

# Recap of last lecture

①. Calculate  $\sqrt{a}$  :  $x_{t+1} = \frac{1}{2} \left( x_t + \frac{a}{x_t} \right)$

②. Computer Arithmetic :  $x \in \mathbb{R} \rightarrow f(x)$

- Float-point numbers, roundoff errors

$$\left\{ \begin{array}{l} \epsilon_m = 2^{-52} \\ \left| \frac{f(x) - x}{x} \right| \leq \frac{1}{2} \cdot \epsilon_m \end{array} \right.$$

double:  $x = (-1)^s \cdot 1.b_1 b_2 \dots b_{52} \times 2^m$ ,  $m = (e_{11})_2 - 1023$

s	e <sub>11:11</sub>	b <sub>1:52</sub>
---	--------------------	-------------------

$$e_{11:11} = \begin{cases} 0 & \left\{ \begin{array}{l} \pm 0. b_{1:52} \times 2^{-1022} \rightarrow \begin{cases} +0 \\ -0 \end{cases} \\ 1-2046 \rightarrow m \in [-1022, 1023] \\ 2047 : \infty \text{ if } b_{1:52} = 0, \text{ NaN if } b_{1:52} \neq 0 \end{array} \right. \end{cases}$$

- loss of significant

③. Solve  $f(x) = 0$

- Bisection  $(f, a, b, \epsilon)$  :  $f(a) \cdot f(b) < 0$ .

$$\left\{ \begin{array}{l} 1. n \text{ iterations : } n+2 \text{ oracles} \\ 2. |r - c_n| \leq \frac{b-a}{2^{n+1}} \end{array} \right.$$

Today's lecture :

①.  $0 \in \{+, -, \times, \div\}$ ,  $x \circ y$  v.s.  $f(x \circ y)$

$\sum_{i=1}^n x_i$  : naive sum, pairwise sum

②. FPI, Newton's, Secant, Sensitivity

# ①. IEEE 754, $x \rightarrow fl(x)$ , $y \rightarrow fl(y)$

$x+y$  will have the following steps: Streaming SIMD Ext. 2

1. load  $x$  and  $y$  into registers  $\leftarrow$  SSE2/AVX <sup>x86</sup>

`MOVSD xmm0, [x]`

`MOVSD xmm1, [y]`

Advanced Vector

2. align the exponents:  $e_{1:n} \rightarrow \min(e_x, e_y)$

3. +/- mantissas  $\leftarrow (b_x + b_y)$

- extra bits of precision x87 FPU uses 80 bits
- Guard bits / extra bits in SSE2/AVX.

$\left\{ \begin{array}{l} 1: \pm 1 \\ 15: \text{exp} \\ 64: \text{mantissa} \end{array} \right.$

4. normalized  $\rightarrow$  rounding  $\rightarrow$  save

$x$              $y$   
 $\downarrow$              $\downarrow$   
 $fl(x)$      $fl(y)$  : rounding (may use extra bits)

$\downarrow$      $\downarrow$   
Reg.  $\rightarrow$  SSE2/AVX:  $fl(x) + fl(y) \rightarrow s$

$\downarrow$   
 $s \rightarrow fl(s)$  : rounding,  $fl(x \cdot y) := fl(s)$

$$fl(x \cdot y) = x \cdot y (1 + \delta), \quad |\delta| \leq \begin{cases} 2^{-24} & \text{single} \\ 2^{-53} & \text{double} \end{cases}$$

• We assume  $fl(x)$ ,  $fl(y)$ ,  $fl(x \cdot y)$  are representable!

Rel. error does not make sense when it is overflow or underflow.

Summation:  $\sum_{i=1}^n x_i$

• Naive sum:

Let  $\{x_i\}_{i=1}^n$  be positive machine numbers in a computer whose unit roundoff error is  $\epsilon$ .

$$S = \sum_{i=1}^n x_i \quad \left| \frac{S - f(S)}{S} \right| \leq (1 + \epsilon)^n - 1 \approx n \cdot \epsilon.$$

$$\hat{S} = 0$$

for  $i=1, 2, \dots, n$ :

$$\hat{S} = x_i$$

return  $\hat{S}$

$$S_k = \sum_{i=1}^k x_i \Leftrightarrow \hat{S}_k$$

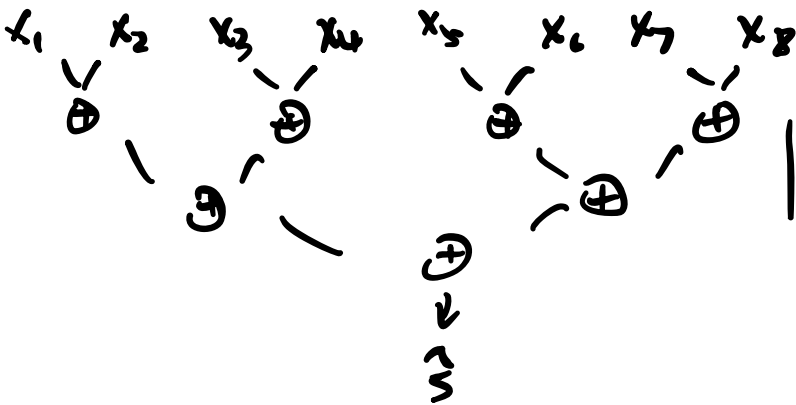
$$\left\{ \begin{array}{l} S_1 = x_1 \\ S_{k+1} = S_k + x_{k+1} \end{array} \right\} \left\{ \begin{array}{l} \hat{S}_1 = x_1 \\ \hat{S}_{k+1} = f(S_k + x_{k+1}) \end{array} \right.$$

