Condition number and limits of accuracy $\boxed{\text{cond}}$

Problems

① find root $r$ of $f(x)$

② given $\vec{u}, \vec{v}$ compute $\vec{u}^T \vec{v}$

③ given $A\vec{x} = \vec{b}$, compute solution $\vec{x}$

Model problem: evaluate $y = h(t)$ (scalar function)

① "$r = h(f)$"

② $\vec{u}^T \vec{v} = h(\vec{u}, \vec{v})$

③ $\vec{x} = h(A, \vec{b})$

$\boxed{\text{Important to understand inputs/outputs}}$

Each problem has a <u>condition #</u> which measures how sensitively output depends on input. In practice, problems w/ large cond #'s are hard to solve numerically.

<u>Model problem</u>

$$y = h(t)$$

$$y + \delta y = h(t + \delta t)$$

so $\delta y = \dfrac{\delta y}{\delta t} \delta t = \dfrac{h(t+\delta t) - h(t)}{\delta t} \delta t \simeq h'(t)\, \delta t$

$$\hat{K}(t) = \max_{\delta t}\left|\frac{\delta y}{\delta t}\right| = |h'(t)|$$

if $h(t)$ differentiable

absolute condition #.

if $h(t)$ differentiable

$$K(t) = \max_{\delta t}\frac{|\delta y/y|}{|\delta t/t|} = \left|\frac{t\,h'(t)}{h(t)}\right|$$

relative condition #

BACKWARD SIDE                    FORWARD SIDE

input $\longrightarrow$ | solution process | $\longrightarrow$ output

$t \longrightarrow$ | process to evaluate $h(t)$ | $\longrightarrow y$

We speak of <u>forward</u> and <u>backward</u> errors.

We want $y = h(t)$, we get $y_A = h_A(t)$

Accuracy

$$\left| \frac{y_A - y}{y} \right| = O(\varepsilon_{mach})$$

$$\left| \frac{h_A(t) - h(t)}{h(t)} \right|$$

Backward stability

$$h_A(t) = h(t + \delta t) \quad \text{for small } \delta t.$$

$$\left| \frac{\delta t}{t} \right| = O(\varepsilon_{mach})$$

Example $h(t) = at$ $a \in \mathbb{R}$ fixed

$$h_A(t) = fl(a) \odot fl(t)$$

$$= a(1+\alpha) \odot t(1+\beta) \qquad \alpha, \beta = O(\varepsilon_{mach})$$

$$\left| \frac{f(t) - t}{t} \right| \leq \frac{1}{2} \varepsilon_{mach}$$

$$= at(1+\alpha)(1+\beta)(1+\gamma)$$

$$\left| \frac{y_A - y}{y} \right| = \left| \frac{at(1+\alpha)(1+\beta)(1+\gamma) - at}{at} \right| \qquad \underline{accurate}$$

$$= \left| (1+\alpha)(1+\beta)(1+\gamma) - 1 \right| = O(\varepsilon_{mach})$$

Also backward stable.

$$h_A(t) = at(1+\alpha)(1+\beta)(1+\gamma)$$

$$= a\left[t(1+\alpha)(1+\beta)(1+\gamma) - t + t\right]$$

$$\underbrace{\phantom{t(1+\alpha)(1+\beta)(1+\gamma) - t}}_{\delta t}$$

$$= a(t + \delta t)$$

$$= h(t + \delta t) \qquad \text{"exact answer for perturbed input"}$$

where $\delta t = t\left[(1+\alpha)(1+\beta)(1+\gamma) - 1\right]$

$$\underbrace{\phantom{(1+\alpha)(1+\beta)(1+\gamma) - 1}}_{O(\varepsilon_{mach})}$$

so $|\delta t / t| = O(\varepsilon_{mach})$

# Backward stable algorithms are accurate

$$\left| \frac{h_A(t) - h(t)}{h(t)} \right| = \left| \frac{h(t + \delta t) - h(t)}{h(t)} \right|$$

$$= |\delta y / y|$$

$$= \frac{|\delta y / y|}{|\delta t / t|} \left| \frac{\delta t}{t} \right| \qquad O(\varepsilon_{mach})$$

$$\leq \kappa(t)$$

$$= \kappa(t) \, O(\varepsilon_{mach})$$

<u>root finding</u>   Given $f(x)$ find $r$ s.t.
$f(r) = 0$. <u>Solve $f(x) = 0$</u>.
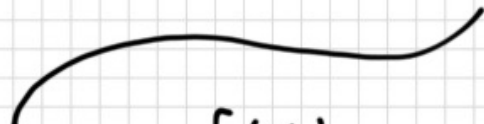
$r = h(f)$

$r_A = h_A(f) \underset{?}{=\!=} h(f + \delta f)$

always possible

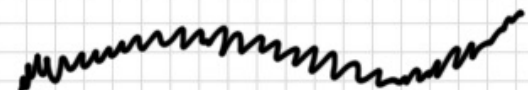$f(x) + \delta f(x) = f(x) - \underbrace{f(r_A)}_{\delta f}$

