_{\infty}^{\frac{x}{2\sqrt{t}}}$$

$$\begin{aligned} & -\frac{1}{2\sqrt{\pi t}} \int_{\infty}^{\frac{x}{2\sqrt{t}}} \exp(-\hat{s}^2) 2\sqrt{t} d\hat{s} \\ &= \frac{1}{2} \frac{2}{\sqrt{\pi}} \int_{\frac{x}{2\sqrt{t}}}^{\infty} \exp(-\hat{s}^2) d\hat{s} = \\ &= \frac{1}{2} \operatorname{erfc} \left(\frac{x}{2\sqrt{t}} \right) \end{aligned}$$

so now we have a solution in terms of the *complimentary error function*.

The Laplace Transform

Given a function defined $\forall t > 0$ then the *Laplace Transform* (LT) of $f(t)$, written $\mathcal{L}\{f(t)\}$ is defined by

$$\mathcal{L}\{f(t)\} = \int_0^{\infty} e^{-st} f(t) dt = F(s)$$

provided this integral exists.

$\mathcal{L}\{f(t)\}$ is a function of s which we denote by $F(s)$.

<u>Lower case letters</u>	<u>Upper case letters</u>
$f(t)$	$F(s)$
$g(t)$	$G(s)$
$h(t)$	$H(s)$

Example: If $f(t) = 1$ then $\mathcal{L}\{1\} = \int_0^\infty 1 \cdot e^{-st} dt$. Look at rhs

$$\begin{aligned} \lim_{T \rightarrow \infty} \int_0^T e^{-st} dt &= \lim_{T \rightarrow \infty} \left[-\frac{1}{s} e^{-st} \right]_0^T \quad (s \neq 0) \\ &= \lim_{T \rightarrow \infty} \left(\frac{1}{s} [1 - e^{-sT}] \right) \\ &\rightarrow \frac{1}{s} \quad \text{if } s > 0 \end{aligned}$$

The limit does not exist for $s < 0$.

What about $s = 0$?

$$\int_0^\infty e^{0t} dt = \int_0^\infty 1 dt$$

which does not exist. So for $f(t) = 1$ $F(s) = \frac{1}{s}$ ($s > 0$) and L.T. does not exist for $s \leq 0$.

Laplace Transforms of Standard Functions

$$1. \mathcal{L}\{1\} = \frac{1}{s} \quad (s > 0)$$

$$2. \mathcal{L}\{e^{at}\} = \int_0^{\infty} e^{at} e^{-st} dt =$$

$$\begin{aligned} \lim_{T \rightarrow \infty} \int_0^T e^{-(s-a)t} dt &= \lim_{T \rightarrow \infty} \left[\frac{1}{s-a} \left(1 - e^{-(s-a)T} \right) \right] \\ &= \frac{1}{s-a} \quad (s > a) \end{aligned}$$

$$3. \mathcal{L}\{t\} = \int_0^{\infty} t e^{-st} dt \quad \text{integration by parts}$$

$$= -\frac{t}{s} e^{-st} \Big|_0^{\infty} + \frac{1}{s} \int_0^{\infty} e^{-st} dt$$

Now $te^{-st} \rightarrow 0$ as $t \rightarrow \infty \because s > 0$ therefore

$$\begin{aligned}\mathcal{L}\{t\} &= \frac{1}{s} \int_0^{\infty} e^{-st} dt = \frac{1}{s} \left[-\frac{e^{-st}}{s} \right]_0^{\infty} \\ &= \frac{1}{s^2}\end{aligned}$$

4. $\mathcal{L}\{t^n\} = \int_0^{\infty} t^n e^{-st} dt \quad (s > 0)$, and put $u = st$ to give
 $dt = \frac{1}{s} du$, hence

$$\int_0^{\infty} t^n e^{-st} dt = \int_0^{\infty} \left(\frac{u}{s}\right)^n e^{-u} \frac{du}{s} = \frac{1}{s^{n+1}} \int_0^{\infty} u^n e^{-u} du$$

where

$$\int_0^{\infty} u^n e^{-u} du = \Gamma(n+1)$$

so we get

$$\frac{\Gamma(n+1)}{s^{n+1}} = \frac{n!}{s^{n+1}}$$

hence

$$\mathcal{L}\{t^n\} = \frac{n!}{s^{n+1}}.$$

5. Similarly $\mathcal{L}\{t^\alpha\} = \frac{\Gamma(\alpha+1)}{s^{\alpha+1}} \quad (s > 0, \alpha > -1)$

$$\begin{aligned} \mathcal{L}\{t^{-1/2}\} &= \frac{\Gamma(-1/2+1)}{s^{-1/2+1}} = \frac{\Gamma(1/2)}{s^{1/2}} \\ &= \frac{\sqrt{\pi}}{s^{1/2}} = \sqrt{\frac{\pi}{s}} \end{aligned}$$

$$\mathcal{L}\{t^{1/2}\} = \frac{\Gamma(1/2+1)}{s^{1/2+1}} = \frac{\Gamma(3/2)}{s^{3/2}}$$

recall $\Gamma(\alpha + 1) = \alpha\Gamma(\alpha)$, so we can write $\Gamma(3/2) = (1/2)\Gamma(1/2) = (1/2)\sqrt{\pi}$

$$= \frac{\sqrt{\pi}}{2s^{3/2}} = \frac{1}{2s}\sqrt{\frac{\pi}{s}}.$$

6. $\mathcal{L}\{\sin t\} = \int_0^\infty e^{-st} \sin kt \, dt$

Using Eulers Identity is fairly quick.

$$\begin{aligned}\int_0^{\infty} e^{-st} \sin kt .dt &= \int_0^{\infty} e^{-st} \operatorname{Im} e^{ikt} dt \\&= \operatorname{Im} \int_0^{\infty} e^{-(s-ik)t} dt \\&= \operatorname{Im} \frac{1}{-(s-ik)} e^{-(s-ik)t} \Big|_0^{\infty} = \operatorname{Im} \frac{1}{(s-ik)} \\&= \operatorname{Im} \frac{1}{(s-ik)} \frac{(s+ik)}{(s+ik)} = \operatorname{Im} \frac{(s+ik)}{s^2+k^2} \\&= \frac{k}{k^2 + s^2}\end{aligned}$$

As an exercise repeat this for the Laplace Transform of $\cos kt$.

Simple Properties

$$1. \mathcal{L}\{f(t) + g(t)\} = \mathcal{L}\{f(t)\} + \mathcal{L}\{g(t)\}$$

$$2. \mathcal{L}\{\lambda f(t)\} = \lambda \mathcal{L}\{f(t)\}$$

$$3. \mathcal{L}\{f(t) \cdot g(t)\} \neq \mathcal{L}\{f(t)\} \cdot \mathcal{L}\{g(t)\}$$

A more elegant way of obtaining the result for trigonometric functions is to consider

$$\begin{aligned} \mathcal{L}\{e^{ikt}\} &= \frac{1}{s - ik} = \frac{s + ik}{s^2 + k^2} \equiv \mathcal{L}\{\cos kt + i \sin kt\} \\ &= \frac{s}{s^2 + k^2} + ik \frac{1}{s^2 + k^2} \end{aligned}$$

Hence

$$\mathcal{L}\{\cos kt\} = \frac{s}{s^2 + k^2} \quad , \quad \mathcal{L}\{\sin kt\} = \frac{k}{s^2 + k^2}$$

Shift Theorem

If $\mathcal{L}\{f(t)\} = F(s)$, then

$$\mathcal{L}\{e^{at}f(t)\} = F(s - a)$$

or $f(t) = t^n$ with $F(s) = \frac{n!}{s^{n+1}}$, then

$$\mathcal{L}\{e^{at}t^n\} = F(s - a) = \frac{n!}{(s - a)^{n+1}}$$

Table of Standard Forms

$f(t)$	$F(s)$
1	$\frac{1}{s}$
$e^{\pm at}$	$\frac{1}{s \mp a}$
t^n	$\frac{n!}{s^{n+1}}$
t^α	$\frac{\Gamma(\alpha+1)}{s^{\alpha+1}}$
$\sin kt$	$\frac{k}{k^2 + s^2}$
$\cos kt$	$\frac{s}{s^2 + k^2}$
$t \sin kt$	$\frac{2ks}{(k^2 + s^2)^2}$
$t \cos kt$	$\frac{s^2 - k^2}{(s^2 + k^2)^2}$

The Inverse Laplace Transform

Given $F(s)$:

$$\mathcal{L}^{-1}(F(s)) = \text{the function } f(t) \text{ for which } \mathcal{L}\{f(t)\} = F(s).$$

Therefore we can read the table of Laplace Transforms in two ways

$$\begin{array}{ccc} f(t) & & F(s) \\ \mathcal{L} & & \\ \longrightarrow & & \\ & \mathcal{L}^{-1} & \\ & \longleftarrow & \end{array}$$

For example reading the earlier table from right to left we have

$$\begin{aligned} \mathcal{L}^{-1}\left(\frac{1}{s}\right) &= 1 \\ \mathcal{L}^{-1}\left(\frac{1}{s^{n+1}}\right) &= \frac{t^n}{n!} \end{aligned}$$

Example 1 : Evaluate $\mathcal{L}^{-1} \left(\frac{1}{s^2 + 2s + 2} \right)$.