def.  $P_k = \frac{\hat{S}_k - S_k}{S_k} \quad |P_k| \leq (1 + \epsilon)^n - 1.$

$$= 1 + \binom{n}{1} \cdot \epsilon + \binom{n}{2} \cdot \epsilon^2 + \dots$$

$$= O(n \cdot \epsilon).$$

• Pairwise sum:



$$\left| \frac{S - \hat{S}}{S} \right| \leq \frac{\epsilon \cdot \log_2 n}{1 - \epsilon \log_2 n} \cdot \left( \frac{\sum_{i=1}^n |x_i|}{\sum_{i=1}^n |x_i|} \right)$$

$$= O(\epsilon \cdot \log_2 n).$$

If we assume rounding errors is random

• naive sum:  $O(\epsilon \sqrt{n})$

• pairwise sum:  $O(\epsilon \sqrt{\log n})$ .

Further reading: Handbook of floating-point Arithmetic, Muller, 2018.

# Lecture-02. Solving nonlinear equations

09/11/2024

In this lecture, we want to solve

$$f(x) = 0, \quad x \in \mathbb{R}, \quad f \in C(\mathbb{R}) \text{ is continuous.}$$

We will assume  $f(a) \cdot f(b) < 0$  and root  $r \in [a, b]$ .

## ② Fixed Point Iteration (FPI)

(1). FPI:  $x_{t+1} = g(x_t) \rightarrow$  find fix point  $x = g(x)$

$$g(x) = \begin{cases} \frac{f(x) + x}{2} \\ x - \frac{f(x)}{f'(x)} \\ \vdots \end{cases} \text{ many ways}$$

Main idea: Guess  $x_0$ , compute  $g(x_0)$ . If  $x_0$  is near to  $r$ , then  $g(x_0)$  is near to  $g(r)$ .

$$x_{t+1} = g(x_t) \text{ until } \begin{cases} |x_{t+1} - x_t| \leq \text{tol} \\ |g(x_t) - x_t| \leq \text{tol} \end{cases}$$

Convergence analysis:  $e_t = x_t - r$ . ( $r$  is a fix point)

$$e_{t+1} = x_{t+1} - r = g(x_t) - g(r) \quad \text{v.s.} \quad x_t - r = e_t.$$

By mean value thm.  $g \in C^1[a, b]$ ,  $\exists \delta \in (a, b)$  s.t.

$$g'(\delta) = \frac{g(b) - g(a)}{b - a} \text{ . let } \begin{cases} b = x_t \\ a = r \\ \delta = \delta_t \end{cases}$$

$$|g(x_{t+1}) - g(r)| = |g'(s_t)| \cdot |x_t - r|, \quad s_t \in \text{conv}(r, x_t)$$

$$|e_{t+1}| = |g'(s_t)| \cdot |e_t|, \quad t \geq 0.$$

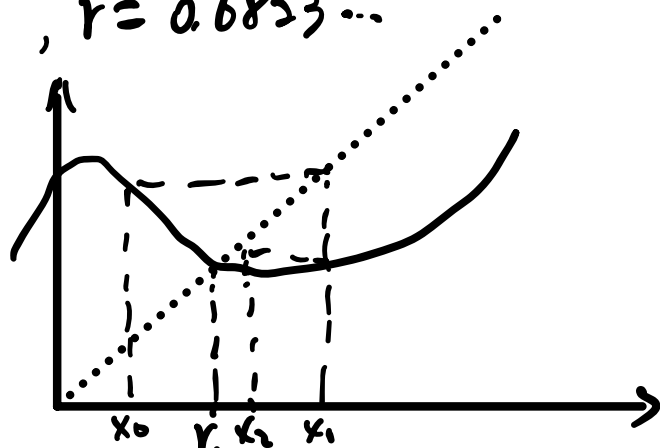
1.  $|g'(s_t)| > 1$ ,  $|e_{t+1}| > |e_t|$ , diverge.
2.  $|g'(s_t)| < 1$ ,  $|e_{t+1}| < |e_t|$ , converge.
3.  $|g'(s_t)| = 1$ , potentially stuck.

• Example:  $f(x) = x^3 + x - 1 = 0$ ,  $r = 0.6823 \dots$

$$\cdot g_1(x) = 1 - x^3 \Rightarrow g'_1(x) > 1$$

$$\cdot g_2(x) = (1 - x)^{\frac{1}{3}} \Rightarrow g'_2(x) < 1$$

$$\cdot g_3(x) = \frac{1 + 2x^3}{1 + 3x^2} \Rightarrow g'_3(x) = 0$$



Convergence:  $\exists I \subset [r-c, r+c]$ , for some  $c > 0$ , s.t.  $|g'(x)| < 1$

on  $I$  and  $x_0 \in I$ , then FPI will converge!

$$\lim_{t \rightarrow \infty} \left| \frac{e_{t+1}}{e_t} \right| = \lim_{t \rightarrow \infty} |g'(s_t)| = g'(r) \neq b < 1.$$

$|e_{t+1}| = c \cdot |e_t|$  : linear convergence!

Iteration complexity:

We want to find  $t$ , s.t.

$$|e_t| = |x_t - r| < \text{tol.} \quad \text{Assume } |g'(x)| \leq m < 1$$

and  $x_0 \in I = [x-c, x+c]$ .

$$|e_{t+1}| \leq m \cdot |e_t| \dots \leq m^{t+1} \cdot |e_0|$$

But  $e_0 = x_0 - r$  is unknown!

