

MACHINE EPSILON: eps

$$\text{eps}_{\text{SP}} = 2^{-24} \sim 6 \cdot 10^{-8}$$

$$\text{eps}_{\text{DP}} = 2^{-53} \sim 10^{-16}$$

$f_l: \mathbb{R} \rightarrow \mathbb{F}$ Maps real #'s onto the floating point set
countably infinite set

For all $x \in \mathbb{R} \exists x' \in \mathbb{F}$ st

$$\frac{|x - x'|}{|x|} \leq \text{eps}$$

\therefore For all $x \in \mathbb{R} \exists \varepsilon \text{ st } |\varepsilon| \leq \text{eps}$

$$f_l(x) = x(1 + \varepsilon).$$

The difference between a floating point number and a real number closest to it is always smaller than eps .

FLOATING POINT ARITHMETIC

The operations $+, -, \times, \div$ have

Computer analogues: $\oplus, \ominus, \otimes, \div \equiv *$

$$x \otimes y = \text{fl}(x * y)$$

for all $x, y \in \mathbb{F} \exists \varepsilon | \varepsilon | < \text{eps}$ st

$$x' \otimes y' = (x * y)(1 + \varepsilon)$$

where $x, y \in \mathbb{R}$

\therefore every fl arithmetic operation is exact up to a relative error of size, at most, eps .

Rule: Complex floating point arithmetic is done via software. //

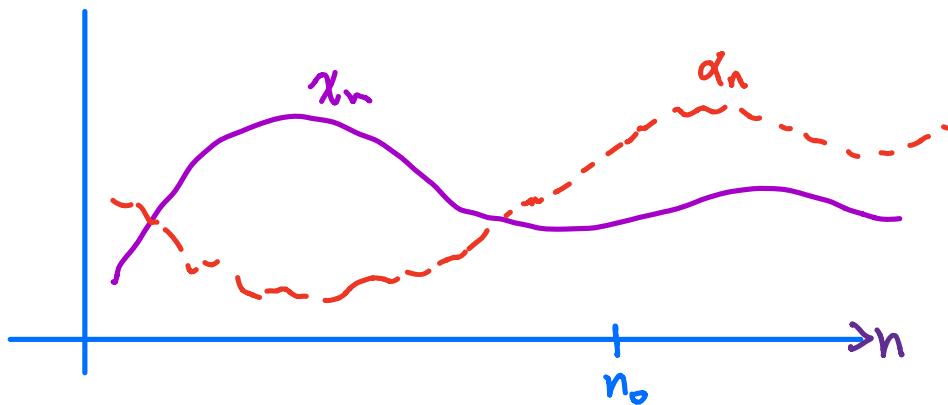
LECTURE 14 Stability

Big "Oh" & Little σ Notation:

Suppose $[x_n], [d_n]$ 2 sequences

if $x_n = O(d_n)$

Big "Oh"



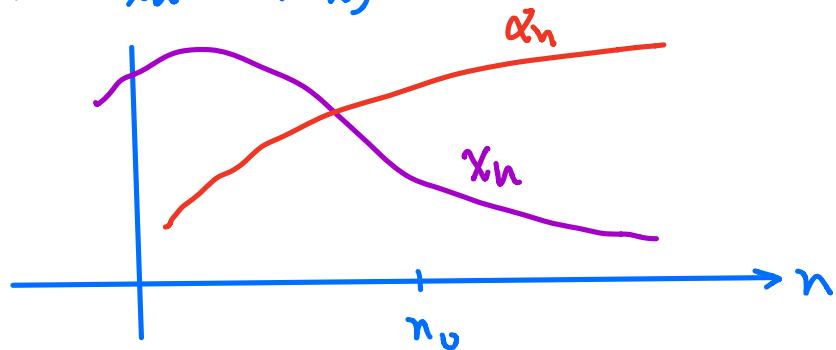
$\exists K, \text{constant and } n_0 \text{ constant s.t.}$

$$|x_n| \leq K |\alpha_n| \text{ for } n \geq n_0$$

$$\lim_{n \rightarrow \infty} \frac{|x_n|}{|\alpha_n|} \leq K$$



IF $x_n = G(\alpha_n)$ "little oh"



$$\text{if } \lim_{n \rightarrow \infty} \frac{|x_n|}{|\alpha_n|} = 0$$

Every G is G but not the other way around.



$$\text{ex)} \quad h(n) = 1 + 4e^{2n}$$

$$g(n) = e^{2n}$$

$$\lim_{n \rightarrow \infty} \frac{h(n)}{g(n)} = \lim_{n \rightarrow \infty} \frac{1 + 4e^{2n}}{e^{2n}} = 4$$

$$h(n) = G(g(n))$$

$$\text{ex)} \quad f(n) = e^n \quad g(n) = e^{2n}$$

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = \lim_{n \rightarrow \infty} \frac{e^n}{e^{2n}} = 0$$

$$\therefore f(n) = G(g(n))$$

//

We might be looking at what happens when $\varepsilon \rightarrow 0$

$$\text{say } g(\varepsilon) = G(f(\varepsilon))$$

$$\lim_{\varepsilon \rightarrow 0} \frac{g(\varepsilon)}{f(\varepsilon)} = 0$$

$$\text{e.g. } \lim_{\varepsilon \rightarrow 0} \frac{\sqrt{\varepsilon} + 3\varepsilon^3}{1 + 4\varepsilon^2} = 0$$

