

Sep 16, 2020

Last time: floating point "1 + δ " model

$$f(x \odot y) = (x \odot y)(1 + \delta) \quad |\delta| \leq \epsilon \begin{cases} 10^{-16} & 64\text{-bit} \\ 10^{-7} & 32\text{-bit} \end{cases}$$

$$\odot \in \{+, -, *, /\}$$

$$\frac{|f(x \odot y) - x \odot y|}{|x \odot y|} = |\delta| \leq \epsilon$$

Given two FP numbers x, y , $f_{\odot}(x, y)$ is accurate

Alg. \tilde{f} for f is accurate if

$$\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} = O(\epsilon)$$

What if
ill-conditioned?

$$\frac{\|\delta f(x)\|}{\|f(x)\|} / \frac{\|\delta x\|}{\|x\|} \text{ could be large}$$

Example: subtraction with FP reps

$$\hat{x} = f(x) = x(1 + \delta_1) \quad |\delta_1| \leq \epsilon$$

$$\hat{y} = f(y) = y(1 + \delta_2) \quad |\delta_2| \leq \epsilon$$

$$f(\hat{x} - \hat{y}) = (\hat{x} - \hat{y})(1 + \delta_3)$$

$$= x(1 + \delta_1)(1 + \delta_3) - y(1 + \delta_2)(1 + \delta_3)$$

$$= \boxed{x(1 + 2\delta_4) - y(1 + 2\delta_5)}$$

$$(|\delta_4|, |\delta_5| \leq \epsilon + O(\epsilon^2) = O(\epsilon))$$

$$= x - y + x2\delta_4 - y2\delta_5$$

$$\underbrace{|f(\hat{x} - \hat{y}) - (x - y)|}_{|x - y|} = \underbrace{|x2\delta_4|}_{O(\sqrt{\epsilon})} - \underbrace{|y2\delta_5|}_{O(\sqrt{\epsilon})} \quad \begin{array}{l} x, y = O(\frac{1}{\sqrt{\epsilon}}) \\ x - y = O(1) \end{array}$$

$$= O(\sqrt{\epsilon}) \quad (\text{not accurate})$$

Problem comes from FP representation error

$$f(\hat{x} - \hat{y}) = \hat{x} - \hat{y} \quad \text{if } 1/2 < \hat{x}/\hat{y} < 2 \quad (\text{HW2})$$

Still, we had: $f(\hat{x} - \hat{y}) = x(1+2\delta_4) - y(1+2\delta_5) = \tilde{x} - \tilde{y}$

$$\frac{|\tilde{x} - x|}{|x|} = O(\epsilon) \quad \frac{|\tilde{y} - y|}{|y|} = O(\epsilon)$$

This is called backward stability

Alg \tilde{f} for f is backward stable if,

for input x : $\tilde{f}(x) = f(\tilde{x}) \quad \frac{\|x - \tilde{x}\|}{\|x\|} = O(\epsilon)$

\tilde{f} always solves a nearby problem exactly

Slightly weaker: plain stability

$$\frac{\|\tilde{f}(x) - f(\tilde{x})\|}{\|f(\tilde{x})\|} = O(\epsilon), \quad \frac{\|x - \tilde{x}\|}{\|x\|} = O(\epsilon)$$

Last time: sums

$$(f) \quad s = \sum_{i=1}^n x_i$$

(f) $s = 0$
for $i = 1:n$
 $s += x_i$
return s

$\hat{s}_k = k^{\text{th}}$ partial sum in alg
returns \hat{s}_n

$$\hat{s}_0 = 0, \quad \hat{s}_1 = x_1, \quad \hat{s}_k = (\hat{s}_{k-1} + x_k)(1 + \delta_k)$$

$$|\delta_k| \leq \epsilon$$

$$\hat{s}_2 = (x_1 + x_2)(1 + \delta_2)$$

$$\hat{s}_3 = (\hat{s}_2 + x_3)(1 + \delta_3) = (x_1 + x_2)(1 + \delta_2)(1 + \delta_3) + x_3(1 + \delta_3)$$

$$\hat{s}_n = (x_1 + x_2) \prod_{j=2}^n (1 + \delta_j) + x_3 \prod_{j=3}^n (1 + \delta_j) + \dots + x_n \prod_{j=n}^n (1 + \delta_n)$$

$$\tilde{x}_1 (1 + (n-1)\delta'_1) + \tilde{x}_2 (1 + (n-1)\delta'_2) + x_3 (1 + (n-2)\delta'_3) + \dots + x_n (1 + \delta'_n)$$

$$= f(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)$$

$$\|x - \tilde{x}\|_\infty \leq \|x\|_\infty \cdot \delta \quad |\delta| < \epsilon$$

backward stable

$$\|x - \tilde{x}\|_\infty / \|x\|_\infty = O(\epsilon) \quad (\text{dimension in big } O)$$

$$f(A, b) \rightarrow z = A^{-1}b \quad (Az = b)$$

$$\star \sup \frac{\| \delta z \|}{\| z \|} / \frac{\| \delta A \|}{\| A \|} = \| A \| \cdot \| A^{-1} \| = \kappa(A)$$

Suppose we have a backward stable alg \tilde{f} for $z = A^{-1}b$

$$\tilde{f}_b(A) = f_b(\tilde{A}) \quad \text{for } \frac{\| \tilde{A} - A \|}{\| A \|} = O(\epsilon)$$

$$\| \tilde{f}_b(A) - f_b(A) \| \not\sim \| f_b(A) \|$$

$$= \| f_b(\tilde{A}) - f_b(A) \| / \| f_b(A) \|$$

$$= \underbrace{\| f_b(A + \delta A) - f_b(A) \|}_{\delta z} / \underbrace{\| f_b(A) \|}_z$$

$$\star \leq \kappa(A) \cdot \frac{\| \delta A \|}{\| A \|} = O(\kappa(A) \epsilon)$$

indep. of alg

backward stability

Can we always have backward stable algos?

$$\text{svd: } A \xrightarrow{f} (U, \Sigma, V^T)$$

$$U^T U = I, \Sigma = (\diagdown), V^T V = I$$

$$\boxed{A} \rightarrow \boxed{U} \boxed{\Sigma} \boxed{V^T}$$

carry out some arithmetic ...

$$\hat{U} \hat{\Sigma} \hat{V}^T \stackrel{?}{=} \text{svd}(\tilde{A})$$

$$\frac{\|\tilde{A} - A\|}{\|A\|} = O(\epsilon)$$

Need: $\hat{U}^T \hat{U} = I$ $\hat{U} = U + E$ $(U+E)^T(U+E) = \frac{U^T U}{I} + E^T U + U^T E + E^T E$
 $\|E\| = O(\epsilon)$

Alternative: forward error analysis?

① problem

② work thru flops

③ show that solution is close

Solving $Ax = b$ depends on $k(A)$