Note that  $|e_0| = |x_0 - r|$  und  $|e_1| = |x_1 - r| < |x_0 - r|$

$$\Rightarrow |e_0| = |x_0 - r| = |x_0 - x_1 + x_1 - r|$$

$$\leq |x_0 - x_1| + |e_1|$$

$$\leq |x_0 - x_1| + m \cdot |e_0|$$

$$\Rightarrow |e_0| \leq |x_0 - x_1| / (1 - m). \text{ So,}$$

$$\Rightarrow |e_t| \leq m^t \cdot \frac{|x_0 - x_1|}{1 - m}. \text{ Given } |e_0| \leq \xi.$$

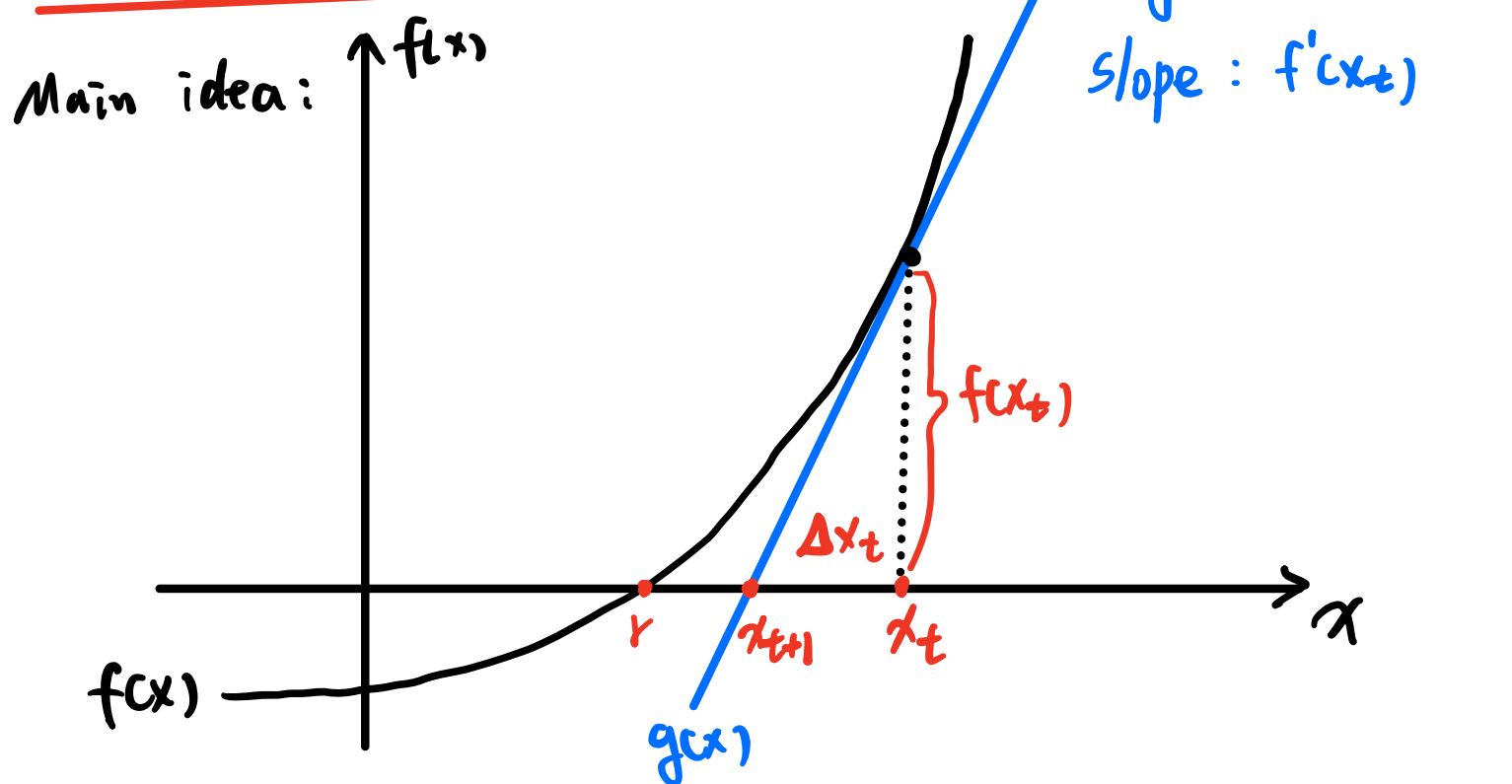
$$\Rightarrow m^t \cdot \frac{|x_0 - x_1|}{1 - m} \leq \xi \Rightarrow k \geq \left\lceil \ln \frac{\xi(1-m)}{|x_0 - x_1|} / \ln m \right\rceil.$$

• Example:  $x = g(x) = \cos x$   $|g'(x)| = |\sin x|$ .

$$m \approx 0.845, \quad x_1 = \cos 1 \approx 0.5403.$$

$$\xi = 0.5 \times 10^{-5} \Rightarrow k \geq 72. \text{ (In worst case).}$$

### ③ Newton's method



Approximate  $f(x)$  by linearization at  $x_t$  to obtain  $x_{t+1}$ .

