

Condition number \rightarrow good/bad problem

Stability \rightarrow good/bad algorithm

Floating point representation

0.3 \rightarrow $0.01001100110011001 \dots$ truncate!

$= 2^{-2} = 1.00110011001 \dots$ mantissa



64 bits

0



Theorem: For each $x \in \pm [10^{-308}, 10^{308}]$, there is an exactly representable floating point number \tilde{x} such that

$$\frac{|x - \tilde{x}|}{|x|} \leq 2^{-52} \approx 2 \cdot 10^{-16} = \epsilon, \text{ "machine precision"}$$

Approximieren \rightarrow Inaccuracies

e.g. $\frac{1}{3} + \frac{1}{3} + \frac{1}{3} = 0.333\dots3 + 0.333\dots3 + 0.333\dots3 = 0.999\dots9 \neq 1$

$\gg a = 0.2$ ↗ error

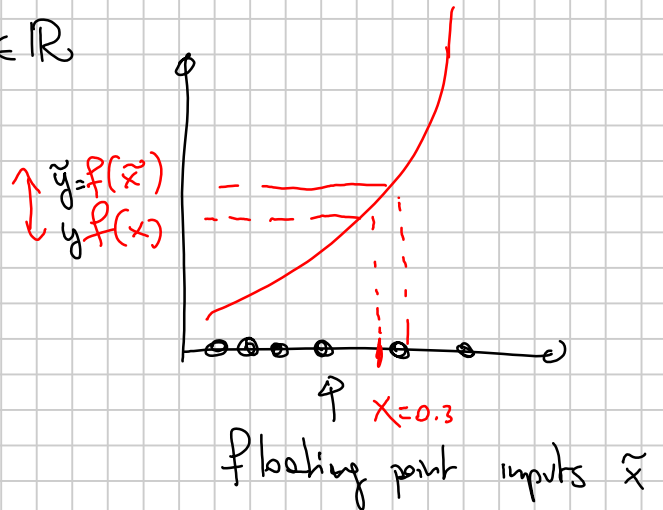
$\gg b = 0.3$ ↗ error

$\gg c = a + b$ ↗ error

$y = f(x)$ computational problem

for most possible input $x \in \mathbb{R}$

The distance $\tilde{y} - y$ depends on $f'(x)$



$$\frac{|\tilde{y} - y|}{|y|} \leq K \cdot \frac{|\tilde{x} - x|}{|x|}$$

\uparrow condition number $K_{f,x}$

The smallest error bound that I can ensure for all inputs x

is $K_{f,x} \cdot u$

\uparrow cond. number \uparrow machine prec

intrinsic error

Vector input:

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \rightarrow \begin{bmatrix} \tilde{x}_1 \\ \tilde{x}_2 \\ \tilde{x}_3 \end{bmatrix}$$

$$\frac{|\tilde{x}_i - x_i|}{|x_i|} \leq u \text{ for each } i$$

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq u \cdot O(u)$$

Usually, one writes $O(u)$ to denote $u \cdot \text{polynomial in the dimension}$.

$f(x: \text{double}) : \text{double}$

Def: an algorithm is stable (on some input x) if the computed

\tilde{y} satisfies

$$\frac{|\tilde{y} - y|}{|y|} \approx K_{f,x} \cdot O(u)$$

Proving stability directly involves keeping track of errors in all operations

$$\begin{aligned} a \oplus b &= \text{floating point appr. } (a+b) \\ &= (a+b)(1+\delta) \quad |\delta| \leq u \end{aligned}$$

Indeed, the floating point approx. \tilde{y} of y satisfies

$$\frac{\tilde{y} - y}{y} = \delta \quad |\delta| \leq u$$

$$\tilde{y} = y(1+\delta)$$

Ex: Keep track of the algorithmic error in computing

$$y = [a_1 \ a_2 \ a_3] \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} = a_1 b_1 + a_2 b_2 + a_3 b_3$$

Assume a_i, b_i already floating point numbers

$$\tilde{y} = (a_1 \oplus b_1 \oplus a_2 \oplus b_2) \oplus a_3 \oplus b_3$$

$$= (a_1 b_1 (1+\delta_1) \oplus a_2 b_2 (1+\delta_2)) \oplus a_3 b_3 (1+\delta_3)$$

$$= \left((a_1 b_1 (1+\delta_1) + a_2 b_2 (1+\delta_2)) (1+\delta_4) + a_3 b_3 (1+\delta_3) \right) (1+\delta_5) \quad |\delta_i| \leq u$$

$$= \underbrace{a_1 b_1 + a_2 b_2 + a_3 b_3}_y + a_1 b_1 (\delta_1 + \delta_4 + \delta_5) + a_2 b_2 (\delta_2 + \delta_4 + \delta_5) + a_3 b_3 (\delta_3 + \delta_5) + O(u^2)$$

$$|\tilde{y} - y| \leq |a_1 b_1| \cdot 3u + |a_2 b_2| \cdot 3u + |a_3 b_3| \cdot 2u$$

$$\leq \exists u \left(|a_1| |b_1| + |a_2| |b_2| + |a_3| |b_3| \right) = \begin{bmatrix} |a_1| & |a_2| & |a_3| \end{bmatrix} \begin{bmatrix} |b_1| \\ |b_2| \\ |b_3| \end{bmatrix}$$

