SINGLE precision macheps $= 2^{-23} \approx 1.2 \times 10^{-7}$

DOUBLE _____ $= 2^{-52} \approx 2.2 \times 10^{-16}$

The gap between a number and the next larger floating point number is:

$$ulp(x) = (0.\underbrace{00\cdots 01}_{P-2})_2 \times 2^E = 2^{-(P-1)} \times 2^E = \varepsilon 2^E$$

$\uparrow$ unit in last position

(we assume $x > 0$, if $x < 0$ then there is gap between $x$ and next smaller floating point number)

<u>Bottom line</u>: There is an error committed by just representing the floating point number in a computer!

One can show that:

$$\boxed{fl(x) = x(1+\delta)}$$

where $|\delta| \leq \varepsilon$ = machine epsilon

(floating point rep. of a number $x$)

( in the case of round to nearest $|\delta| < \frac{\varepsilon}{2}$ )

This means the <u>relative</u> error of rounding is

$$\left| \frac{x - fl(x)}{x} \right| = |\delta| \leq \varepsilon$$

A key feature of the IEEE standard is that it guarantees correctly rounded operations:

$$\boxed{fl(x \odot y) = (x \odot y)(1+\delta)}$$

where $|\delta| \leq \varepsilon$ = machine epsilon

where $\odot = +, -, \times, \div$ operations, $x, y$ are machine numbers

$\Rightarrow$ any single operation is guaranteed to give a result w/ rel-error $\leq$ macheps.

This is achieved in hardware by working with higher precision intermediate results and then rounding to desired precision. (sometimes called "guard bits")

However the result of a sequence of two or more arithmetic operations is not guaranteed to be correctly rounded.

Take for example:

$$x = 1 \qquad \text{and} \qquad (x+y)-z = 2^{-25}$$
$$y = 2^{-25}$$
$$z = 1$$

(all representable with no error in say IEEE single precision)

$$x+y = 1+2^{-25} = 1.\underbrace{0\cdots0}_{24\ zeros}1 \qquad \text{(in base 2)}$$

$$fl(x+y) = 1 \qquad (\text{because } \varepsilon = 2^{-23} \text{ in single precision})$$

thus:

$$fl\big(fl(x+y)-z\big) = 0 \qquad \text{whereas the exact result } 2^{-25} \text{ can be represented with no error!}$$

(note: • result would be correct if IEEE double precision is used
• we used round to nearest number. Other rounding method may give different results. )

Other aspects of IEEE standard you may encounter are:

signed 0: $\qquad +0 = -0 \qquad$ (sign bit does mean something in this case)

infinity: $\qquad \dfrac{1}{+0} = +\infty, \quad \dfrac{1}{-0} = -\infty$

NaN = result of an invalid operation (i.e. bug)
$\qquad$ = Not A Number

When do you encounter NaNs?

it makes sense to define:

$$a \times 0 = 0 \quad \text{if} \quad a \text{ finite}$$

$$\frac{a}{0} = \infty \quad \text{if} \quad a > 0$$

$$a \times \infty = \infty \quad \text{if} \quad a > 0$$

etc...

However the operations:

$$\left. \begin{array}{c} 0 \times \infty \\ \frac{0}{0} \\ \infty - \infty \end{array} \right\} \text{give NaN.}$$

<u>Question</u>: what about

$$\frac{\infty}{0} = \infty$$

$$\frac{0}{\infty} = 0$$

$$\frac{\infty}{\infty} = \text{NaN} \quad ( \text{I could have} \\ x_k \to \infty, y_k \to \infty \text{ but} \\ \frac{x_k}{y_k} \text{ undefined limit.} \\ \text{if I don't know more} )$$

+ parallel resistor example

---

If $x$ approximates a number $x_*$:

$$x - x_* = \underline{\text{error}}$$

$$|x - x_*| = \underline{\text{absolute error}}$$

$$\frac{|x - x_*|}{|x|} = \underline{\text{relative error}} \quad ( \text{if } x \text{ and } x_* \text{ are close it doesn't affect} \\ \text{not that much to change denominator} \\ \text{by } |x_*| )$$