To do this, draw the tangent line at  $x_t$  and use the intersection point of this line and  $x$ -axis as an approximation root. To measure the slope of  $g(x)$ :

$$\frac{f(x_t)}{\Delta x_t} = f'(x_t), \text{ where } \Delta x_t := x_t - x_{t+1}.$$

Then we have the following updates:

$$(2). \quad x_{t+1} = x_t - \Delta x_t = x_t - \frac{f(x_t)}{f'(x_t)} \quad \text{Newton's}$$

An alternative way: the tangent line of  $g(x)$ :

$$y - f(x_t) = f'(x_t) \cdot (x - x_t). \text{ letting } y = 0, \text{ we will get}$$

the intersection point  $\Rightarrow x = x_t - \frac{f(x_t)}{f'(x_t)} := x_{t+1}$ .

Newton's algo: Given  $f$  and  $f'$ , with  $t = 0, 1, 2, \dots$ , it generates  $\{x_t\} \rightarrow r$ , defined by (2).

Remark: Newton's algo. can be viewed as "optimal" fixed point iteration. Recall that FPI:  $x_{t+1} = g(x_t)$ , with

$$g(x) = x - \frac{f(x)}{f'(x)}, \quad g'(x) = 1 - \frac{f'(x)^2 - f(x) \cdot f''(x)}{f'(x)^2} = \frac{f(x) \cdot f''(x)}{f'(x)^2}.$$

So, assume that  $f'(r) \neq 0$ .  $\Rightarrow g'(r) = \frac{f(r) \cdot f''(r)}{f'(r)^2} = 0$ .

Note  $r$  is the root of  $f(x)$ :  $f(r) = 0$ .

Error analysis: Recall  $f(r) = 0$ ,  $e_t = x_t - r$ ,  $g(x) = x - \frac{f(x)}{f'(x)}$ .

We assume  $f''$  is continuous and  $r$  is a simple root of  $f$ , i.e.,  $f(r) = 0 \neq f'(r)$ .

From Newton's Iteration (2), we have

$$\begin{aligned} e_{t+1} &= x_{t+1} - r = x_t - \frac{f(x_t)}{f'(x_t)} - r \\ &= e_t - \frac{f(x_t)}{f'(x_t)} = \frac{e_t \cdot f'(x_t) - f(x_t)}{f'(x_t)} \end{aligned}$$

We hope this part is small!

$$= \frac{f''(\xi_t)}{2 f'(x_t)} \cdot e_t^2, \text{ where the last eqn. follows by}$$

Taylor's theorem. Recall if  $f$  is two times differentiable,

$$\text{then } f(x) = f(a) + f'(a) \cdot (x-a) + \frac{f''(\xi)}{2} (x-a)^2,$$

$$\xi \in \text{conv}(x, a). \text{ Here } \begin{cases} x := r & e_t = x_t - r \\ a := x_t & \xi_t \in \text{conv}(r, x_t). \end{cases}$$



$$0 = f(r) = f(x_t) + f'(x_t) \cdot (-e_t) + \frac{f''(\xi_t)}{2} (-e_t)^2.$$

It then leads to  $e_{t+1} = \frac{f''(\xi_t)}{2f'(x_t)} \cdot e_t^2.$

If we assume further that

1.  $f'(x) \neq 0$  for all  $x \in [r - |e_0|, r + |e_0|] := I$
2.  $f'(x)$  is continuous for all  $x \in I$ .
3.  $\exists M = \frac{1}{2} \left( \sup_{x \in I} |f''(x)| \right) \left( \sup_{x \in I} \frac{1}{|f'(x)|} \right)$  s.t.  $M|e_0| < 1$ .

If the above conditions hold, then

$$|e_{t+1}| \leq M \cdot |e_t|^2.$$

• To estimate  $M$ , given  $\delta > 0$ , let's define  $M(\delta)$ :

$$M(\delta) = \frac{1}{2} \cdot \frac{\max_{|x-r| \leq \delta} |f''(x)|}{\min_{|x-r| \leq \delta} |f'(x)|}. \quad \text{Since } f' \in C(\mathbb{R}), f'(r) \neq 0.$$

One can choose  $\delta$  s.t.  $\delta \cdot M(\delta) < 1$ . This is possible as

$$\delta \rightarrow 0, \quad M(\delta) \rightarrow \frac{|f''(r)|}{|2f'(r)|}. \quad \text{Hence, } \delta \cdot M(\delta) \rightarrow 0.$$

Assume we start a point  $x_0$  satisfy  $|x_0 - r| \leq \delta$ . Then

$$|e_0| \leq \delta. \quad |e_1| = |x_1 - r| = \frac{1}{2} \frac{|f''(\xi_0)|}{|f'(\xi_0)|} \cdot e_0^2 \leq M(\delta) \cdot e_0^2$$

$$= |e_0| \cdot |e_0| \cdot M(\delta)$$

$$\leq |e_0| \cdot \delta \cdot M(\delta) \leq |e_0| \leq \delta.$$

$$\Rightarrow |x_1 - r| \leq \delta, \quad x_1 \in (r - \delta, r + \delta).$$

$$\Rightarrow |e_t| \leq [\delta \cdot M(\delta)]^t \cdot |e_0|. \quad \Rightarrow \lim_{t \rightarrow \infty} |e_t| = 0.$$

Convergence of Newton's method: Let's assume  $f' \in C(\mathbb{R})$ , and let  $r$  be a simple root of  $f$ . Then there is a neighbor of  $r$  and a constant  $M$  s.t. if Newton's starts in that neighbor, the points  $\{x_t\}$  become steadily closer to  $r$  and satisfy  $|e_{t+1}| \leq M \cdot |e_t|^2$ , ( $t \geq 0$ ).

Compute  $\sqrt{a}$  in our last lecture:

$$x_{t+1} = \frac{1}{2} \left( x_t + \frac{a}{x_t} \right).$$

Define  $f(x) = x^2 - a$ , to find  $\sqrt{a}$  is to find  $f(r) = 0$ .