In dimension n ,

$$|\tilde{y} - y| \leq nu \cdot \underbrace{|a^T| \cdot |b|}_{\text{not } y}$$

$|a^T| \cdot |b|$ can be much larger than $|y|$, if there is cancellation, e.g.

$$y = \begin{bmatrix} 1 & -1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 10^6 \\ 10^6 \\ 1 \end{bmatrix} = 1$$

$$|a^T| \cdot |b| = 2 \cdot 10^6 + 1$$

Is this large error due to an unstable algorithm or an ill-conditioned problem?

Backward stability Idea: sometimes, you can see the exactly computed \tilde{y} as the exact result of a computation with perturbed inputs

$$\tilde{y}(a, b) = y(\hat{a}, \hat{b})$$

$$\tilde{y} = a_1 \underbrace{b_1(1 + \delta_1 + \delta_4 + \delta_5)}_{\hat{b}_1} + a_2 \underbrace{b_2(1 + \delta_2 + \delta_4 + \delta_5)}_{\hat{b}_2} + a_3 \underbrace{b_3(1 + \delta_3 + \delta_5)}_{\hat{b}_3} + O(u^2)$$

$$= a_1 \hat{b}_1 + a_2 \hat{b}_2 + a_3 \hat{b}_3$$

$$\left| \frac{\hat{b}_1 - b_1}{b_1} \right| = |\delta_1 + \delta_4 + \delta_5| \leq 3|u| \quad \frac{\|\hat{b} - b\|}{\|b\|} \leq 3u + O(u^2) = O(u)$$

We can use condition number bounds:

$$\frac{|\tilde{y} - y|}{|y|} \leq K_{y,b} \cdot \frac{\|\hat{b} - b\|}{\|b\|} = K_{y,b} \cdot O(u)$$

of the same order of magnitude as the intrinsic error $K_{y,b} \cdot u$

\Rightarrow we can conclude that this algorithm is stable!

Def: An algorithm is called backward stable (on input x)

if the result $\tilde{y}(x)$ it computes in machine arithmetic can be seen as the exact result of a slightly perturbed input, $y(\hat{x})$, $\frac{\|\hat{x} - x\|}{\|x\|} = O(u)$

Backward stability \Rightarrow stability!

$\triangle!$ We can't apply this trick to every algorithm!

Ex: outer product $Y = ab^T = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} [b_1 \ b_2 \ b_3]$

$$= \begin{bmatrix} a_1 b_1 & a_1 b_2 & a_1 b_3 \\ a_2 b_1 & a_2 b_2 & a_2 b_3 \\ a_3 b_1 & a_3 b_2 & a_3 b_3 \end{bmatrix}$$

$$\tilde{Y} = \begin{bmatrix} a_1 \odot b_1 & a_1 \odot b_2 & a_1 \odot b_3 \\ \vdots & \vdots & \vdots \\ \dots & \dots & a_3 \odot b_3 \end{bmatrix} \neq \begin{bmatrix} \hat{a}_1 \\ \hat{a}_2 \\ \hat{a}_3 \end{bmatrix} [\hat{b}_1 \ \hat{b}_2 \ \hat{b}_3],$$

n general, because the errors do not always make that matrix rank 1

$$\text{Ex: } Y = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} [4 \ 5 \ 6] = \begin{bmatrix} 4 & 5 & 6 \\ 8 & 10 & 12 \\ 12 & 15 & 18 \end{bmatrix}$$

$$\tilde{V} = \begin{bmatrix} 3.99 & 4.01 & 5.99 \\ 7.99 & 9.99 & 12.01 \\ 12.01 & 14.99 & 18.02 \end{bmatrix}$$

↑ ↑ ↑

n general, the columns will not be multiples of the same vector
 \Rightarrow we cannot write it as $\tilde{V} = \begin{bmatrix} \hat{a}_1 \\ \hat{a}_2 \\ \hat{a}_3 \end{bmatrix} [b_1 \ b_2 \ b_3]$

"A posteriori" backward stability

Ex: linear system $Ax=y$, $A \in \mathbb{R}^{n \times n}$

Given a computed solution \tilde{x} , let us consider the residual $\tilde{r} = A\tilde{x} - y$