# Chapter 2   Nonlinear equations

We will study algorithms to find the <u>roots</u> or <u>zeros</u> of some function $f: \mathbb{R} \to \mathbb{R}$, i.e. the points $x$ for which $f(x) = 0$. Numerical algorithms are essential as many times it is not possible to find an explicit formula for the root of a function. (without going too far, it is known that for polynomials of degree $\geq 5$ there is no all purpose formula for finding the roots, as we have for e.g. a quadratic).

Other examples of non-linear equations:

$$f(x) = x - \tan x = 0$$

$$f(\lambda) = \alpha e^{\lambda} + \beta \frac{(e^{\lambda t} - 1)}{\lambda} - \gamma = 0 \qquad (t \text{ is given})$$

(population model, see book's intro for chap §2)

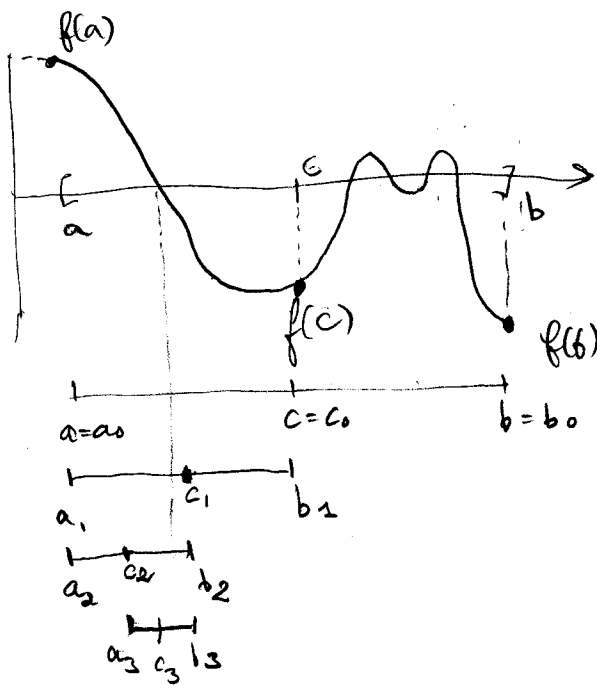## §2.1  The bisection or interval halving method

<u>key idea</u>:   By the intermediate value theorem, any <u>continuous</u> function on $[a, b]$ such that

$$f(a) f(b) < 0 \qquad (\text{i.e. function has opposite signs at end pts})$$

must have at least one zero in $[a, b]$.

This principle can be used to progressively narrow down an interval containing a root of $f(x)$.

Here $f(a)f(b) < 0$
and $f(a)f(c) < 0$ means
there is a root in $[a,c]$ interval.
then reiterate.

Note: it is clear that if there
are several roots in interval,
bisection will converge to one of
them. If we need a particular
root then initial interval must
contain only this root...

Here is how the algorithm looks like in pseudo code:

inputs: $a, b$, stopping criterion parameters, Maxit.

outputs: $a, b$

$u = f(a); \quad v = f(b); \quad e = b - a;$

if sign$(u)$ = sign$(v)$ then stop  (* bisection cannot be applied,
user needs to narrow down
location of root*)

for $k = 1 : Maxit$.
  $e = e/2$
  $c = a + e$      (* midpoint *)
  $w = f(c)$
  if $|e| < \delta$ or $|w| < \epsilon$ then stop    (* stopping criteria ... can be
                                                                      something else *)
  if sign$(w)$ sign$(u) < 0$
    $b = c; \quad v = w$     (* root is in left interval *)
  else
    $a = c; \quad u = w$     (* root is in right interval *)

## Stopping criterion:

Here we have 3 conditions that terminate iterations:

- # of iterations > Max iterations
  ( good to have this safeguard)

- $|b-a| < \delta$    error small enough i.e. we know our root within $|b-a|$

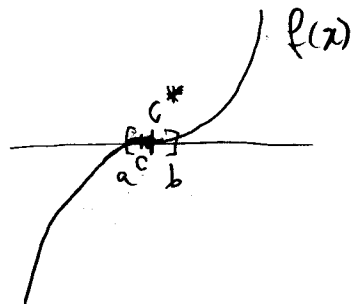- $|f(c)| < \delta$    c is sufficiently good a loop.