$f', f''$  are continuous and  $f'(r) = 2r = 2\sqrt{a} \neq 0$ .

$$\text{Newton's: } x_{t+1} = x_t - \frac{f(x_t)}{f'(x_t)} = x_t - \frac{x_t^2 - a}{2x_t} = \frac{1}{2} \left( x_t + \frac{a}{x_t} \right).$$

To estimate  $M$ :

$$M = \left| \frac{f''(r)}{2f'(r)} \right| = \frac{1}{2r} = \frac{1}{2\sqrt{a}}.$$

Handle multiple roots:

Def.  $r \in \mathbb{R}$  is a root of multiplicity  $m$  of  $f(x)$  if there is

a polynomial  $s(x)$  s.t.  $s(r) \neq 0$  and  $f(x) = (x-r)^m \cdot s(x)$ .

$\left\{ \begin{array}{l} m=1 : r \text{ is a simple root.} \\ m \geq 2 : r \text{ is a multiple root.} \end{array} \right.$

$$\Leftrightarrow f(r) = 0, f^{(i)}(r) = 0, i=1, 2, \dots, m-1, f^{(m)}(r) \neq 0.$$

Consider  $f(x) = x^m$  ( $m \geq 2$ ),  $f$  has multiple roots at  $x=0$ .  
 $f'(x) = m \cdot x^{m-1} = 0$ . Above analysis for simple root does not work for this case!

$$x_{t+1} = x_t - \frac{f(x_t)}{f'(x_t)} = x_t - \frac{x_t^m}{m \cdot x_t^{m-1}} = \frac{m-1}{m} \cdot x_t$$

$$\Rightarrow e_{t+1} = \frac{m-1}{m} \cdot e_t \Rightarrow e_t = \left(\frac{m-1}{m}\right)^t \cdot e_0. \text{ It is slow.}$$

Q: In the multiple case, can we fix Newton's so that it converges faster?

If  $f(x) = (x-r)^m \cdot s(x)$ , then we have

$$|e_{t+1}| \lesssim \frac{m-1}{m} \cdot |e_t|. \text{ To show this,}$$

$$x_{t+1} = x_t - \frac{f(x_t)}{f'(x_t)} \Leftrightarrow e_{t+1} = e_t - \frac{f(x_t)}{f'(x_t)}$$

$$e_{t+1} = e_t - \frac{(x_t - r)^m \cdot s(x_t)}{m(x_t - r)^{m-1} \cdot s(x_t) + (x_t - r)^m \cdot s'(x_t)}$$

$$= e_t - \frac{e_t \cdot s(x_t)}{m \cdot s(x_t) + s'(x_t) \cdot e_t}$$

$$= \left[ \frac{(m-1) \cdot s(x_t) + s'(x_t) \cdot e_t}{m \cdot s(x_t) + s'(x_t) \cdot e_t} \right] \cdot e_t$$

$$= \frac{(m-1) \cdot f(x_t)}{f'(x_t)} + \frac{s'(x_t)}{m \cdot s(x_t) + e_t \cdot s'(x_t)} \cdot e_t^2$$

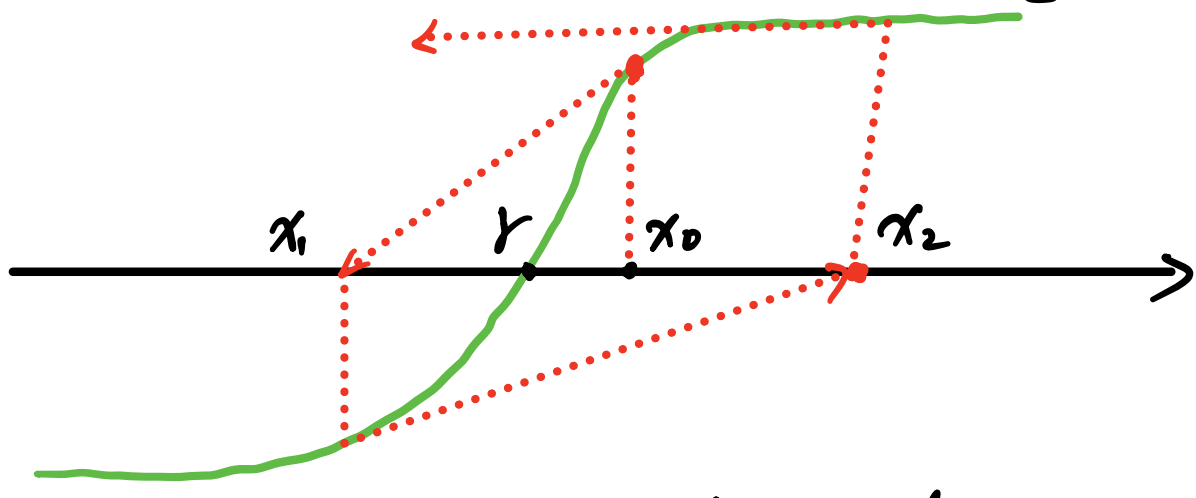
Fixing it: Just need to add  $-\frac{(m-1) \cdot f(x_t)}{f'(x_t)}$  on right side

$$\Rightarrow e_{t+1} = O(e_t^2). \text{ So, a fixed version is:}$$

$$x_{t+1} = x_t - \frac{m \cdot f(x_t)}{f'(x_t)}$$

## Summary of Newton's:

- For simple root, it converges quadratically when  $x_0$  is close enough to  $r$ .
- For multiple roots, it converges linearly but can be fixed.
- When  $x_0$  is far from  $r$ , it may diverge!