The structure should suggest some variation of

$$\mathcal{L}^{-1} \left(\frac{1}{s^2} \right) \quad \text{or} \quad \mathcal{L}^{-1} \left(\frac{1}{s^2 + k^2} \right).$$

Now $s^2 + 2s + 2 = (s + 1)^2 + 1^2$. Hence

$$\mathcal{L}^{-1} \left(\frac{1}{s^2 + 2s + 2} \right) \equiv \mathcal{L}^{-1} \left(\frac{1}{(s + 1)^2 + 1^2} \right) \quad \text{c.f.} \quad \mathcal{L}^{-1} \left(\frac{1}{s^2 + k^2} \right)$$

with s replaced by $s + 1$.

From the table:

$$\mathcal{L} \left\{ e^{at} \sin kt \right\} = \frac{k}{(s + a)^2 + k^2}$$

So $a = -1$, $k = 1$

$$\therefore \mathcal{L} \left\{ e^{-t} \sin t \right\} = \frac{1}{(s+1)^2 + 1^2} \text{ and } \mathcal{L}^{-1} \left(\frac{1}{s^2 + 2s + 2} \right) = e^{-t} \sin t$$

Example 2 : $\mathcal{L}^{-1} \left(\frac{1}{(s+1)(s+2)} \right)$

Using partial fractions this becomes $\mathcal{L}^{-1} \left\{ \frac{1}{(s+1)} - \frac{1}{(s+2)} \right\}$

$$= \mathcal{L}^{-1} \left\{ \frac{1}{(s+1)} \right\} - \mathcal{L}^{-1} \left\{ \frac{1}{(s+2)} \right\}$$

We know

$$\mathcal{L} \left\{ e^{at} \right\} = \frac{1}{(s-a)} \therefore \mathcal{L} \left\{ e^{-t} \right\} = \frac{1}{s+1} \quad \& \quad \mathcal{L} \left\{ e^{-2t} \right\} = \frac{1}{s+2}$$

Hence we have

$$e^{-t} - e^{-2t}$$

Use of Laplace Transforms For Solving Differential Equations

Laplace Transforms provide us with a very useful technique for solving IVP's.

This is based on two results (which are also theorems):

$$1. \mathcal{L} \{f'(t)\} = sF(s) - f(0)$$

$$2. \mathcal{L} \{f''(t)\} = s^2 F(s) - sf(0) - f'(0)$$

$\mathcal{L} \{f'(t)\} = \int_0^{\infty} e^{-st} f'(t) dt$ and we integrate by parts to give

$$f(t) e^{-st} \Big|_0^{\infty} + s \int_0^{\infty} e^{-st} f(t) dt = -f(0) + sF(s)$$

$\mathcal{L}\{f''(t)\} = \int_0^\infty e^{-st} f''(t) dt$ which requires integration by parts twice (and includes use of the earlier result)

$$\begin{aligned} f'(t) e^{-st} \Big|_0^\infty + s \int_0^\infty e^{-st} f'(t) dt &= -f'(0) + s(-f(0) + sF(s)) \\ &= s^2 F(s) - sf(0) - f'(0) \end{aligned}$$

We can use these to transform an IVP such as

$$\left. \begin{aligned} ay'' + by' + cy &= f(t) \\ y(0) &= \alpha; \quad y'(0) = \beta \end{aligned} \right\}$$

into an algebraists problem by taking Laplace Transforms of both sides of the ode. Let $Y(s) = \mathcal{L}\{y(t)\}$ then

$$a(s^2 Y(s) - sy(0) - y'(0)) + b(sY(s) - y(0)) + cY(s) = F(s)$$

and simplifying

$$Y(s) = \frac{F(s) + ay'(0) + (as + b)y(0)}{as^2 + bs + c}.$$

The final step in the working entails taking inverse Laplace Transforms to have $Y(s) \rightarrow y(t)$.

Example: Solve the IVP

$$\left. \begin{array}{l} y'' + 2y' + y = 0 \\ y(0) = 2; \quad y'(0) = -1. \end{array} \right\}$$

Start by taking LT of the whole equation, with $Y(s) = \mathcal{L}\{y(t)\}$

$$\begin{aligned} \mathcal{L}\{y'' + 2y' + y\} &= \mathcal{L}\{0\} \\ \mathcal{L}\{y''\} + 2\mathcal{L}\{y'\} + \mathcal{L}\{y\} &= 0 \end{aligned}$$

The initial condition simplifies the problem considerably to give

$$\begin{aligned} s^2 Y(s) - sY(0) - Y'(0) + 2(sY(s) - Y(0)) + Y(s) &= 0 \\ Y(s)(s^2 + 2s + 1) - 2s - 3 &= 0 \\ Y(s)(s^2 + 2s + 1) &= 2s + 3 \\ Y(s) &= \frac{2s + 3}{(s + 1)^2} \end{aligned}$$

The right hand side can be written as

$$\frac{2s + 3}{(s^2 + 2s + 1)} \equiv \frac{A}{(s + 1)} + \frac{B}{(s + 1)^2}$$

for $A = 2$, $B = 1$ to give

$$Y(s) = \frac{2}{(s + 1)} + \frac{1}{(s + 1)^2}$$

Now taking inverse LT's

$$\mathcal{L}^{-1}\{Y(s)\} = \mathcal{L}^{-1}\left\{\frac{2}{(s + 1)} + \frac{1}{(s + 1)^2}\right\}$$

therefore

$$y(t) = 2\mathcal{L}^{-1}\left\{\frac{1}{(s + 1)}\right\} + \mathcal{L}^{-1}\left\{\frac{1}{(s + 1)^2}\right\}$$

Use the Shift Theorem for the second term

$$\mathcal{L}\{e^{at}f(t)\} = F(s - a)$$

we get a te^{-t} .

So the solution is

$$y(t) = 2e^{-t} + te^{-t}$$

Example: Solve the IVP

$$\left. \begin{aligned} y'' + 4y' + 6y &= 1 + e^{-t} \\ y(0) = y'(0) &= 0. \end{aligned} \right\}$$

Start by taking LT of the whole equation, with $Y(s) = \mathcal{L}\{y(t)\}$

$$\mathcal{L}\{y'' + 4y' + 6y\} = \mathcal{L}\{1 + e^{-t}\}$$

$$\mathcal{L}\{y''\} + 4\mathcal{L}\{y'\} + 6\mathcal{L}\{y\} = \mathcal{L}\{1\} + \mathcal{L}\{e^{-t}\}$$

The initial condition simplifies the problem considerably to give

$$\begin{aligned} s^2 Y(s) + 4sY(s) + 6Y(s) &= \left(\frac{1}{s} + \frac{1}{s+1}\right) \\ Y(s) &= \frac{1}{s^2 + 4s + 6} \left(\frac{1}{s} + \frac{1}{s+1}\right) \end{aligned}$$

Now taking inverse LT's

$$\mathcal{L}^{-1} \{Y(s)\} = \mathcal{L}^{-1} \left\{ \frac{1}{s(s^2 + 4s + 6)} + \frac{1}{(s^2 + 4s + 6)(s + 1)} \right\}$$

after firstly using partial fractions, i.e. decompose as

$$\frac{1}{6s} + \frac{1}{3(s + 1)} - \frac{s + 2}{2(s^2 + 4s + 6)} - \frac{2}{3(s^2 + 4s + 6)}$$

therefore

$$\begin{aligned} y(t) = & \frac{1}{6} \mathcal{L}^{-1} \left\{ \frac{1}{s} \right\} + \frac{1}{3} \mathcal{L}^{-1} \left\{ \frac{1}{(s + 1)} \right\} - \frac{1}{2} \mathcal{L}^{-1} \left\{ \frac{s + 2}{(s^2 + 4s + 6)} \right\} \\ & - \frac{2}{3} \mathcal{L}^{-1} \left\{ \frac{1}{(s^2 + 4s + 6)} \right\} \end{aligned}$$

We know $s^2 + 4s + 6 \equiv (s + 2)^2 + 2$. So use the Shift Theorem

$$\mathcal{L} \{e^{at} f(t)\} = F(s - a)$$

$$\mathcal{L}\{\cos kt\} = \frac{s}{s^2 + k^2}; \quad \mathcal{L}\{\sin kt\} = \frac{k}{s^2 + k^2}; \quad \mathcal{L}\{1\} = \frac{1}{s}$$

So the solution is

$$y(t) = \frac{1}{6} + \frac{1}{3}e^{-t} - \frac{1}{2}e^{-2t} \cos \sqrt{2}t - \frac{2}{3}e^{-2t} \sin \sqrt{2}t.$$

Theoretically, we can also extend this to an n^{th} order IVP (with constant coefficients) by finding $\mathcal{L}\{f^{(n)}(t)\}$.

Introduction to Numerical Methods

Numerical Integration

The need frequently arises for evaluating the definite integral of a function that has no explicit antiderivative or whose antiderivative is not easily obtainable.