⚠️ Never test for a number == 0 in algorithms this is very unlikely!
and you will get very easily breakable code
( unfortunately this is the first stopping criterion suggestion in
algorithm 2.1 in book... please avoid doing it.
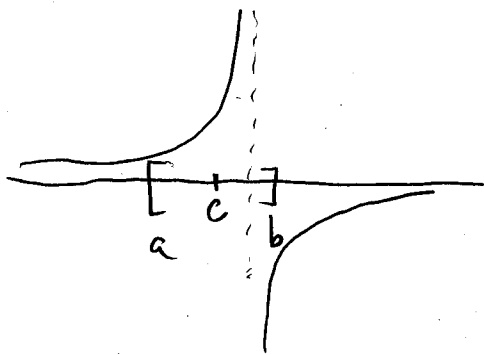
Why do we need all these criteria?

    that we should limit # of iter is clear
    also we can have the following pathological cases:



$c* =$ multiple root so $f$ is very flat
in a neighborhood of $c*$
~ hard to give root location with
a lot of precision

$|f(c)| < \varepsilon$ but $|b-a| > \delta$



Nobody stops user to supply a
discontinuous function in which
case it is better to end iterations
gracefully.

$|b-a| < \delta$ but $|f(c)| > \varepsilon$

Another stopping criterion that is suggested by your textbook is:

$$\frac{|c_n - c_{n-1}|}{|c_n|} < \varepsilon \quad , \quad c_n \neq 0$$

which means we stop when the relative difference in two consecutive approximations of the root become small.

The use of signs instead of values of $f(a), f(b), f(c)$ is to avoid potential overflow: when multiplying two floating point numbers, the result could be larger than largest representable number.

Error analysis: denote by $[a_0, b_0]$, $[a_1, b_1]$, ... the intervals that bracket the root, as generated by algorithm.

We of course have:

$$a_0 \leq a_1 \leq a_2 \leq \cdots \leq b_0$$
$$a_0 \leq \cdots \leq b_2 \leq b_1 \leq b_0$$
$$b_{n+1} - a_{n+1} = \frac{1}{2}(b_n - a_n) \quad \text{for } n \geq 0.$$

Since $\{a_n\}$ is an increasing sequence bdd above it converges
$\{b_n\}$ ———— decreasing ———— below ————

$$b_{n+1} - a_{n+1} = \frac{1}{2}(b_n - a_n) = \frac{1}{2^2}(b_{n-1} - a_{n-1})$$
$$\cdots$$
$$= \frac{1}{2^{n+1}}(b_0 - a_0)$$

$$\Rightarrow \lim_{n \to \infty} b_n - a_n = \lim \frac{1}{2^n}(b_0 - a_0) = 0$$

$\Rightarrow$ the limit of $\{a_n\}$ and $\{b_n\}$ is the same, call it $r$.

If we take the limit:

$$0 \geq f(a_n) f(b_n) \longrightarrow 0 \geq (f(r))^2$$

(and we can since $f$ is continuous)

$$\Rightarrow f(r) = 0$$
$r$ is root.

If we are at step $n$ in algorithm which point in $[a_n, b_n]$ should we take as our approx to the root? The best point (if we do not know anything else about root) is the mid point:

$$c_n = \frac{a_n + b_n}{2}$$

$$|r - c_n| \leq \frac{1}{2}(b_n - a_n) = \frac{1}{2}\frac{1}{2^n}(b_0 - a_0) = 2^{-(n+1)}(b_0 - a_0)$$

$\overset{\shortparallel}{\text{error}}$

Here is a theorem summarizing our findings.