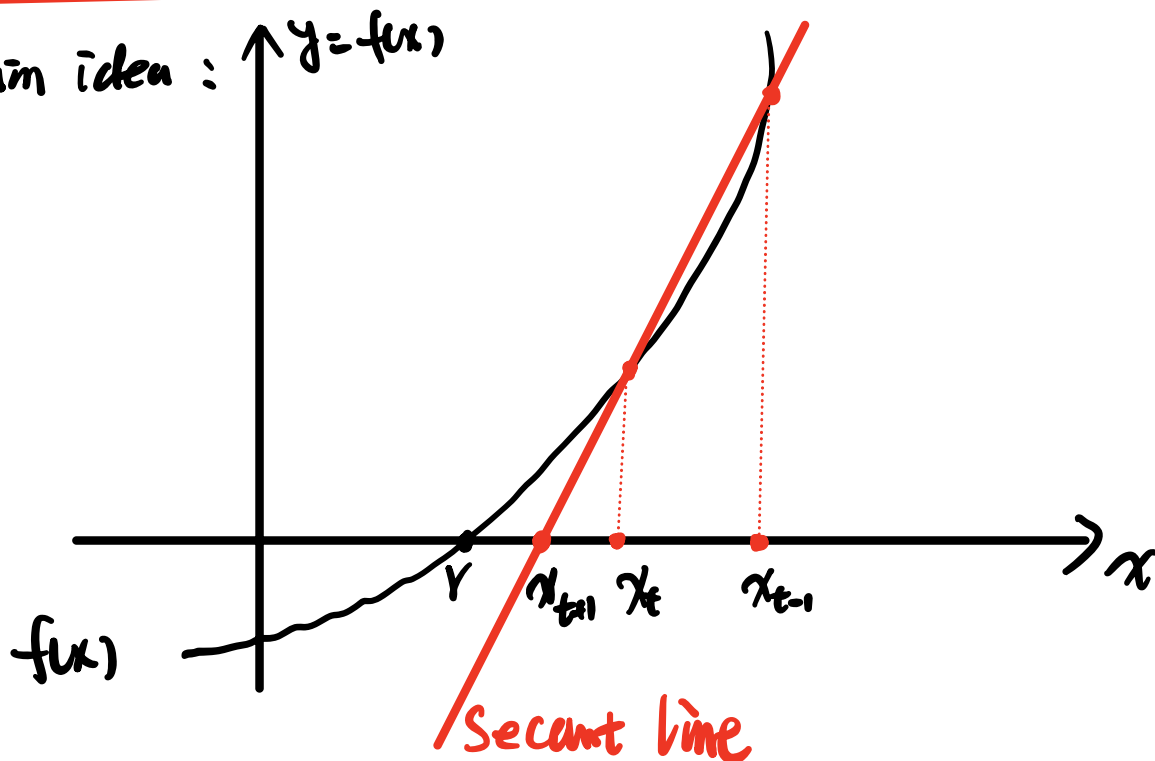


- One can use Bisection to find a good  $x_0$  and then apply Newton's.
- One needs to have  $f'(x)$  and  $f'(x) \neq 0$ !

Q: If  $f'(x)$  is hard to obtain, can we find a good method s.t. it is superlinear?

## ④ Quasi-Newton method: Secant method

Main idea:



To approximate the derivative at the point  $x_t$  with previous point  $x_{t-1}$ .

$$f'(x_t) \approx \frac{f(x_t) - f(x_{t-1})}{x_t - x_{t-1}} : \text{difference quotient.}$$

(Note  $f'(x) = \lim_{u \rightarrow x} \frac{f(x) - f(u)}{x - u}$ ). Hence, we have

$$x_{t+1} = x_t - f(x_t) \cdot \left[ \frac{x_t - x_{t-1}}{f(x_t) - f(x_{t-1})} \right] \quad \text{Secant}$$

$$= \frac{x_{t-1} \cdot f(x_t) - x_t \cdot f(x_{t-1})}{f(x_t) - f(x_{t-1})}$$

Per-iteration only requires 1 function evaluation.

Error analysis: Consider  $e_t = x_t - r$  and  $f(r) = 0$ ,  $r$  is simple.

$$e_{t+1} = x_{t+1} - r = \frac{x_{t-1} \cdot f(x_t) - x_t \cdot f(x_{t-1})}{f(x_t) - f(x_{t-1})} - r = \frac{e_{t-1} \cdot f(x_t) - e_t \cdot f(x_{t-1})}{f(x_t) - f(x_{t-1})}$$

It leads to  $e_{t+1} = \frac{f(x_t)/e_t - f(x_{t-1})/e_{t-1}}{f(x_t) - f(x_{t-1})} \cdot e_t e_{t-1}$

$$\text{Def. } h(x) = \begin{cases} \frac{f(x) - f(r)}{x - r}, & x \neq r \\ f'(r), & x = r \end{cases} \Rightarrow h'(x) = \begin{cases} \frac{f'(x)(x-r) - f(x)}{(x-r)^2} & x \neq r \\ \frac{1}{2} f''(r) & x = r. \end{cases}$$

$$= \frac{h(x_t) - h(x_{t-1})}{f(x_t) - f(x_{t-1})} \cdot e_t e_{t-1}$$

$$\frac{h(x_t) - h(x_{t-1})}{f(x_t) - f(x_{t-1})} = \frac{h(x_t) - h(x_{t-1})}{x_t - x_{t-1}} \cdot \frac{x_t - x_{t-1}}{f(x_t) - f(x_{t-1})} = \frac{h'(\xi_t)}{f'(\xi_t)}$$

