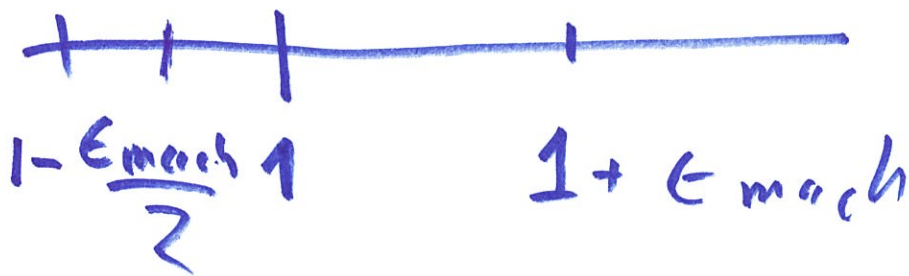


A

$$f(x) = \frac{1 - (1-x)}{x}$$

$$x = \epsilon_{mach} / 4$$



$$1-x \approx 1; f(x) = \frac{1-1}{1} = 0$$

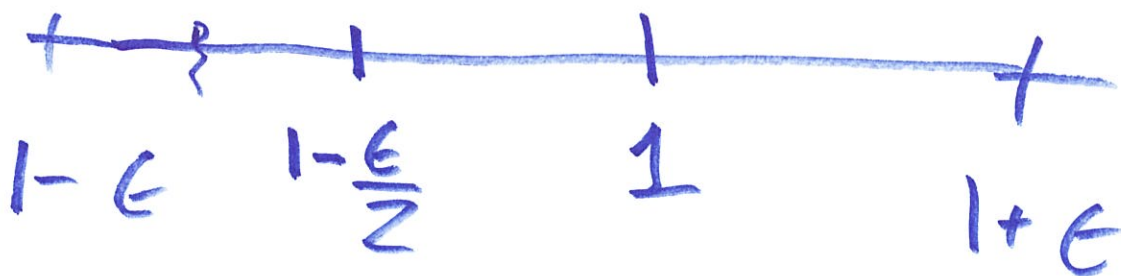
$$x = \frac{\epsilon_{mach}}{4} = 1.0001;$$

$$1-x \approx 1 - \frac{\epsilon_{mach}}{4}; 1 - (1-x) = \frac{\epsilon_{mach}}{4}$$

$$f(x) = 2$$

$$x = \epsilon_{\text{mach}} \cdot \frac{3}{4}$$

B



$$1-x = 1-\epsilon/2$$

$$1-(1-\epsilon) = \epsilon/2$$

$$f(x) = \frac{\epsilon}{2 \cdot \frac{3}{4}\epsilon} = \frac{4}{6} = \frac{2}{3} = 0.66$$

$$x = \epsilon_{\text{mach}} \cdot \frac{3}{4} = 1.001$$

$$1-x = 1-\epsilon; f(x) = \frac{\epsilon}{\frac{3}{4}\epsilon} = \frac{4}{3} = 1.33$$

C

Chopping  $\left| \frac{x - \text{float}(x)}{x} \right| < \epsilon_{\text{mach}}$   
round to nearest

$$\left| \frac{x - fl(x)}{x} \right| < \frac{\epsilon_{\text{mach}}}{2}$$

$\pm$  shift mantissa to match leading exponent

- Cancellation

$N=3$        $x = 1.32$   
 $y = 1.31$  ← trailing

$$x - y = 0.01$$

①  $x, y$

D

Since Toy floating pointing

$$L = -1 \quad U = 1$$

$$N = 3 \quad \beta = 2$$

write #, see the distances

Double

$$N = 5$$

②  ~~$S_1 = \sum \frac{1}{e}$~~

$$S_1 = \sum_{k=1}^{1000} \left( \frac{e^k}{1+e^k} \right) \frac{e^{-k}}{e^{-n}}$$

$$= \sum_{k=1}^{1000} \frac{1}{e^{-k} + 1} \approx 992$$

$$\begin{aligned}\sqrt{a} - \sqrt{b} &= \frac{(\sqrt{a} - \sqrt{b})(\sqrt{a} + \sqrt{b})}{\sqrt{a} + \sqrt{b}} \\ &= \frac{a - b}{\sqrt{a} + \sqrt{b}}\end{aligned}$$