## Theorem ( Bisection method)

If $[a_0, b_0], [a_1, b_1], [a_2, b_2]$ are the intervals generated by the bisection method then the limits $\lim_{n \to \infty} a_n$, $\lim_{n \to \infty} b_n$ exist and are equal, the limit being a root of $f$.

If $r = \lim_{n \to \infty} c_n = \frac{a_n + b_n}{2}$ Then:

$$|r - c_n| \leq \frac{1}{2^{n+1}}(b_0 - a_0)$$

With this explicit error bound it is easy to estimate how many iterates are needed to achieve a desired precision.

# § 2.3 Newton's method :

extremely important algorithm that can be generalized to systems of non-linear equations and can be used for example in optimization.

**key idea :** Let $x^*$ be a root of $f(x)$ and assume $f \in C^2$. Then Taylor's theorem says:

$$0 = f(x^*) = f(x+h) = f(x) + h f'(x) + O(h^2)$$

where $h = x^* - x$.

If we solve for $h$ : $\quad h = - \dfrac{f(x)}{f'(x)}$.

Thus if $x$ is close to $x^*$ we can expect $x + h = x - \dfrac{f(x)}{f'(x)}$

to be even closer (we will prove this rigorously).

Putting everything together we get the iteration:

$$x_0 = \text{given}$$
$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} \qquad (\text{breaks down if } f'(x^*) = 0)$$

## Newton's method

$$v = f(x_0)$$
if $|v| < \varepsilon$ then stop $\qquad$ (* $\varepsilon \equiv$ tolerance *)

for $k = 1 : Maxit$
$\qquad x_{new} = x - v / f'(x)$
$\qquad v = f(x_{new})$
$\qquad$ if $|x_{new} - x| < \delta$ or $|v| < \varepsilon$ stop
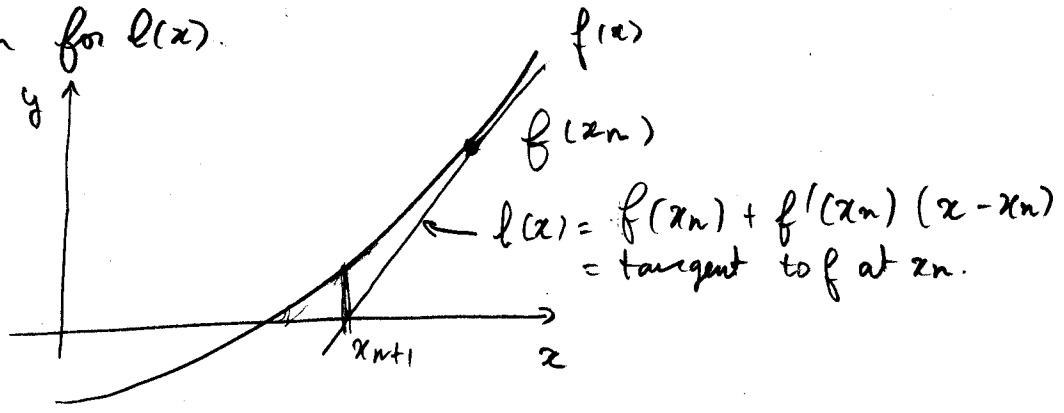$\qquad x = x_{new}$