$\eta_t, \xi_t \in \text{conv}(x_t, x_{t-1})$ . We assume  $f(x) \in C^2[\mathbb{R}]$ .

$$f(x) = f(a) + f'(a) \cdot (x-a) + \frac{f''(\xi)}{2} (x-a)^2$$

$\xi \in \text{conv}(x, a)$ . Here  $\begin{cases} x := r \\ a := \eta_t, \xi_t \in \text{conv}(r, \eta_t). \end{cases}$

$$\Rightarrow f(r) = f(\eta_t) + f'(\eta_t)(r - \eta_t) + \frac{(r - \eta_t)^2}{2} \cdot f''(\xi_t)$$

$$\frac{1}{2} f''(\xi_t) = \frac{f'(\eta_t)(\eta_t - r) - f(\eta_t)}{(\eta_t - r)^2} = h'(\eta_t)$$

$$\Rightarrow e_{t+1} = \frac{f''(\xi_t)}{2 f'(\xi_t)} \cdot e_t \cdot e_{t-1}, \quad \begin{cases} \xi_t \in \text{conv}(x_t, x_{t-1}, r) \\ \eta_t \in \text{conv}(x_t, x_{t-1}). \end{cases}$$

Note  $x_t - x_{t-1} = e_t - e_{t-1}$ .

If  $x_0, x_1, x_2$  are close enough to  $r$  and

$\left| \frac{f''(\xi_t)}{2 f'(\xi_t)} \right| \leq M$  is properly bounded. Suppose that

$|e_1| \leq \frac{1}{2M}$  and  $|e_0| \leq \frac{1}{2M}$ . Then  $|e_2| \leq M \cdot |e_1| \cdot |e_0|$

$$\leq \frac{1}{2} |e_0| \text{ and } \leq \frac{1}{2} |e_1| \Rightarrow$$

$$|e_2| \leq \frac{1}{2} \min \{ |e_0|, |e_1| \} \leq \frac{1}{2^2 \cdot M}$$

By induction,  $|e_{t+1}| \leq \dots \leq \frac{1}{2^{t+1} \cdot M}$ ,  $\{e_t\} \rightarrow$ .

$$\left| \frac{f''(\xi_t)}{2f'(\xi_t)} \right| \rightarrow \left| \frac{f''(\eta)}{2f'(\eta)} \right| \stackrel{\Delta}{=} C.$$

$|e_{t+1}| \sim C \cdot |e_t| \cdot |e_{t-1}|$ . We assume that  $|e_{t+1}| \sim A |e_t|^d$ ,

$A > 0$ , meaning  $\lim_{t \rightarrow \infty} \frac{|e_{t+1}|}{A |e_t|^d} = 1 \Rightarrow$

$$\Leftrightarrow A |e_t|^d \sim C \cdot |e_t| \cdot |e_{t-1}|$$

$$\Rightarrow |e_t|^{d-1} \sim C \cdot A^{-1} \cdot |e_{t-1}|$$

$$\Rightarrow |e_t| \sim (C \cdot A^{-1})^{\frac{1}{d-1}} \cdot |e_{t-1}|^{\frac{1}{d-1}}$$

$$\Rightarrow A^{1+\frac{1}{d-1}} = C^{\frac{1}{d-1}}, \quad d = \frac{1}{d-1} \quad (d > 0, \Rightarrow d = \frac{1+\sqrt{5}}{2})$$

$$\approx 1.6180\dots$$

$$|e_{t+1}| \sim C \cdot |e_t| \cdot |e_{t-1}| \Leftrightarrow A = C^{\frac{1}{d}} = C^{\frac{d-1}{d}}$$

That is  $|e_{t+1}| \sim \left( \frac{\sqrt{5}-1}{2} \right)^{\frac{1+\sqrt{5}}{2}} \cdot |e_t|^{\frac{1+\sqrt{5}}{2}}$

$$\approx \left| \frac{f''(\eta)}{2f'(\eta)} \right|^{\frac{\sqrt{5}-1}{2}} \cdot |e_t|^{\frac{1+\sqrt{5}}{2}}$$

Multiple roots:

One can reduce to the following error iteration:

$$e_{t+1} = \frac{e_{t-1} e_t^m - e_t e_{t-1}^{m-1} + O(|e_{t-1}|^{m+2})}{e_t^m - e_{t-1}^m + O(|e_{t-1}|^{m+1})} \quad (m \geq 2)$$

- $\Leftrightarrow$ 
  1.  $m=2$ , it is a Fibonacci sequence:  $|e_{t+1}| \sim (\lambda+1) \cdot |e_t|$
  2.  $m \geq 3$ ,  $|e_{t+1}| \sim \lambda |e_t|$  with  $\lambda \cdot (\lambda-1) \cdot (\lambda^m + \lambda^{m-1} - 1) = 0$ ,  $\lambda \in (0, 1)$ .

### f) Stability of solving nonlinear equation:

To design the stop conditions of above methods, one can consider two types of errors:

- Forward error (FE):  $|x_t - r|$
- Backward error (BE):  $|f(x_t) - f(r)| = |f(x_t)|$ .

To get approximate  $x_t$ , the forward error comes from the algorithm while the backward error comes from  $f$  itself.