$$\sin\left(\pi \cdot \left(k^{10} + \frac{1}{k}\right)\right)$$

$$\sin(x + 2\pi) = \sin x$$

$$\sin(x + 2\pi \cdot n) = \sin x$$

$n$  is integer

$$\sin(x + \pi) = -\sin(x)$$

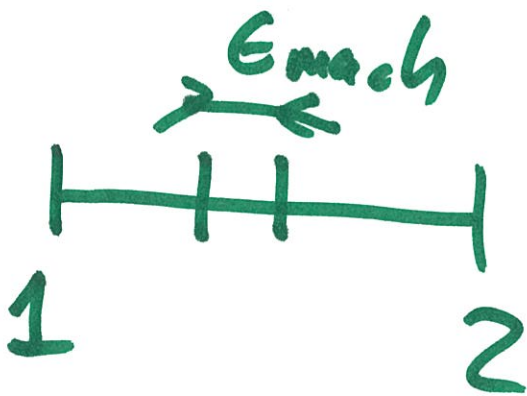
F

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

$$S_2 = \sum_{k=1}^{\infty} \frac{1}{k^2}$$

Analytical answer:  $\frac{\pi^2}{6}$

$$1 + \frac{1}{4} + \frac{1}{9} + \frac{1}{16} + \dots$$



Small

$$S = \sum_{k=1000}^{\infty} \frac{1}{k^2}$$

Sum, from large to small <sup>G</sup>

$$K = K_{\max}$$

$$K = K_{\max} + 1 \dots$$

What is the smallest  
positive # representable  
in Matlab (Double

M

LIFL

precision)

$$\text{realmin} = \beta^L = 2^{-1022}$$

subnormals

$$x_{\text{smallest}} = \underbrace{0.00\dots01}_{51} \cdot \beta^L =$$
$$= 2^{-1022} \cdot 2^{-52}$$



I

- 1
- 10
- 11
- 100
- 101
- 110
- 111
- 1000

$$52 = 32 + 16 + 4$$

$$= 110100$$

$$x_1 = \underbrace{1111 \dots 1}_5 =$$

$$53$$

$$= 1.111 \cdot 10^{110100}$$

$$\underbrace{\hspace{1.5cm}}_{52}$$

$$x_2 = x_1 + 1 = 1.000 \dots 10^{110101}$$

$$x_3 = \underbrace{11111}_{56 \text{ 'one's'}}$$

$$x = \pm \left( d_1 + \frac{d_2}{\beta} + \frac{d_3}{\beta^2} + \dots + \frac{d_N}{\beta^{N-1}} \right) \beta^E$$

$$= \pm \left( d_1 \cdot \beta^{N-1} + d_2 \cdot \beta^{N-2} + d_3 \cdot \beta^{N-3} \right.$$

$$\left. + \dots + d_{N-1} \cdot \beta + d_N \right) \cdot \frac{\beta^E}{\beta^{N-1}}$$

Find  $x^*$  s.t.

10

$$f(x^*) = \emptyset$$

$$x^6 + 5x^5 + 2x^4 + 3x^3 +$$

$$x = \cos x + 2x^2 + x + 5 = \emptyset$$

g

$\cos g$

$\cos(\cos(g))$

$x = \cos x$

Guess  $x_0$  as a possible  $\hookrightarrow$   
solution to  $f(x) = 0$   
iterations

$$x_{k+1} = g(x_k)$$

$\hookrightarrow$  iterative  
function

$$x_1 = g(x_0)$$

$$x_2 = g(x_1) = g(g(x_0))$$

$$x_3 = g(x_2) = g(g(g(x_0))) \dots$$

Hope  
 $x_n$  - "final" results

How fast the correct solution<sup>y</sup>  
is obtained?