<u>Graphical interpretation</u> : What we are doing is approximating the problem of finding a root of a non-linear function by a sequence of linear root-finding problems. Indeed if we replace $f(x)$ by its first two terms in Taylor expansion around a point $c$:
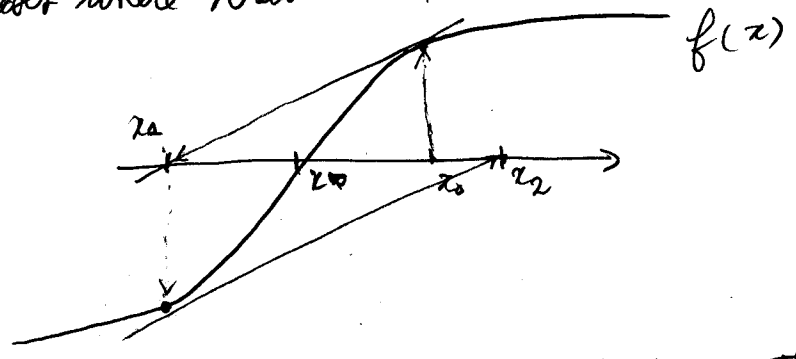
$$f(x) = f(c) + f'(c)(x-c) + \frac{f''(c)}{2}(x-c)^2 + \cdots$$

$$\ell(x) = f(c) + f'(c)(x-c) \quad = \quad \text{tangent of } f(x) \text{ at } c.$$

What Newton's method does is to solve the (trivial) root finding problem for $\ell(x)$.



$f(x)$

$f(x_n)$

$\ell(x) = f(x_n) + f'(x_n)(x - x_n)$
= tangent to $f$ at $x_n$.

$x_{n+1}$

Just by looking at this graphical construction one can come up with cases where Newton's method does not converge!



$f(x)$

$x_1 \quad x_0 \quad x_0 \quad x_2$

So any convergence result for Newton's method works only if initial guess is sufficiently close to $x_{**}$.

We shall prove in 2 ways the convergence of N.M.

The first way is straight forward and can be generalized easily to more dimensions.

For the second way we have to wait until we see the powerful contracting mapping theorem.

Let us first introduce the following notion of continuity.

<u>Def</u> (Lipschitz continuity) A function $f$ is said to be Lipschitz continuous with Lipschitz constant $L > 0$ if :

$$|f(x) - f(y)| \leq L |x - y|$$

Note that Lipschitz continuity $\implies$ continuity.

for some $\varepsilon > 0$ $\quad \exists \, \delta = \dfrac{\varepsilon}{L}$ s.t. $\quad |x - y| < \delta = \dfrac{\varepsilon}{L}$

$$\implies \quad |f(x) - f(y)| \leq L \frac{\varepsilon}{L} = \varepsilon .$$

Also if $f \in C^1[a,b]$ ($[a,b]$ bdd interval) then $f$ is Lipschitz cont.

Using Taylor's theorem :

$$f(x) = f(y) + f'(\xi)(x-y) \quad \text{for some } \xi \in [a,b]$$

$$\implies \quad |f(x) - f(y)| \leq |f'(\xi)| \, |x - y| \leq L \, |x - y|$$

where $L = \max\limits_{x \in [a,b]} |f'(x)|$, which exists and is attained because $f'(x) = $ cont. on closed & bdd = compact set.

We are now ready to show the following theorem. (slightly more powerful than books' theorem)

## Theorem (convergence of Newton's method)

Let $x_*$ be a simple zero of $f(x)$ i.e. s.t. $f'(x_*) \neq 0$, and assume $f'(x)$ is Lipschitz continuous on a nbd of $x_*$ with Lipschitz constant $L$. Then:

$$\exists \, \varepsilon > 0, \; c > 0 \; \text{s.t.}$$
$$|x_0 - x_*| < \varepsilon \implies |x_{n+1} - x_*| \leq c \, |x_n - x_*|^2$$

translated in english: Newton's method converges **quadratically** provided initial guess $x_0$ (iterate) is close enough to $x_*$.