Recall the fundamental theorem of calculus

$$\int_a^b f(x) dx = F(b) - F(a)$$

Consider the usual (Riemann) definition of a definite integral

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

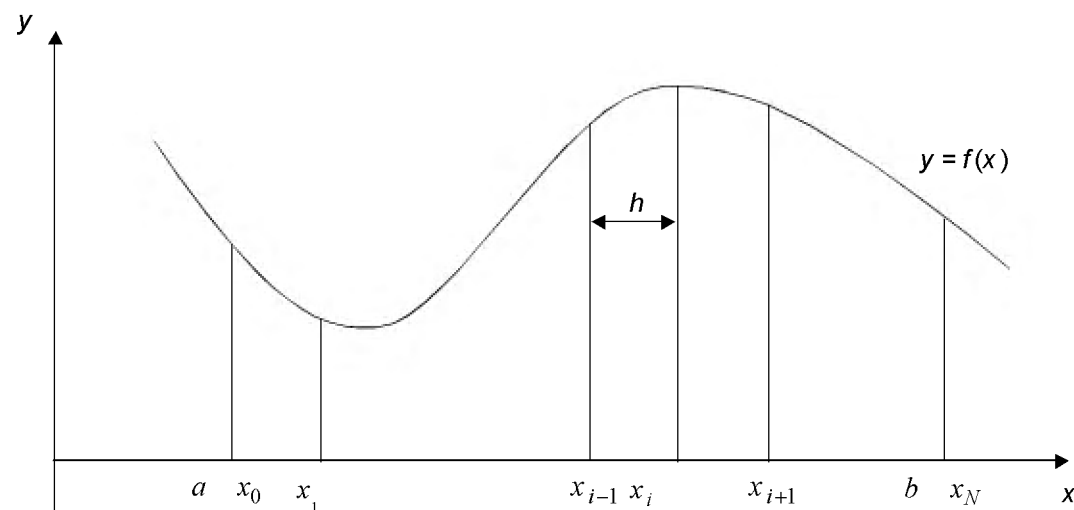
as the area under the curve where the curve is the graph of $f(x)$ plotted against x .

Assuming f is a "well behaved" function on $[a, b]$, we can define this in many different ways (which all lead to the same value for the definite integral).

Divide the interval $[a, b]$ into N intervals with end points $x_0 = a < x_1 < x_2 < \dots < x_{N-1} < x_N = b$.

Then $h = (b - a) / N$ is the constant length of an interval $x_{i+1} - x_i$, which tends to zero as $N \rightarrow \infty$.

$x_i = a + ih$ for each $i = 0, 1, \dots, N$.



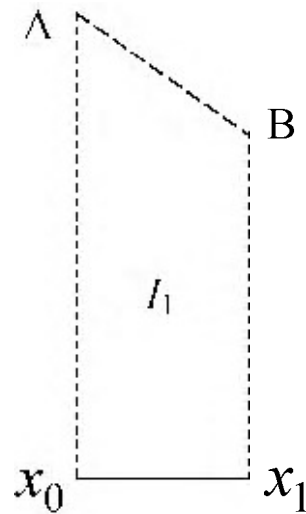
We could approximate the definite integral by *The trapezium rule*;

$$\int_a^b f(x) dx = \frac{h}{2} \left[f(a) + 2 \sum_{i=1}^{N-1} f(x_i) + f(b) \right].$$

Each strip is labelled I_i for $1 \leq i \leq N$. So the area I under the curve between $x = a$ and $x = b$ is given by

$$I = I_1 + I_2 + \dots + I_N \quad (1)$$

Put $f(x_0) = A$ and $f(x_1) = B$. Consider the first strip (a trapezium). That is, points A and B are joined with a straight line and the area of the trapezium ABx_1x_0 is assumed to be a good approximation to the area under the curve.



The area ABx_1x_0 is simply

$$I_1 \cong \frac{h}{2}(f_0 + f_1)$$

An approximation to

$$I_1 = \int_{x_0}^{x_1} f(x) dx$$

is now given by

$$I_1 \cong \frac{h}{2} (f_0 + f_1) \tag{2}$$

We can extend this to consider N strips. We can write (1) in integral form as

$$\int_{x_0}^{x_N} f(x) dx = \int_{x_0}^{x_1} f(x) dx + \int_{x_1}^{x_2} f(x) dx + \dots + \int_{x_{N-1}}^{x_N} f(x) dx$$

From (2)

$$I_1 = \frac{h}{2} (f_0 + f_1)$$

and by extension

$$I_2 = \frac{h}{2} (f_1 + f_2)$$

$$I_3 = \frac{h}{2} (f_2 + f_3), \dots$$

So summing these quantities gives the following $I =$

$$\frac{h}{2} [(f_0 + f_1) + (f_1 + f_2) + (f_2 + f_3) + \dots (f_{N-1} + f_N)]$$

OR

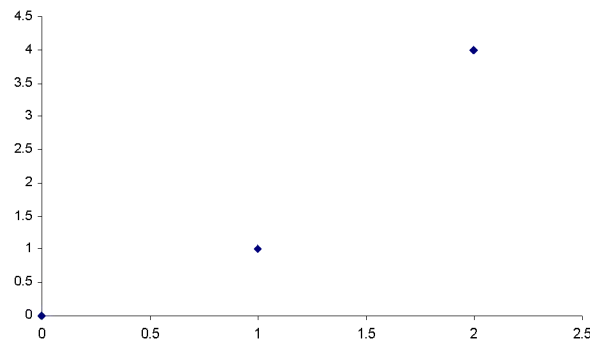
$$I = \frac{h}{2} [f_0 + 2f_1 + 2f_2 + \dots + 2f_{N-1} + f_N] \quad (3)$$

This is the *Composite Trapezoidal Rule*.

Lagrangian Interpolation

In general the term **interpolation** is the term used to describe the process in which a function, usually a polynomial function is constructed to pass through a given set of data points.

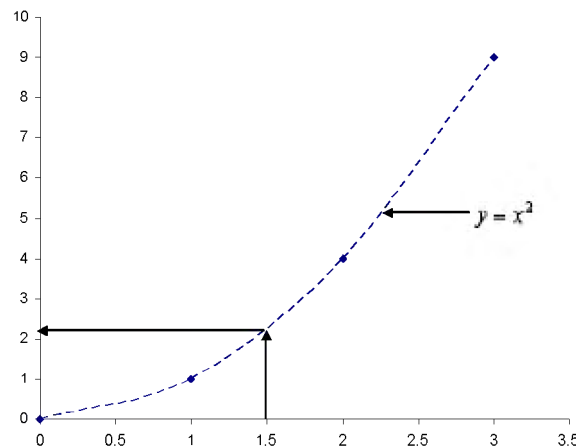
As an example consider the following data points (x_i, y_i) as follows:
 $(0, 0)$; $(1, 1)$; $(2, 4)$



In this example the function $y = x^2$ passes through these 3 data points.

Why is interpolation important and how is it used?

Suppose an estimate is required for the value of y at some x other than $x = 0, 1, 2$. In this case, the interpolating polynomial $y = x^2$ can be used to estimate such intermediate values, e.g. $x = 1.5$, which gives $y = 2.25$



Given the data points $(0, 0)$; $(1, 1)$; $(2, 4)$ the values of y corresponding to any value of x between $x = 0$ and $x = 2$ can be estimated. This is known as *interpolation*.

Lagrange Polynomials

For a given set of $n + 1$ nodes x_i ,

$$\begin{array}{ccccccc} i & 0 & 1 & 2 & \dots\dots\dots & n \\ x_i & x_0 & x_1 & x_2 & \dots\dots\dots & x_n \\ y_i & y_0 & y_1 & y_2 & \dots\dots\dots & y_n \end{array}$$

the Lagrange polynomials are the $n + 1$ polynomials l_i defined by

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

where

$$l_i(x_j) = \delta_{ij} = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j \end{cases}$$

Then we define the interpolating polynomial as

$$p_n(x) = \sum_{i=0}^n y_i l_i(x)$$

If each Lagrange Polynomial is of degree at most n , then p_n also has this property. The Lagrange Polynomials can be characterized as follows:

By evaluating this product for each x_j , we see that this is indeed a characterization of the Lagrange Polynomials. Moreover, each polynomial is clearly the product of n "mononomials", and thus has degree no greater than n .

Example

Construct the polynomial interpolating the data

$$\begin{array}{cccc} x & 1 & \frac{1}{2} & 3 \\ y & 3 & -10 & 2 \end{array}$$

by using Lagrange Polynomials.