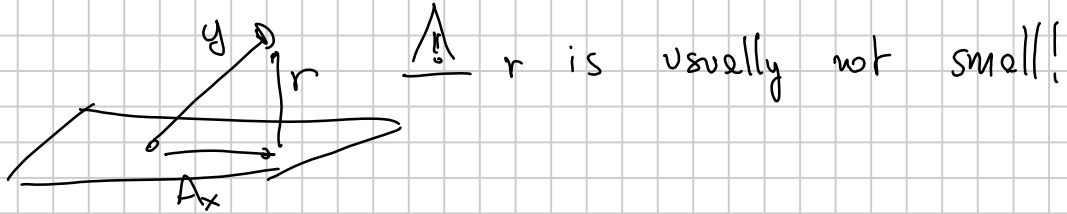
Note that $A\tilde{x} = y + \tilde{r} \Rightarrow \tilde{x}$ is the exact solution of the linear system $A\tilde{x} = \hat{y}$ with $\hat{y} = y + \tilde{r}$

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq K(A) \frac{\|\hat{y} - y\|}{\|y\|} = K(A) \frac{\|\tilde{r}\|}{\|y\|}$$

(we can usually ignore the fact that \tilde{r} and the other quantities are computed in machine arithmetic, their order of magnitude is still reliable).

Can we do it also for least-squares problems?

$$\min_{x \in \mathbb{R}^n} \|Ax - y\|_2^2$$



However, the gradient $g(\tilde{x}) = A^T A \tilde{x} - A^T y$ is zero in the exact solution; we can use the previous α -posteriori bound on the normal equations:

$$\frac{\|\tilde{x} - x\|}{\|x\|} = \underbrace{\kappa(A^T A)}_{\uparrow} \frac{\|A^T A \tilde{x} - A^T y\|}{\|A^T y\|}.$$

$\hat{!}$ $\kappa(A)^2$ larger than the cond. number of our problem, if $\alpha \approx 0$

$$\Rightarrow \text{even if } \frac{\|A^T A \tilde{x} - A^T y\|}{\|A^T y\|} = \mathcal{O}(u),$$

this does not mean that the algorithm is stable.

obs: small residual usually come from algorithms that are backward stable.

ex: $Ax = y, A \in \mathbb{R}^{n \times n}$

$$\hat{A} \hat{x} = \hat{y}$$

$$\hat{A} = A + E, \frac{\|E\|}{\|A\|} = \mathcal{O}(u)$$

$$\hat{y} = y + f, \frac{\|f\|}{\|y\|} = \mathcal{O}(u)$$

$$\begin{aligned} (A+E)\tilde{x} &= y+f \Rightarrow \|A\tilde{x} - y\| = \|f - E\tilde{x}\| \leq \|f\| + \|E\| \cdot \|\tilde{x}\| \\ &= \mathcal{O}(u) (\|y\| + \|A\| \cdot \|\tilde{x}\|) \end{aligned}$$

\Rightarrow the residual $\|A\tilde{x} - y\|$ is always small with respect to

to $\|y\|$ but to $\|y\| + \|A\| \cdot \|\tilde{x}\|$

$$\|y\| = \|Ax\| \leq \|A\| \cdot \|x\|$$

Variant of the residual test for least-squares problems:

$$x \text{ solves } \min_x \|Ax - y\| \quad \text{with } r = Ax - y$$

$$\Leftrightarrow \begin{bmatrix} I_m & A \\ A^T & 0 \end{bmatrix} \begin{bmatrix} -r \\ x \end{bmatrix} = \begin{bmatrix} y \\ 0 \end{bmatrix} \quad \Leftrightarrow \begin{cases} -r + Ax = y \\ -A^T r = 0 \end{cases} \Leftrightarrow A^T(Ax - y) = 0$$

$(m+n) \times (m+n)$

This can be used for a residual test:

given \tilde{x} with residual \tilde{r}

$$\frac{\left\| \begin{bmatrix} -\tilde{r} \\ \tilde{x} \end{bmatrix} - \begin{bmatrix} -r \\ x \end{bmatrix} \right\|}{\left\| \begin{bmatrix} -r \\ x \end{bmatrix} \right\|} \leq \kappa \left(\begin{bmatrix} I_m & A \\ A^T & 0 \end{bmatrix} \right) \frac{\left\| \begin{bmatrix} 0 \\ -A^T \tilde{r} \end{bmatrix} \right\|}{\left\| \begin{bmatrix} y \\ 0 \end{bmatrix} \right\|}$$

$$\Rightarrow \frac{\|\tilde{x} - x\|}{\left\| \begin{bmatrix} -r \\ x \end{bmatrix} \right\|^2} \leq \kappa \left(\begin{bmatrix} \alpha I & A \\ A^T & 0 \end{bmatrix} \right) \cdot \frac{\frac{1}{\alpha} \|A^T \tilde{r}\|}{\|y\|}$$

"A posteriori" stability tests, using the computed result.