**Proof:** We work by induction

$$x_{k+1} - x_* = x_k - \frac{f(x_k)}{f'(x_k)} - x_*$$

$$= \frac{1}{f'(x_k)} \left( \underbrace{f(x_*)}_{=0} - f(x_k) + f'(x_k)(x_k - x_*) \right)$$

Using fundamental theorem of calculus:

$$f(x_k) - f(x_*) = \int_{x_*}^{x_k} dx \, f'(x) \qquad \text{thus:}$$

$$A = f(x_*) - f(x_k) + f'(x_k)(x_k - x_*) = \int_{x_*}^{x_k} (f'(x_k) - f'(x)) \, dx$$

$$\overset{c.o.v.}{=} \int_1^0 \left( f'(x_k) - f'(x_k + t(x_* - x_k)) \right)(x_* - x_k) \, dt$$

$$x = x_k + t(x_* - x_k)$$

$$|A| \leq \int_0^1 |f'(x_k) - f'(x_k + t(x_* - x_k))| \, |x_* - x_k| \, dt$$

$$\leq \int_0^1 L t \, |x_* - x_k|^2 \, dt \qquad \text{by Lipschitz cont. of } f.$$

$$= \frac{L}{2} \, |x_k - x_*|^2$$

Therefore:

$$|x_{k+1} - x*| \leq \frac{1}{|f'(x_k)|} \frac{L}{2} |x_k - x*|^2$$

which screams for quadratic convergence if we can show
$|f'(x_k)|^{-1}$ is bounded in some neighborhood of $x*$.
Incidentally this justifies us dividing by $f'(x_k)$.

We will show that for $\varepsilon$ sufficiently small:

$$(*) \quad 2|f'(x_k)| \geq |f'(x*)| > 0$$

$\uparrow$ since $f'(x*) \neq 0$

From $(*)$ we get:

$$|x_{k+1} - x*| \leq |f'(x*)|^{-1} L |x_k - x*|^2$$

So in addition if $f'(x*)$ is small convergence is not as good as if $f'(x*)$ was large (makes sense from graphical prospective).

<u>Proof of $(*)$</u>: Since $f'(x)$ is continuous there is $\varepsilon_1$ for which:

$$|x - x*| < \varepsilon_1 \implies sgn\, f'(x) = sgn\, f'(x*)$$

if $f'(x*) > 0$:

$$f'(x) = f'(x*) + f'(x) - f'(x*)$$
$$\geq f'(x*) - |\qquad|$$
$$\geq f'(x*) - L|x - x*|$$

Choosing $\varepsilon = \min\left(\varepsilon_1, \frac{|f'(x*)|}{2L}\right)$:

$$f'(x) \geq f'(x*) - \frac{1}{2}|f'(x*)|$$
$$= \frac{f'(x*)}{2} \quad \text{since } f'(x*) > 0$$
$$\implies |f'(x)| \geq \frac{1}{2}|f'(x*)|$$

if $f'(x*) < 0$

$$f'(x) = f'(x*) + f'(x) - f'(x*)$$
$$\leq f'(x*) + |f'(x) - f'(x*)|$$
$$\leq f'(x*) + L|x - x*|$$

Choosing $\varepsilon = \min\left(\varepsilon_1, \frac{|f'(x*)|}{2}\right)$

$$\underbrace{f'(x)}_{-|f'(x)|} \leq \underbrace{f'(x*)}_{-|f'(x*)|} + \frac{|f'(x*)|}{2}$$

$$-|f'(x)| \leq -\frac{|f'(x*)|}{2}$$
$$|f'(x)| \geq \frac{1}{2}|f'(x*)|$$

Local convergence of Newton's method become global in some particular cases. For example:

## Theorem (Newton's method for convex functions)

Let $f$ be $C^2(\mathbb{R})$, convex and increasing. If $f$ has a zero it is unique and Newton's method converges regardless of initial iterate $x_0$.

## Secant method:

key idea: For Newton's method we need to know derivative of $f$, and this may not be available or be very costly. (this is especially true for high dimensions)

$\leadsto$ replace derivative by secant:

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

$$\approx \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}.$$

Therefore the secant update is:

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

and algorithm is similar to Newton's method, with the only difference being that we need to keep two iterates at any time.

The convergence for the secant method is <u>local</u> like for Newton's method.

The convergence rate is of $\frac{1 + \sqrt{5}}{2} \approx 1.6$

(do you recognize this rate from computer lab #2?)