Solution:

$$\begin{array}{cccc} i & 0 & 1 & 2 \\ x_i & 1 & \frac{1}{2} & 3 \\ y_i & 3 & -10 & 2 \end{array}$$

We construct the Lagrange Polynomials:

So we start with $i = 0$, so need x_1, x_2

$$\begin{aligned} l_0(x) &= \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} \\ &= \frac{\left(x - \frac{1}{2}\right)(x - 3)}{\left(1 - \frac{1}{2}\right)(1 - 3)} \\ &= -\left(x - \frac{1}{2}\right)(x - 3) \end{aligned}$$

need x_0, x_2

$$\begin{aligned} l_1(x) &= \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} \\ &= \frac{(x - 1)(x - 3)}{\left(\frac{1}{2} - 1\right)\left(\frac{1}{2} - 3\right)} \\ &= \frac{4}{5}(x - 1)(x - 3) \end{aligned}$$

need x_0, x_1

$$\begin{aligned} l_2(x) &= \frac{(x - x_0)(x - x_1)}{(x_2 - x_1)(x_2 - x_1)} \\ &= \frac{(x - 1)\left(x - \frac{1}{2}\right)}{(3 - 1)\left(3 - \frac{1}{2}\right)} \\ &= \frac{1}{5}(x - 1)\left(x - \frac{1}{2}\right) \end{aligned}$$

Then the interpolating polynomial, in "Lagrange Form" is

$$\begin{aligned}
p_2(x) &= \sum_{i=0}^n y_i l_i(x) = y_0 l_0(x) + y_1 l_1(x) + y_2 l_2(x) \\
&= 3l_0(x) - 10l_1(x) + 2l_2(x) \\
&= -3\left(x - \frac{1}{2}\right)(x - 3) - 8(x - 1)(x - 3) \\
&\quad + \frac{2}{5}(x - 1)\left(x - \frac{1}{2}\right)
\end{aligned}$$

A useful check on the correctness of the calculations performed is

$$\sum_{i=0}^n l_i(x) = 1$$

Simpson's Rule

Simpson's Rule is obtained by considering two strips together, fitting a quadratic polynomial through the points on the curve and then integrating the quadratic

(diagram 1)

Let

$$I = \int_{x_0}^{x_2} f(x) dx \quad (2)$$

To find I first fit a quadratic to the points ABC . By Lagrangian Interpolation of degree 2

$$f(x) \cong p_2(x) = f_0 l_0 + f_1 l_1 + f_2 l_2 \quad (3)$$

$$\begin{aligned} l_0 &= \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)}; \quad l_1(x) = \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} \\ l_2(x) &= \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)} \end{aligned}$$

Now

$$(x_0 - x_1) = -h ; (x_0 - x_2) = -2h$$

$$(x_1 - x_0) = h ; (x_1 - x_2) = -h$$

$$(x_2 - x_1) = h ; (x_2 - x_0) = 2h$$

So from (3)

$$f(x) \cong p_2(x) = \frac{(x - x_1)(x - x_2)}{2h^2} f_0 - \frac{(x - x_0)(x - x_2)}{h^2} f_1 + \frac{(x - x_0)(x - x_1)}{2h^2} f_2$$

$f(x)$ is now substituted into (2) to give

$$I \cong \frac{1}{2h^2} \int_{x_0}^{x_2} \{ (x - x_1)(x - x_2) f_0 - 2(x - x_0)(x - x_2) f_1 + (x - x_0)(x - x_1) f_2 \} dx \quad (4)$$

(4) is integrated by parts to give

$$I \cong \frac{h}{3} (f_0 + 4f_1 + f_2) \quad (5)$$

This is *Simpson's Rule*.

Composite Form of Simpson's Rule

This is similar in principle to the Composite Trapezoidal Rule but here the region of integration is split into pairs of strips as shown in the diagram

(diagram 2)

$$I = \int_{x_0}^{x_{2n}} f(x) dx$$

is now split into pairs of strips such that

$$I = \int_{x_0}^{x_2} f(x) dx + \int_{x_2}^{x_4} f(x) dx + \dots + \int_{x_{2n-2}}^{x_{2n}} f(x) dx$$

Each integral on the right hand side of this expression above is now evaluated by Simpson's Rule.

This gives from (5) applied in the appropriate interval

$$I \cong \frac{h}{3} \left(\begin{array}{c} f_0 + 4f_1 + 2f_2 + 4f_3 + 2f_4 + \dots \\ + 2f_{2n-2} + 4f_{2n-1} + f_{2n} \end{array} \right)$$

so that

$$I \cong \frac{h}{3} \left(f_0 + f_{2n} + 4 \sum_{i=1}^n f_{2i-1} + 2 \sum_{i=1}^{n-1} f_{2i} \right) \quad (6)$$

This is the composite form of Simpson's Rule.

Numerical Linear Algebra

Definition An **upper triangular** matrix U has zero elements below the principal diagonal, i.e.

$$u_{ij} = 0 \quad \forall i > j.$$

For an upper triangular matrix U consider the system $U\underline{x} = \underline{b}$ in component form

$$\begin{array}{rcll} u_{11}x_1 + u_{12}x_2 + \dots + u_{1n}x_n & = & b_1 \\ u_{22}x_2 + \dots + u_{2n}x_n & = & b_2 \\ u_{33}x_3 + \dots + u_{3n}x_n & = & b_3 \\ \vdots & & \vdots \\ \vdots & & \vdots \\ \vdots & & \vdots \\ u_{n-1,n-1}x_{n-1} + u_{n-1,n}x_n & = & b_{n-1} \\ u_{nn}x_n & = & b_n \end{array}$$

The determinant of the matrix U , can now easily be obtained from

$$\det(U) = \prod_{i=1}^n u_{ii}.$$

Recall that U has an inverse if $\det(U) \neq 0$, hence for U to be non-singular requires the strict condition $u_{ii} \neq 0 \quad \forall i$. This linear system can now easily be solved by **back-substitution**. The algorithm for this scheme becomes

$$\begin{aligned} x_n &= \frac{b_n}{u_{nn}} \\ x_i &= \frac{1}{u_{ii}} \left(b_i - \sum_{j=i+1}^n u_{ij} x_j \right) \text{ for } i = n-1, n-2, \dots, 1 \end{aligned}$$

Definition A **lower triangular** matrix L has zero elements above the principal diagonal, i.e.

$$l_{ij} = 0 \quad \forall i < j.$$

Now suppose L is a lower triangular matrix and consider the system $L\underline{x} = \underline{b}$.

A *lower triangular* matrix L has zero elements above the main diagonal, i.e.

$$l_{ij} = 0 \quad \forall i < j.$$

Now suppose L is a lower triangular matrix and consider the system $L\underline{x} = \underline{b}$.

$$\begin{array}{rcl}
 l_{11}x_1 & & = b_1 \\
 l_{21}x_1 + l_{22}x_2 & & = b_2 \\
 \vdots & \vdots & \ddots \\
 \vdots & \vdots & \ddots \\
 \vdots & \vdots & \ddots \\
 l_{n1}x_1 + l_{n2}x_2 + \dots + l_{nn}x_n & = & b_n
 \end{array}$$

Again, because

$$\det(L) = \prod_{i=1}^n l_{ii},$$

L is non-singular iff $l_{ii} \neq 0 \quad \forall i$. We solve this system by **forward-substitution**. The scheme is given by

$$\begin{aligned} x_1 &= \frac{b_1}{l_{11}} \\ x_i &= \frac{1}{l_{ii}} \left(b_i - \sum_{j=1}^{i-1} l_{ij} x_j \right) \quad \text{for } i = 2, 3, \dots, n \end{aligned}$$

Definition An $n \times n$ matrix is called a *band matrix* if integers p and q exist such that $p, q \in (1, n)$ with the property that $a_{ij} = 0$ whenever $i + p \leq j$ or $j + q \leq i$. The band width w of a banded matrix is defined to be $w = p + q - 1$.

So for example, the matrix

$$A = \begin{pmatrix} 7 & 2 & 0 \\ 3 & 5 & -1 \\ 0 & -5 & -6 \end{pmatrix}$$

is a banded matrix with $p = q = 2$ and band width $w = 3$.

So the definition of a banded matrix forces such matrices to have all its non-zero elements to be clustered close to the main diagonal.

p is then called the *upper band-width* and q the *lower band-width*.

Solving equations with banded matrices can be much more efficient than with full matrices because the triangular factors in the LU decomposition are also banded. This results in greater efficiency in terms of storage and amount of computations required. Banded matrices arise

from the discretization of differential equations, because the finite difference method use to approximate derivatives only involve values at nearby mesh points.