Better model
$f(x) + \delta f(x) = f(x) + \mathcal{E}_{mach} g(x)$

$f(x)$

$f(x) + \varepsilon_{mach}\, g(x)$

$$f(r_A) + \varepsilon_{mach}\, g(r_A) = 0.$$

so we find root of $f + \delta f$

Sensitivity formula for roots

forward error (relative) $\left|\dfrac{r_A - r}{r}\right|$

backward error

To compute forward error consider

$$f(r_A) + \varepsilon g(r_A) = 0$$

$$f(r + \delta r) + \varepsilon g(r + \delta r) = 0$$

$$\underbrace{f(r)}_{0} + f'(r)\delta r + \varepsilon g(r) + \underbrace{\varepsilon g'(r)\delta r}_{small} + \underbrace{O(\delta r^2)}_{small} = 0$$

small

$$\delta r \simeq -\frac{\varepsilon g(r)}{f'(r)}$$

$$\left| \frac{r_A - r}{r} \right| \simeq \varepsilon \left| \frac{g(r)}{r f'(r)} \right|$$

error magnification = $\dfrac{\text{relative forward error}}{\text{backward error}}$

$$f(r_A) + \varepsilon \, g(r_A) = 0$$

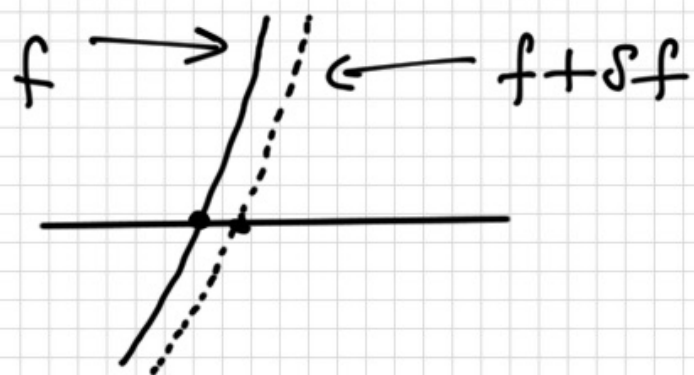so $|f(r_A)| = \varepsilon \, |g(r_A)|$

$$\simeq \varepsilon \, |g(r)|$$

$$= \left| \dfrac{1}{r \, f'(r)} \right|$$

result makes sense w/ pictures

$f \longrightarrow$      $\longleftarrow f + \delta f$

Ex   $f(x) = (x-1)(x-2)(x-3)$

$\delta f(x) = \varepsilon x^4$ where $\varepsilon = 10^{-7}$

formula ok
for any $\varepsilon$ but   $\varepsilon_{mach} \approx 2 \times 10^{-16}$

What happens to $r = 2$.

$$f'(x) = \frac{d}{dx}\left[(x-1)(x-2)(x-3)\right]$$

$$f'(2) = -1 \; ; \; g(x) = x^4, \; g(2) = 16$$

$$\delta r = -\varepsilon\frac{g(2)}{f'(2)} = 1.6 \times 10^{-6}$$

$$r_A = r + \delta r \simeq 2.0000016 \qquad \text{error mag}$$

$$= \left|\frac{1}{2f'(2)}\right| = \frac{1}{2}$$

Problem: find root $r$ of $f(x) = 0$

$$r = h(f) \qquad \text{"function point of view"}$$

output $\uparrow$ $\qquad$ $\uparrow$ input

Actually compute $r_A = h_A(f)$

$\uparrow$ alogorithm to evaluate $h$ (solve EQN)

(relative) forward error

$$\left| \frac{r_A - r}{r} \right| \qquad \longleftarrow \text{what we want small}$$

backward error $|f(r_A)| \longleftarrow$ what we can test

EX  $f(x) = (x - 4/3)^3$

say  $r_A = \frac{4}{3} + 10^{-P}$
$$\left( \begin{array}{l} \text{so}\ \ r_A = 4/3 + 6.1 \\ \qquad\ = 4/3 + 0.01 \\ \quad\ \text{etc} \end{array} \right)$$

Backward error  $|f(r_A)| = 10^{-3P}$

Forward error
(relative)  $\dfrac{|10^{-P}|}{|4/3|} = \dfrac{3}{4} 10^{-P}$

error mag  $= \dfrac{3}{4} 10^{2P}$  $\left( \begin{array}{l} \text{so}\ \ \frac{3}{4} 100 \\ \ \frac{3}{4} 10,000 \end{array} \right)$

Numerically $\quad x^3 - 4x^2 + \frac{16}{3}x - \frac{64}{27}$

$[r_A \; h] = \text{bisection} (@(x)(x-4/3)^3, 1, 2, 1e-8)$

$\phantom{[r_A \; h] = \text{bisection} (@(x)(x-4/3)^3,} a \; b \; \text{tol}$

$\phantom{[r_A \; h] =} = 1.333333335581686$

$\text{abs}((r_A - r)/r) = 1.8626 \, e-9$

$|(r_A - 4/3)^3| = 1.5318 \, e-26$

$\text{error mag} = 1.2160 \, e+17$

Created with Doceri

Numerically

$$[r_A, k] = \text{bisection}(@\sin, 3pi/4, 3pi/2, 1e-8)$$

$$\frac{3\pi}{4} \quad \frac{3\pi}{2} \quad \text{tol}$$

$$= 3.14159265663956$$

$$\text{abs}((r_A - pi)/pi) = 9.3132e-10$$

$$\text{abs}(\sin(r_A)) = 2.9258e-9$$

$$\text{error mag} = 0.3183$$

Created with Doceri

The approximate root $r_A = h_A(f)$

$$= h(f + \delta f)$$

results from two $\longrightarrow$ Perturbation of function

Sources ① round-off errors in function evaluations

$\underline{f + \delta f_① + \delta f_②}$
$\tilde{f}$

② details and roundoff errors in algorithm.

Can't really see ①. Backward error $|f(r_A)|$

will be nonzero due to ②