//

ACCURACY = FIDELITY

Let \bar{f} be an algorithm that computes f

Definition:

$$\text{Abs err (absolute error)} \equiv \|\bar{f}(x) - f(x)\|$$

$$\text{Rel err (relative error)} \equiv \frac{\|\bar{f}(x) - f(x)\|}{\|f(x)\|}$$

for some x .

$$\text{Relative Accuracy} \equiv \frac{\|\bar{f}(\bar{x}) - f(x)\|}{\|f(x)\|}$$

want this to be $O(\epsilon_{ps})$ //

STABILITY \bar{f} is relatively STABLE if

$$\frac{\|\bar{f}(x) - f(\bar{x})\|}{\|\bar{f}(\bar{x})\|} = O(\epsilon_{ps})$$

$$\text{for some } \bar{x} \quad \frac{\|\bar{x} - x\|}{\|x\|} = O(\epsilon_{ps}) //$$

We note that we compute $\bar{f}(\bar{x}) \neq \bar{f}(x) \neq f(\bar{x})$ generally.

BACKWARD STABILITY

The algorithm \tilde{f} is backward stable if

$$\tilde{f}(x) = f(\tilde{x}) \text{ for some } x \text{ with } \frac{\|\tilde{x} - x\|}{\|x\|} = O(\epsilon_{\text{ps}})$$

Gives the right answer using a nearby input.

Rank: Since all norms are "equivalent" on finite dimensional problems, an algorithm will be stable (or otherwise) regardless of the (unweighted) norm used.

LECTURE 15

Example of backward stability:

Subtraction

$$f = x_1 - x_2$$

$$\tilde{f}(\tilde{x}) = f_l(x_1) \ominus f_l(x_2)$$

$$f_l(x_1) = x_1(1 + \varepsilon_1) \quad |\varepsilon_1| < \epsilon_{\text{ps}}$$

$$f_l(x_2) = x_2(1 + \varepsilon_2) \quad |\varepsilon_2| < \epsilon_{\text{ps}}$$

$$\bar{f}(\bar{x}) = [x_1(1+\varepsilon_1) - x_2(1+\varepsilon_2)](1+\varepsilon_3)$$

$|\varepsilon_3| < \text{eps}$

$$\begin{aligned}\therefore f(x_1) \ominus f(x_2) &= x_1(1+\varepsilon_1)(1+\varepsilon_3) \\ &\quad - x_2(1+\varepsilon_2)(1+\varepsilon_3) \\ &= x_1(1+\varepsilon_1 + \varepsilon_3 + \varepsilon_1\varepsilon_3) \\ &\quad - x_2(1+\varepsilon_2 + \varepsilon_3 + \varepsilon_2\varepsilon_3) \\ &\varepsilon_1\varepsilon_3 \text{ } \& \varepsilon_2\varepsilon_3 = O(\text{eps}^2) \\ \text{we ignore}\end{aligned}$$

$$\therefore f(x_1) \ominus f(x_2) = x_1(1+\varepsilon_4) - x_2(1+\varepsilon_5)$$

$\therefore \bar{f}(x) = f(x_1) \ominus f(x_2)$ is equal to the difference
of \bar{x}_1 and \bar{x}_2 where

$$\frac{|\bar{x}_1 - x_1|}{|x_1|} \leq \text{eps} \quad \frac{|\bar{x}_2 - x_2|}{|x_2|} \leq \text{eps}$$

\therefore Backward stable operation. //

ex) $A = xy^*$ $x \in \mathbb{C}^m$ $y \in \mathbb{C}^n$

$$A = \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \quad \begin{matrix} \uparrow \\ m \\ \downarrow \\ \leftarrow n \rightarrow \end{matrix} \quad \bar{A}(x, y) = f(x) \otimes f(y)^*$$

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

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

$$f(x) \otimes f(y)^*$$

$$= (x + \delta x) \otimes (y + \delta y)^*$$

$$= (x + \delta x) (y + \delta y)^* (I + \delta)$$

$$\text{so let } \bar{x} = x + \delta x \quad \bar{y} = y + \delta y$$

$\bar{x} \bar{y}^* (I + \delta)$ could find neighboring

$$\begin{array}{c} \uparrow \\ \bar{x} \text{ & } \bar{y} \end{array}$$

but generally this δ destroys the rank 1 nature of the exact product.

AXIOMS OF FLOATING POINT ARITHMETIC

$$\begin{array}{l}
 \boxed{\begin{array}{l} \text{(I)} \\ \forall x \in \mathbb{R} \exists \varepsilon \quad |\varepsilon| < \text{eps} \\ \text{st} \quad f_1(x) = x(1+\varepsilon) \end{array}} \\
 \boxed{\begin{array}{l} \text{(II)} \\ \text{For all } x, y \in \mathbb{R} \exists \varepsilon \quad |\varepsilon| < \text{eps} \text{ st} \\ x \otimes y = x * y (1+\varepsilon) \end{array}}
 \end{array}$$

ACCURACY OF BACKWARD STABLE ALGORITHMS

Thm: Suppose BW stable algorithm. Apply it to solve a problem $f: X \rightarrow Y$ w/ condition number k_f on a computer satisfying Floating Point Axioms:

$$\begin{aligned}
 \Rightarrow \text{Rel Err} & \frac{\|\bar{f}(x) - f(x)\|}{\|f(x)\|} \\
 & = O(k_f \text{eps})
 \end{aligned}$$

This means that algorithmic choices that improve the conditioning, by making it smaller, can have a significant impact on accuracy, as well as on stability.