From an applied mathematics perspective, the special case occurring when $p = q = 2$ and band width $w = 3$ is the most useful for us. These matrices are called *tridiagonal* since they have the form

$$\begin{pmatrix} a_{11} & a_{12} & 0 & \cdots & \cdots & 0 \\ a_{21} & a_{22} & a_{23} & \ddots & & \vdots \\ 0 & a_{32} & a_{33} & a_{34} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & 0 \\ \vdots & & \ddots & \ddots & \ddots & a_{n-1,n} \\ 0 & \cdots & \cdots & 0 & a_{n,n-1} & a_{nn} \end{pmatrix}$$

The LU Decomposition

It is often advantageous to think of Gaussian elimination as constructing a lower tridiagonal matrix L and an upper triangular matrix U , so that $LU = A$. To illustrate this method consider the following *tridiagonal* linear system

$$\begin{pmatrix} b_0 & 0 & 0 & \cdots & \cdots & 0 \\ a_1 & b_1 & 1 & & & \vdots \\ 0 & a_2 & \ddots & \ddots & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & 0 \\ \vdots & & & a_{n-1} & b_{n-1} & \vdots \\ 0 & \cdots & \cdots & 0 & a_n & b_n \end{pmatrix} \begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ y_{n-1} \\ y_n \end{pmatrix} = \begin{pmatrix} p_0 \\ p_1 \\ p_2 \\ \vdots \\ p_{n-1} \\ p_n \end{pmatrix} \quad (1)$$

We can write this in the form

$$\mathbf{M} \cdot \mathbf{y} = \mathbf{p}$$

and then consider the matrix factorisation $\mathbf{M} = \mathbf{L}\mathbf{U}$

$$= \begin{pmatrix} 1 & 0 & \cdots & \cdots & \cdots & 0 \\ l_1 & 1 & 0 & & & \vdots \\ 0 & l_2 & 1 & \ddots & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \vdots \\ \vdots & & & l_{n-1} & 1 & 0 \\ 0 & \cdots & \cdots & 0 & l_n & 1 \end{pmatrix} \begin{pmatrix} d_0 & u_0 & 0 & \cdots & \cdots & 0 \\ 0 & d_1 & u_1 & & & \vdots \\ 0 & 0 & d_2 & u_2 & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & 0 \\ \vdots & & & 0 & d_{n-1} & u_{n-1} \\ 0 & \cdots & \cdots & 0 & 0 & d_n \end{pmatrix}$$

$$= \begin{pmatrix} d_0 & u_0 & 0 & \cdots & \cdots & 0 \\ l_1 d_0 & l_1 u_0 + d_1 & u_1 & & & \vdots \\ 0 & l_2 d_1 & l_2 u_1 + d_2 & u_2 & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & 0 \\ \vdots & & & l_{n-1} d_{n-2} & l_{n-1} u_{n-2} + d_{n-1} & u_{n-1} \\ 0 & \cdots & \cdots & 0 & l_n d_{n-1} & l_n u_{n-1} + d_n \end{pmatrix}$$

where L is the lower triangular matrix above and U is the upper triangular matrix above.

Equating the elements of the original tridiagonal matrix M with the elements of the product matrix LU , we find that

$$d_0 = b_0 \tag{2}$$

$$u_i = \frac{a_i}{d_{i-1}}, \quad i = 0, 1, 2, \dots, n-1$$

$$l_i = \frac{a_i}{d_{i-1}}, \quad d_i = b_i - l_i u_{i-1}, \quad i = 1, 2, \dots, n$$

We then solve the original system given by (1) by solving two smaller problems. The tridiagonal system (1) is written in the form

$$\mathbf{M}.\mathbf{y} = (\mathbf{LU}).\mathbf{y} = \mathbf{L}(\mathbf{U}.\mathbf{y}) = \mathbf{p}$$

We introduce an intermediate vector $\mathbf{z} = \mathbf{U}\mathbf{y}$ so that (1) becomes

$$\mathbf{L}\mathbf{z} = \mathbf{p}, \quad \mathbf{U}\mathbf{y} = \mathbf{z}.$$

First we solve the problem

$$\mathbf{L}\mathbf{z} = \mathbf{p}$$

for the intermediate vector \mathbf{z} , i.e. we solve the system

$$\begin{pmatrix} 1 & 0 & \cdots & \cdots & \cdots & 0 \\ l_1 & 1 & 0 & & & \vdots \\ 0 & l_2 & 1 & \ddots & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \vdots \\ \vdots & & & l_{n-1} & 1 & 0 \\ 0 & \cdots & \cdots & 0 & l_n & 1 \end{pmatrix} \begin{pmatrix} z_0 \\ z_1 \\ \vdots \\ \vdots \\ z_{n-1} \\ z_n \end{pmatrix} = \begin{pmatrix} p_0 \\ p_1 \\ p_2 \\ \vdots \\ p_{n-1} \\ p_n \end{pmatrix} \quad (3)$$

Trivially, the z_i can be found by forward substitution so

$$z_0 = p_0, \quad z_i = p_i - l_i z_{i-1}, \quad i = 1, 2, \dots, n$$

This determines the vector \mathbf{z} .

Having obtained \mathbf{z} we find \mathbf{y} by solving $\mathbf{U}\mathbf{y} = \mathbf{z}$;

$$\begin{pmatrix} d_0 & u_0 & 0 & \cdots & \cdots & 0 \\ 0 & d_1 & u_1 & & & \vdots \\ 0 & 0 & d_2 & u_2 & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & 0 \\ \vdots & & & 0 & d_{n-1} & u_{n-1} \\ 0 & \cdots & \cdots & 0 & 0 & d_n \end{pmatrix} \begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ \vdots \\ y_{n-1} \\ y_n \end{pmatrix} = \begin{pmatrix} z_0 \\ z_1 \\ \vdots \\ \vdots \\ z_{n-1} \\ z_n \end{pmatrix} \quad (4)$$

The solution of (4) is trivially obtained by backward substitution

$$y_n = \frac{z_n}{d_n}, \quad y_i = \frac{z_i - u_i y_{i+1}}{d_i}, \quad i = n-1, n-2, \dots, 2, 1$$

This gives us the solution, \mathbf{y} , of our original problem (1), $\mathbf{M}\mathbf{y} = \mathbf{p}$.

This method can be also extended to decompose non-sparse matrices, i.e. $A = LU$, thus opening up a wider class of associated methods.

$$= \begin{pmatrix} l_{11} & 0 & \cdots & \cdots & \cdots & 0 \\ l_{21} & l_{22} & 0 & & & \vdots \\ \vdots & l_{32} & & \ddots & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \vdots \\ \vdots & & & l_{n-1} & l_{n-1,n-1} & 0 \\ l_{n1} & l_{n2} & \cdots & \cdots & \cdots & l_{nn} \end{pmatrix} \times \begin{pmatrix} u_{11} & u_{12} & \cdots & \cdots & \cdots & u_{1n} \\ 0 & u_{22} & u_{23} & \cdots & \cdots & u_{2n} \\ 0 & 0 & \ddots & \ddots & & \vdots \\ \vdots & & \ddots & \ddots & \ddots & \vdots \\ \vdots & & & 0 & u_{n-1,n-1} & u_{n-1,n} \\ 0 & \cdots & \cdots & 0 & 0 & u_{nn} \end{pmatrix}$$

If $l_{ii} = 1 \quad \forall \quad 1 \leq i \leq n$ as mentioned in the earlier discussion, this ensures a unique solution and is called *Doolittle's method*. *Crout's method* requires $u_{ii} = 1 \quad \forall \quad 1 \leq i \leq n$.

Example:

Consider the matrix $\begin{pmatrix} 2 & -1 & 1 \\ 3 & 3 & 9 \\ 3 & 3 & 5 \end{pmatrix}$ which can be factorised into

$$L = \begin{pmatrix} 1 & 0 & 0 \\ 1.5 & 1 & 0 \\ 1.5 & 1 & 1 \end{pmatrix}, \quad U = \begin{pmatrix} 2 & -1 & 1 \\ 0 & 4.5 & 7.5 \\ 0 & 0 & -4 \end{pmatrix}$$

Iterative Techniques

Seldom used for low dimensional problems - Gaussian elimination preferred.

For large systems with a high percentage of zero entries these are efficient in terms of

1. Computer Storage

2. Computational Time

Consider $(n \times n)$ linear system: $A\underline{x} = \underline{b}$

Initial approximation $\underline{x}^{(0)}$ to solution \underline{x} which generates a sequence of vectors $\{\underline{x}^{(k)}\}_{k=0}^{\infty}$ that converges to \underline{x} .

Most techniques involve a process which converts $A\underline{x} = \underline{b}$ to

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

where T is an $(n \times n)$ matrix and \underline{c} is a vector.

The initial vector $\underline{x}^{(0)}$ is selected and a sequence of approximations generated by computing

$$\underline{x}^{(k)} = T\underline{x}^{(k-1)} + \underline{c}$$

for each $k = 1, 2, 3, \dots$

$$E_1 : \quad 10x_1 - x_2 + 2x_3 \quad = 6$$

$$E_2 : \quad -x_1 + 11x_2 - x_3 + 3x_4 \quad = 25$$

$$E_3 : \quad 2x_1 - x_2 + 10x_3 - x_4 \quad = -11$$

$$E_4 : \quad 3x_2 - x_3 + 8x_4 \quad = 15$$

Now consider $A\underline{x} = \underline{b}$ with

$$A = \begin{bmatrix} 10 & -1 & 2 & 0 \\ -1 & 11 & -1 & 3 \\ 2 & -1 & 10 & -1 \\ 0 & 3 & -1 & 8 \end{bmatrix}, \quad \underline{b} = \begin{bmatrix} 6 \\ 25 \\ -11 \\ 15 \end{bmatrix}$$

which has an exact solution

$$\underline{x} = \begin{bmatrix} 1 \\ 2 \\ -1 \\ 1 \end{bmatrix}$$

To convert $A\underline{x} = \underline{b}$ to the form $\underline{x} = T\underline{x} + \underline{c}$; solve E_i for each x_i ($i = 1, 2, 3, 4$) to give

$$\left. \begin{aligned} x_1^{(k)} &= \frac{1}{10} \left\{ x_2^{(k-1)} - 2x_3^{(k-1)} + 6 \right\} \\ x_2^{(k)} &= \frac{1}{11} \left\{ x_1^{(k-1)} + x_3^{(k-1)} - 3x_4^{(k-1)} + 25 \right\} \\ x_3^{(k)} &= \frac{1}{10} \left\{ -2x_1^{(k-1)} + x_2^{(k-1)} + x_4^{(k-1)} - 11 \right\} \\ x_4^{(k)} &= \frac{1}{8} \left\{ -3x_2^{(k-1)} + x_3^{(k-1)} + 15 \right\} \end{aligned} \right\} \quad (5)$$