convergence rate

$$f(x) = 0;$$

$x^*$  is an exact  
solution

$$f(x^*) = 0$$

$$x_0 \quad ; \quad E_0 = |x_0 - x^*|$$

$$x_1 = g(x_0) \quad ; \quad E_1 = |x_1 - x^*|$$

$$x_2 = g(x_1) \quad ; \quad E_2 = |x_2 - x^*|$$

$$E_{k+1} = C \cdot E_k^r$$

$C$  is a constant;  $r =$  convergence rate

N

# Bisection

$r = 1 \Rightarrow$  linear convergence

$$E_{k+1} = C \cdot E_k$$

~~$$|E_{k+1}| < |E_k|$$~~

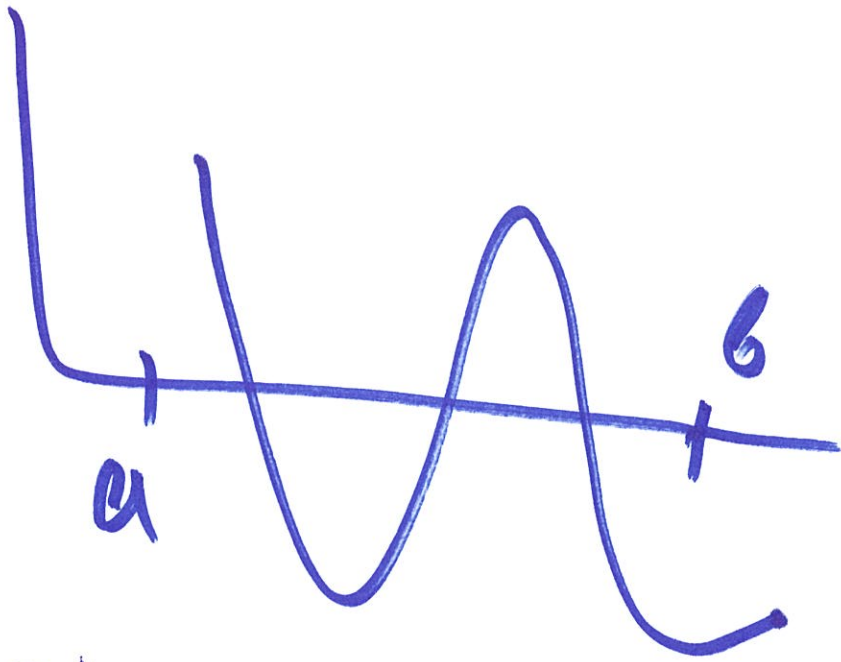
$$0 < C < 1$$

# Newton method

$r = 2$  quadratic

$1 < r < 2$  - super linear

Cubic convergence (HLW)



Positive: it will ALWAYS find a zero

$$E_{k+1} \leq \frac{E_k}{2}; r=1$$

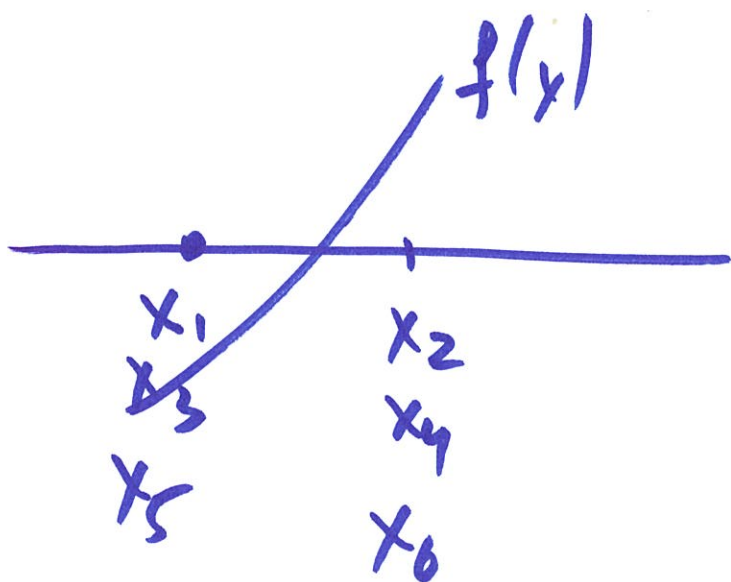
Linear convergence

One bit of a solution  
at a time

# Stopping criteria

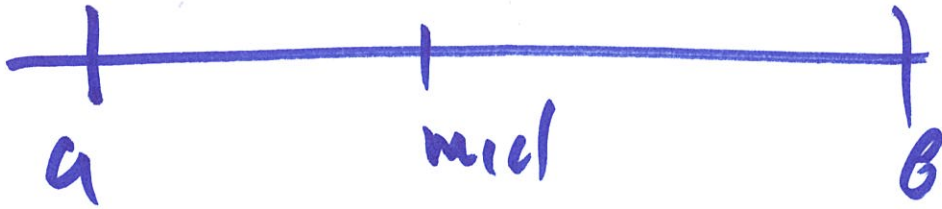
Q

- ① Stop if  $|f(x_c)| < \text{Tolerance}$
- ②  $|x_{k+1} - x_k| < \text{tolerance}$  ✓
- ③ Max # iterations





R



while  $|a - b| < \text{Tolerance}$

$$m = \frac{a + b}{2}$$

if  $f(a) \cdot f(m) < 0$

then  $b = m$

else  $a = m$

endif

end while

$B[a, b] :=$

{

$$m = \frac{a+b}{2}$$

if  $|f(m)| < \text{Tolerance}$   
Return  $m$

~~if  $f(m) < 0$~~

if  $f(a) \cdot f(m) < 0$

then  $B[a, m]$

else  $B[b, m]$

end if

}

—

—