• **Example 1:**  $f(x) = (x - \frac{2}{3})^3 = x^3 - 2x^2 + \frac{4}{3}x - \frac{8}{27}$

After 16th iteration of Bisection:

BE  $\approx 2.0 \times 10^{-6}$  while FE  $\approx 10^{-5}$ .

Why? Note  $|f(x_t)| = (x_t - \frac{2}{3})^3 = e_t^3$ .



• **Example 2:**  $f(x) = \sin(x) - x$ . Let  $x_t = 10^{-3}$ ,  $r = 0$

$$BE: |f(x_t)| = |\sin(0.001) - 0.001|$$

$$\sin(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \frac{x^9}{9!} - \dots$$

$$|f(x_t)| = |\sin(x_t) - x_t| = O(x_t^3) \approx 1.62 \times 10^{-10}$$

$$FE = |x_t - r| = 10^{-3}$$

Sensitivity: Small floating point errors in the function  $f$  translate into large errors in the root! It is called a sensitive numerical problem.

To measure the sensitivity, assume  $f(r) = 0$ . There is a small change is made on  $f(x)$ :  $\epsilon \cdot g(x)$  where  $\epsilon$  is small value. We actually find:

$$f(x) + \epsilon \cdot g(x) = 0. \quad \text{Def. } F(x) = f(x) + \epsilon \cdot g(x)$$

Let  $r + \Delta r$  is the root of  $F(x)$ . That is

$$F(r + \Delta r) = f(r + \Delta r) + \epsilon \cdot g(r + \Delta r) = 0,$$

where  $\Delta r$  is the change of  $r$  due to small error  $\epsilon$ .

Use Taylor polynomial of  $f$ :

$$\left\{ \begin{array}{l} f(r + \Delta r) = f(r) + \Delta r \cdot f'(r) + O(|\Delta r|^2) \\ g(r + \Delta r) = g(r) + \Delta r \cdot g'(r) + O(|\Delta r|^2) \end{array} \right.$$

$$\text{So, } 0 = f(r + \Delta r) + \xi \cdot g(r + \Delta r)$$

$$= \Delta r \cdot f'(r) + \Delta r \cdot \xi \cdot g'(r) + \xi \cdot g(r) + O(|\Delta r|^2)$$

$$\Rightarrow \Delta r \cdot (f'(r) + \xi \cdot g'(r)) \approx -\xi \cdot g(r)$$

$$\Rightarrow \Delta r \approx \frac{-\xi \cdot g(r)}{f'(r) + \xi \cdot g'(r)} \approx -\xi \cdot \frac{g(r)}{f'(r)}$$

Sensitivity of roots.

(assume  $\xi \cdot g'(r) \ll f'(r)$ ).

Example: Say we want to estimate polynomial

$$p(x) = \prod_{i=1}^6 (x-i), \text{ there are system errors at the final}$$

estimate is  $p(x) + \xi \cdot x^2$ . Find sensitivity.

Solution: Let  $f(x) = \prod_{i=1}^6 (x-i)$ ,  $\xi = -10^{-6}$  and  $g(x) = x^2$

$$\Delta r \approx -\xi \cdot \frac{g(r)}{f'(r)} = -\frac{\xi \cdot 6^2}{5!} = -2332.8 \cdot \xi$$

(Note  $f'(x) = \sum_j (x-b_j) \cdot h_j(x) + \prod_{i=1}^5 (x-i)$ , so  $f'(6) = 5!$ )

$$10^{-6} \cdot 6^2 \approx 0.2799$$

the estimated root is  $r + \Delta r = 6.0023328$ .

$\Rightarrow$  6 digits of  $f(x)$  will cause 3 digits change error of root!

**Error Magnification Factor (EMF):**

Idea: when bad cases happen, quantity  $FE/\beta E$  be large.

So, a reasonable way to measure this:

$$EMF = \left| \frac{\text{rel. FE}}{\text{rel. BE}} \right| = \left| \frac{\Delta r / r}{\sum g_{ur} / g_{ur}} \right|, \quad \frac{\Delta r}{r} \approx \left| \frac{\sum g_{ur}}{r \cdot f'(r)} \right|$$

$$= \left| \frac{g_{ur}}{r \cdot f'(r)} \right| \quad \text{Note: rel. BE: } \left| \frac{f(r) + \sum g_{ur} - f(r)}{g_{ur}} \right|$$

is relative to  $g$ .

EMF is highly related to condition number.

• Example:  $W(x) = \prod_{i=1}^{20} (x-i)$ . Use the sensitivity formula to estimate the root change in  $x^{15}$  term of  $W(x)$  on root  $r=16$ . Then find EMF.

$$W_2(x) = W(x) + \epsilon \cdot g_{ur}, \quad g_{ur} = -1,672,280,820 \cdot x^{15}$$

$$W'(x=16) = 15! 4!$$

$$\Delta r \approx \left| \frac{\sum g_{ur}}{f'(r)} \right| \approx 6.1432 \times 10^{13} \cdot \epsilon \quad (\epsilon_m = \pm 2.22 \times 10^{-16})$$

$$\Delta r \approx \pm 0.0136.$$

$$\left| \frac{g_{ur}}{r \cdot f'(r)} \right| = \frac{16^{15} \cdot 1,672,280,820}{15! 4! 16} \approx 3.8 \times 10^{12}$$

$$(W(x) = (x-16) \cdot q_1(x), \quad q_1(x) = \prod_{i=1}^{15} (x-i) \cdot \prod_{j=7}^{20} (x-j))$$

$$W'(x) = q_1(x) + (x-16) \cdot q_1'(x), \Rightarrow W'(16) = q_1(16)$$