So (5) can be written as

$$\left. \begin{aligned} x_1^{(k)} &= \frac{1}{10} \left\{ x_2^{(k-1)} - 2x_3^{(k-1)} \right\} + \frac{3}{5} \\ x_2^{(k)} &= \frac{1}{11} \left\{ x_1^{(k-1)} + x_3^{(k-1)} - 3x_4^{(k-1)} \right\} + \frac{25}{11} \\ x_3^{(k)} &= \frac{1}{10} \left\{ -2x_1^{(k-1)} + x_2^{(k-1)} + x_4^{(k-1)} \right\} - \frac{11}{10} \\ x_4^{(k)} &= \frac{1}{8} \left\{ -3x_2^{(k-1)} + x_3^{(k-1)} \right\} + \frac{15}{8} \end{aligned} \right\}$$

\Rightarrow

$$T = \begin{bmatrix} 0 & \frac{1}{10} & -\frac{1}{5} & 0 \\ \frac{1}{11} & 0 & \frac{1}{11} & -\frac{3}{11} \\ -\frac{1}{5} & \frac{1}{10} & 0 & \frac{1}{10} \\ 0 & -\frac{3}{8} & \frac{1}{8} & 0 \end{bmatrix} \quad - = \begin{bmatrix} \frac{3}{5} \\ \frac{25}{11} \\ -\frac{11}{10} \\ \frac{15}{8} \end{bmatrix}$$

For initial approximation let $\underline{x}^{(0)} = \underline{0}^T$ and generate $\underline{x}^{(1)}$ by substi-

tuting $(0, 0, 0, 0)^T$ in (5) to get

$$(0.6, 2.2727, -1.1000, 1.8750) = \underline{x}^{(1)}$$

and $\underline{x}^{(2)}$ can now be obtained.

In general obtain $\underline{x}^{(k)} = \left(x_1^{(k)}, x_2^{(k)}, x_3^{(k)}, x_4^{(k)} \right)^T$.

By the very nature of iterative techniques, these continue indefinitely, hence one requires a convergence criterion to terminate the computation once the imposed condition is satisfied. So we need to know when to stop iterating, i.e. does the sequence converge to the solution of the given system of equations? To answer this, we need a means of measuring the distance between n — dimensional vectors.

Vector Norms

In two and three dimensions, the size of vectors called the modulus is generally obtained by Pythagoras. So if

$$\underline{v} = (a_1, a_2, a_3)$$

then the modulus of this vector $|\underline{v}| = \sqrt{(a_1^2 + a_2^2 + a_3^2)}$ which gives us the distance from the origin.

Now consider $\underline{x} = (x_1, x_2, \dots, x_n)^T$. Let \mathbb{R}^n denote the set of all n -dimensional vectors (so all vector components are real). To define distance in \mathbb{R}^n we use the notion of a *norm*. A pair of double vertical lines $\|\cdot\|$ is used to denote a norm, i.e. size of an n -dimensional vector.

More formally a *vector norm* on \mathbb{R}^n is a function, $\|\cdot\|$, such that

$$\|\cdot\| : \mathbb{R}^n \longrightarrow \mathbb{R}$$

with the following properties:

(i) $\|\underline{x}\| \geq 0 \quad \forall \underline{x} \in \mathbb{R}^n$

(ii) $\|\underline{x}\| = 0 \quad \text{iff } \underline{x} = \underline{0}$

(iii) $\|\alpha \underline{x}\| = |\alpha| \|\underline{x}\| \quad \forall \alpha \in \mathbb{R} \text{ and } \underline{x} \in \mathbb{R}^n$

(iv) $\|\underline{x} + \underline{y}\| \leq \|\underline{x}\| + \|\underline{y}\| \quad \forall \underline{x}, \underline{y} \in \mathbb{R}^n \quad (\text{Triangle inequality})$

If $\underline{x} = (x_1, x_2, \dots, x_n)^T$, then the vector norm $\|\underline{x}\|_p$ for $p = 1, 2, \dots$ is defined as

$$\|\underline{x}\|_p = \left\{ \sum_{i=1}^n |x_i|^p \right\}^{1/p}$$

Then there are a number of ways of defining a *norm*, depending on the value of p .

The special case $\|\underline{x}\|_\infty$ is given by

$$\|\underline{x}\|_\infty = \max_{1 \leq i \leq n} |x_i| \quad l_\infty \text{ norm.}$$

The most commonly encountered vector norm (also known as the modulus of a vector) is $\|\underline{x}\|_2$ and defined by

$$\|\underline{x}\|_2 = \left\{ \sum_{i=1}^n x_i^2 \right\}^{1/2} \quad l_2 \text{ norm}$$

For obvious reasons this is also called the *Euclidean norm*.

Examples

1. $\underline{x} = (\sin k, \cos k, 2^k)^\top$; k is a positive integer

$$\|\underline{x}\|_2 = \sqrt{\sin^2 k + \cos^2 k + 4^k} = \sqrt{1 + 4^k}$$

$$\|\underline{x}\|_\infty = \max|\sin k, \cos k, 2^k| = 2^k$$

2. $\underline{x} = (2, 1, 3, -4)^\top$

$$\|\underline{x}\|_2 = \sqrt{30} \quad \text{and} \quad \|\underline{x}\|_\infty = 4$$

These and other types of vector norms are further summarized in the following table, using a particular vector $\underline{v} = (1, 2, 3)$

Name	Symbol	value	numerical value
l_1 – norm	$\ \underline{x}\ _1$	6	6.000
l_2 – norm	$\ \underline{x}\ _2$	$\sqrt{14}$	3.742
l_3 – norm	$\ \underline{x}\ _3$	$6^{2/3}$	3.302
l_4 – norm	$\ \underline{x}\ _4$	$2^{1/4}\sqrt{7}$	3.146
l_∞ – norm	$\ \underline{x}\ _\infty$	3	3.00

If $\underline{x} = (x_1, x_2, \dots, x_n)^T$, $\underline{y} = (y_1, y_2, \dots, y_n)^T \in \mathbb{R}^n$, then the definitions for l_2 and l_∞ can be easily extended to consider distances between \underline{x} and \underline{y}

$$\|\underline{x} - \underline{y}\|_2 = \left\{ \sum_{i=1}^n (x_i - y_i)^2 \right\}^{1/2}$$

$$\|\underline{x} - \underline{y}\|_{\infty} = \max_{1 \leq i \leq n} |x_i - y_i|$$

Having introduced the concept of a vector norm, we can now define the most appropriate form for our convergence criteria, by making use of the l_{∞} norm. One such relative condition makes use of the infinity norm l_{∞} as a test for convergence, where

$$\frac{\|\underline{x}^{(k)} - \underline{x}^{(k-1)}\|_{\infty}}{\|\underline{x}^{(k-1)}\|_{\infty}} < \varepsilon.$$

ε is the specified tolerance (> 0) and $\|\underline{x}\|_{\infty} = \max_{1 \leq i \leq n} |x_i|$.

In the numerical example, actual computations yield for $k = 9$ & 10

$k = 9$	$k = 10$	
0.9997	1.0001	based on $\varepsilon = 10^{-3}$
2.0004	1.9998	
-1.0004	-0.9998	
1.0006	0.9998	

$$\begin{aligned}
 \frac{\|\underline{x}^{(10)} - \underline{x}^{(9)}\|_{\infty}}{\|\underline{x}^{(9)}\|_{\infty}} &= \frac{\|(4, -6, 6, -8) 10^{-4}\|_{\infty}}{1.9998} \\
 &= 4 \times 10^{-4} < \varepsilon
 \end{aligned}$$

Also $\|\underline{x}^{(10)} - \underline{x}^{(9)}\|_{\infty} = 2 \times 10^{-4} < \varepsilon$

The iterative technique used is called the *JACOBI* method.

It consists of solving the i^{th} equation in $A\underline{x} = \underline{b}$ for x_i (provided $a_{ii} \neq 0$) to obtain

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

and generating each $x_i^{(k)}$ from components of $\underline{x}^{(k-1)}$ ($k \geq 1$) by

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

We can now refine this method very simply by using the most up-to-date x_i .

To compute $x_i^{(k)}$, components of $\underline{x}^{(k-1)}$ are used.

Since for $i > 1$, $x_1^{(k)}, \dots, x_{i-1}^{(k)}$ have already been computed and are likely to be better approximations to the actual solutions x_1, \dots, x_{i-1} than $x_1^{(k-1)}, \dots, x_{i-1}^{(k-1)}$, it is far better to calculate $x_i^{(k)}$ using the most recently calculated values, i.e.

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

$$(i = 1, \dots, n)$$

This is called the *GAUSS-SEIDEL* method.

Re-doing the previous example using the G-S method iterative technique:

$$10x_1 - x_2 + 2x_3 = 6$$

$$-x_1 + 11x_2 - x_3 + 3x_4 = 25$$

$$2x_1 - x_2 + 10x_3 - x_4 = -11$$

$$3x_2 - x_3 + 8x_4 = 15$$

which is written as

$$\left. \begin{aligned} x_1^{(k)} &= \frac{1}{10} \left\{ x_2^{(k-1)} - 2x_3^{(k-1)} + 6 \right\} \\ x_2^{(k)} &= \frac{1}{11} \left\{ x_1^{(k)} + x_3^{(k-1)} - 3x_4^{(k-1)} + 25 \right\} \\ x_3^{(k)} &= \frac{1}{10} \left\{ -2x_1^{(k)} + x_2^{(k)} + x_4^{(k-1)} - 11 \right\} \\ x_4^{(k)} &= \frac{1}{8} \left\{ -3x_2^{(k)} + x_3^{(k)} + 15 \right\} \end{aligned} \right\}$$

Again letting $\underline{x}^{(0)} = \underline{0}^T$

Here

$$\underline{x}^{(4)} = (1.0009, 2.0003, -1.0003, 0.9999)$$

$$\underline{x}^{(5)} = (1.0001, 2.0000, -1.0000, 1.0000)$$

$$\frac{\|\underline{x}^{(5)} - \underline{x}^{(4)}\|_{\infty}}{\|\underline{x}^{(4)}\|_{\infty}} = 4 \times 10^{-4} \quad (< \varepsilon)$$

We note that Jacobi's method required twice as many iterations for the same level of accuracy.

If A is strictly diagonally dominant then for any choice of $\underline{x}^{(0)}$ the sequence of solutions generated by both Gauss-Seidel and Jacobi converge to the unique solution.

Root Finding - Nonlinear Algebraic Equations

A fundamental problem in numerical analysis consists of obtaining the zero of a function. Given a function $y = f(x)$ obtain the root of $f(x) = 0$, i.e. find the value of $x =$ which satisfies $f(x) = 0$.

e.g.

$$f(x) = x - \sin x$$

or another example:- zero's of the polynomial

$$f(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$$

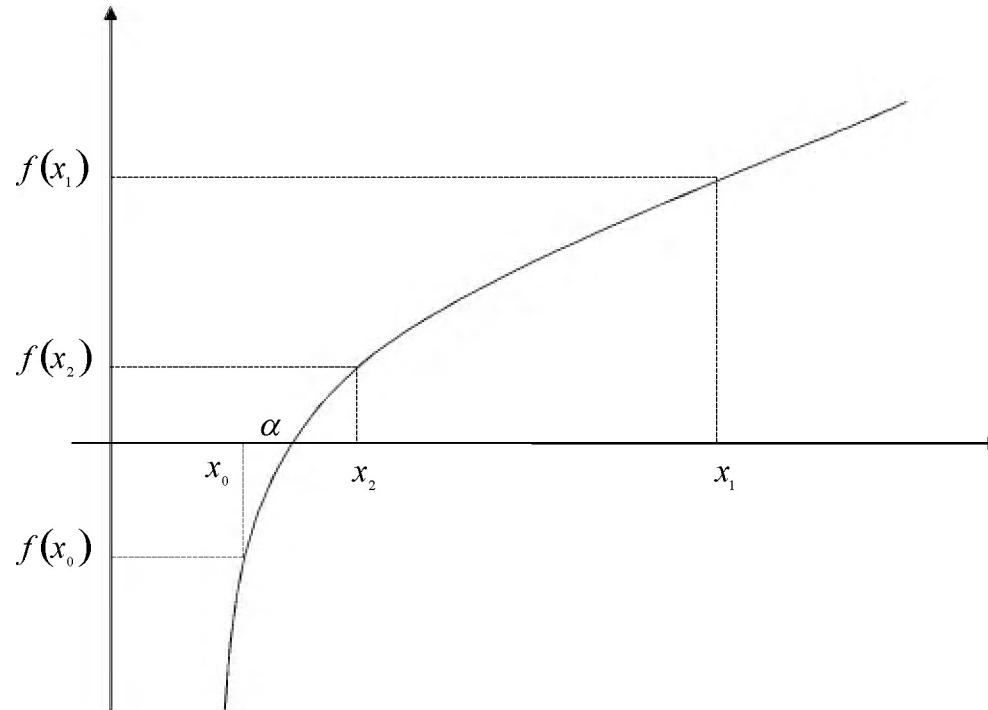
Generally the various schemes can be grouped under four sections: -

- (i) Methods which do not use derivatives of the function
- (ii) Methods which do use $f'(x)$
- (iii) Methods for polynomials
- (iv) Methods which deal with complex roots

Many computer programs (e.g. Mathematica, matlab and excel) have built-in equation solving algorithms, providing very accurate approximations to the solution. Two of the best known root finding algorithms are the bisection method and Newton's method. The latter uses the derivative and provides us with a highly efficient scheme. The one drawback is possible failure when the original estimate is too far from the root. The bisection method does not use any form of calculus and usually takes longer to converge - but is beautifully simple to use.

Methods Which Do Not Use Derivatives - Bisection Method

First we assume that an interval exists which contains the root which is required. We look at this graphically:



In this diagram

$$f(x_2) f(x_1) > 0$$

$$f(x_0) f(x_1) < 0$$

$$f(x_0) f(x_2) < 0$$

Hence there is a root in $[x_0, x_1]$ as indicated by the fact that the product

$$f(x_0) f(x_1) \leq 0$$

The method of bisection continues by finding the midpoint of $[x_0, x_1]$

$$x_2 = \frac{(x_0 + x_1)}{2}$$

Then calculate the function value at x_2 . If the product

$$f(x_0) f(x_2) \leq 0$$

then the root must lie in the interval $[x_0, x_2]$

(This is the case in the diagram given above)

If on the other hand

$$f(x_0) f(x_2) > 0$$

then the root lies in the interval $[x_2, x_1]$

Whichever interval is chosen, the new interval containing the root can be further subdivided. After n steps of bisection the interval containing the root will be reduced in size to

$$\frac{(x_1 - x_0)}{2^n}$$

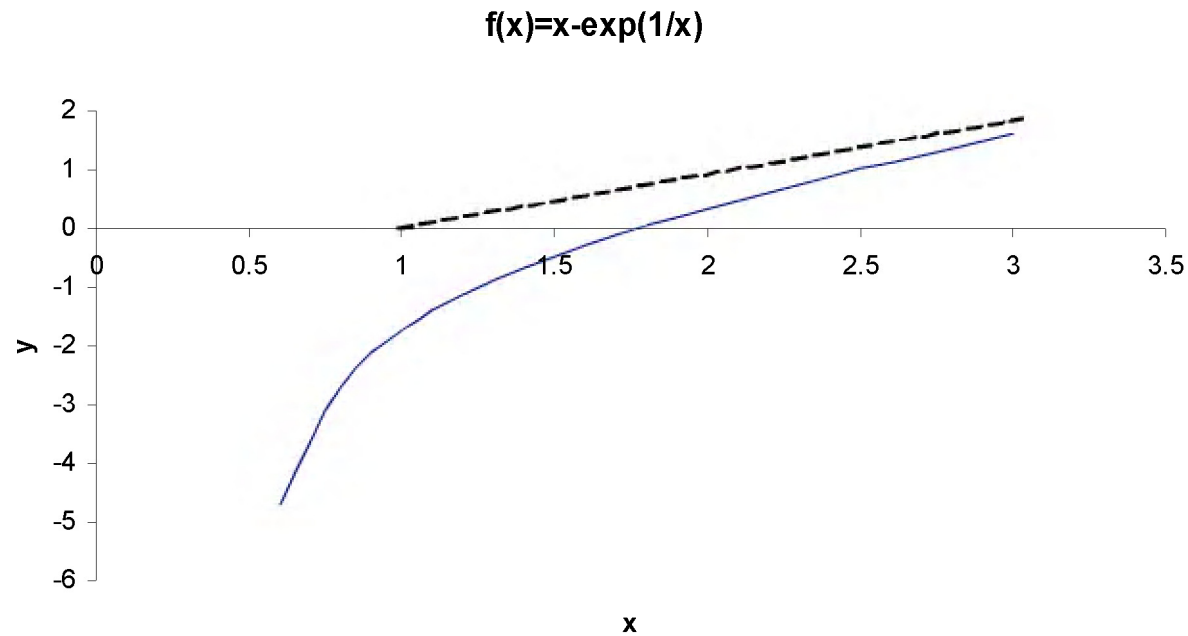
where in the example the value $x_1 = 2$ and $x_0 = 1$.

If the size of the interval becomes smaller than some specified tolerance, t , then the calculation stops.

Example

Consider

$$f(x) = x - e^{1/x}$$



There is a root in $[1, 2]$

Use the bisection method to show that the root of $f(x) = x - e^{1/x}$ in the interval $[1, 2]$ is 1.763 (correct to 3 decimal places)

$$f(x_0) = f(1) = 1 - e^1 < 0$$

$$f(x_1) = f(2) > 0$$

$$x_2 = \frac{x_0 + x_1}{2} = 1.5$$

$$f(1.5) = 1.5 - e^{2/3} = f(x_2) = -0.4477$$

$$f(x_0) f(x_2) > 0$$

\therefore root in $[x_2, x_1]$ i.e. in $[1.5, 2]$

$$x_3 = \frac{x_2 + x_1}{2} = \frac{1.5 + 2}{2} = 1.75$$

$$f(x_3) = 1.75 - e^{1/1.75} = -0.0208$$

$$f(x_3) f(x_2) > 0$$

\therefore root in $[x_3, x_1]$ i.e. in $[1.75, 2]$

$$\begin{aligned}x_4 &= \frac{1.75 + 2}{2} = \frac{3.75}{2} = 1.875 \\f(x_4) &= f(1.875) = 0.1704 \\f(x_3)f(x_4) &< 0\end{aligned}$$

\therefore root in $[x_3, x_4]$

i.e. in $[1.75, 1.875]$

$$\begin{aligned}x_5 &= \frac{x_3 + x_4}{2} = \frac{1.75 + 1.875}{2} = 1.8125 \\f(x_5) &= f(1.8125) = 0.0763 \\f(x_3)f(x_5) &< 0\end{aligned}$$

\therefore root in $[x_3, x_5]$

i.e. in $[1.75, 1.8125]$

$$\begin{aligned}x_6 &= \frac{x_3 + x_5}{2} = \frac{1.75 + 1.8125}{2} = 1.78125 \\f(x_6) &= f(1.78125) > 0 \\f(x_3) f(x_6) &< 0\end{aligned}$$

\therefore root in $[x_3, x_6]$

i.e. in $[1.75, 1.78125]$

$$\begin{aligned}x_7 &= \frac{x_3 + x_6}{2} = 1.765625 \\f(x_7) &= f(1.765625) > 0 \\f(x_3) f(x_7) &< 0\end{aligned}$$

\therefore root in $[x_3, x_7]$

i.e. in $[1.75, 1.765625]$

$$\begin{aligned}x_8 &= \frac{x_3 + x_7}{2} = 1.7578125 \\f(x_8) &= f(1.7578125) < 0 \\f(x_3) f(x_8) &> 0\end{aligned}$$

\therefore root in $[x_8, x_7]$

i.e. in $[1.7578125, 1.765625]$

$$\begin{aligned}x_9 &= \frac{x_8 + x_7}{2} = 1.76178 \\f(x_9) &= f(1.76178) < 0 \\f(x_8) f(x_9) &> 0\end{aligned}$$

\therefore root in $[x_8, x_7]$

i.e. in $[1.761718, 1.765625]$

$$\begin{aligned}x_{10} &= \frac{x_9 + x_7}{2} = 1.7636715 \\f(x_{10}) &= f(1.7636715) > 0 \\f(x_9) f(x_{10}) &< 0\end{aligned}$$

\therefore root in $[x_9, x_{10}]$

i.e. in $[1.761718, 1.7636715]$

$$\begin{aligned}x_{11} &= \frac{x_9 + x_{10}}{2} = 1.76269 \\&= 1.763 \text{ to 3 decimal places}\end{aligned}$$

Methods Which Do Use The Derivative $f'(x)$

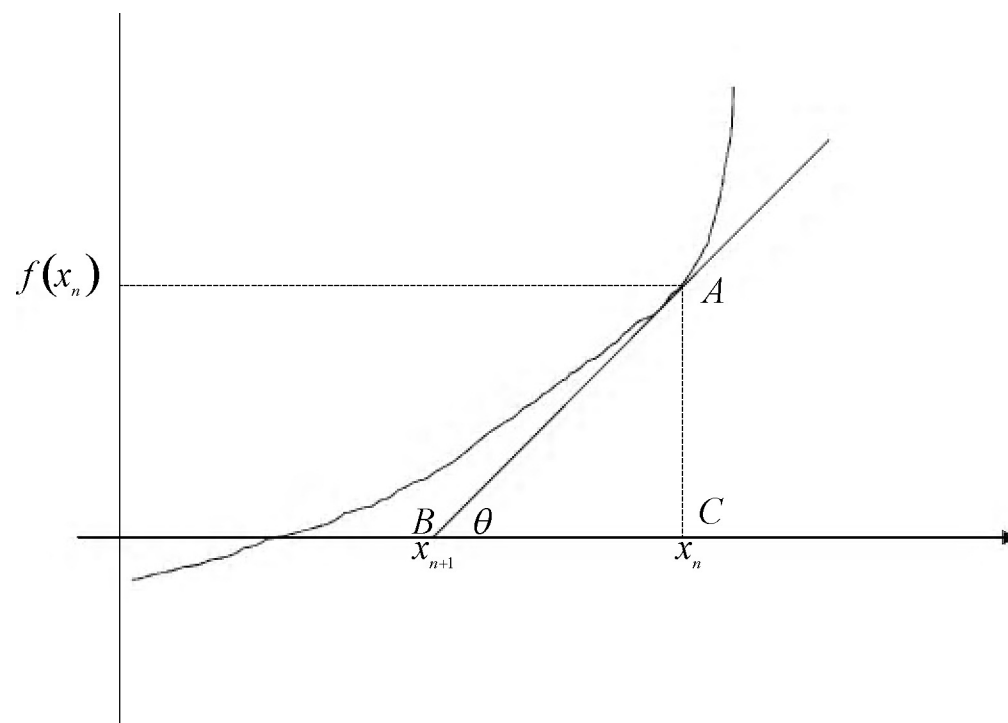
Newton-Raphson Method

For example

$$\begin{aligned}f(x) &= x - e^{1/x} \\f'(x) &= 1 - e^{1/x} \cdot \left(-\frac{1}{x^2}\right)\end{aligned}$$

$$f'(x) = 1 + \frac{e^{1/x}}{x^2}$$

From above we see that $\tan \theta = \frac{AC}{BC} = \frac{f(x_n)}{(x_n - x_{n+1})}$



But

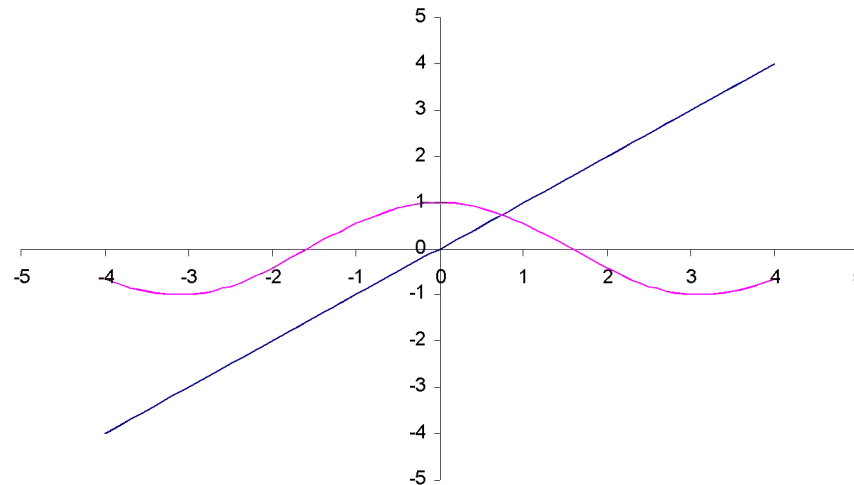
$$\begin{aligned}\tan \theta &= f'(x_n) \\ f'(x_n) &= \frac{f(x_n)}{(x_n - x_{n+1})}\end{aligned}$$

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

This is the *Newton-Raphson Technique*.

Example: Solve for roots, the function $f(x) = x - \cos x$.

Start by considering $x = \cos x$. That is draw $y = x$ and $y = \cos x$ to obtain an initial guess for the root(s).



Clearly the diagram above shows that there is only one root $\alpha \in (0, 1)$.

We use the Newton formula

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}, \quad n = 0, 1, \dots$$

where $n = 0$ is the initial guess. $f(x_n) = x_n - \cos x_n \longrightarrow f'(x_n) = 1 + \sin x_n$

$$\begin{array}{ccc} x & 0 & 1 \\ f(x) & -1 & 0.75 \end{array}$$

so numerically we also see that $f(0)f(1) < 0 \implies \alpha \in (0, 1)$. NR formula for this function becomes

$$x_{n+1} = x_n - \frac{x_n - \cos x_n}{1 + \sin x_n}, \quad x_0 = 1$$

$$x_1 = x_0 - \frac{x_0 - \cos x_0}{1 + \sin x_0} = 0.75036$$

$$x_2 = 0.75036 - \frac{0.75036 - \cos 0.75036}{1 + \sin 0.75036} = 0.73911$$

$$x_3 = 0.73911 - \frac{0.73911 - \cos 0.73911}{1 + \sin 0.73911} = 0.73909$$

$$x_4 = 0.73909 - \frac{0.73909 - \cos 0.73909}{1 + \sin 0.73909} = 0.73909$$

which gives the root $\alpha \approx 0.73909$