

One thing missing last time:

Catastrophic cancellation:

Subtracting two nearly equal numbers cancel the most significant digits but the result can have large relative error

ex 1. evaluate  $\frac{1}{1-x} - 1$  for  $|x| \ll 1$ ,  $x \in \mathbb{F}$

Method 1: Direct evaluation

$$\text{Output}_1 = \left[ \frac{1}{(1-x)(1+\delta_1)} (1+\delta_2) - 1 \right] (1+\delta_3) \quad |\delta_i| \leq \epsilon_{\text{mach}}$$

$$= \frac{[1+\delta_2 - (1-x)(1+\delta_1)] (1+\delta_3)}{(1-x)(1+\delta_1)}$$

$$= \frac{\delta_2 - \delta_1 + x(1+\delta_1)}{1+x} \frac{1+\delta_3}{1+\delta_1}$$

when  $x \sim O(\delta_2 - \delta_1)$ , relative error  $\sim O\left(\frac{\delta_2 - \delta_1}{x}\right) = O(1)$

Method 2: Rearrange calculation

$$\text{from } \frac{1}{1-x} - 1 = \frac{x}{1-x}$$

$$\text{Output}_2 = \frac{x(1+\delta_1)}{(1-x)(1+\delta_2)} (1+\delta_3)$$

relative error  $\sim O(\delta)$  even when  $x \sim O(\delta)$

ex 2.  $\frac{e^x - 1}{x}$ ,  $|x| \ll 1$

assume the exp and log function are both computed with a relative error not exceeding  $\epsilon_{\text{mach}}$

from Taylor expansion

$$\frac{e^x - 1}{x} = \frac{1 + x + \frac{1}{2}x^2 + \dots - 1}{x} \approx 1 + \frac{1}{2}x + O(x^2)$$

Method 1: Direct evaluation

$$\text{Output}_1 = \frac{[e^x(1+\delta_1) - 1](1+\delta_2)}{x(1+\delta_3)} (1+\delta_4)$$

$$= \frac{(1 + x + \frac{1}{2}x^2 + \dots)(1+\delta_1) - 1}{x} \frac{1+\delta_2}{1+\delta_3} (1+\delta_4)$$

$$\approx \left( \frac{\delta_1}{x} + 1 + \frac{1}{2}x \right) \frac{1+\delta_2}{1+\delta_3} (1+\delta_4)$$

↑  
relative error  $\sim O(\frac{\delta_1}{x})$

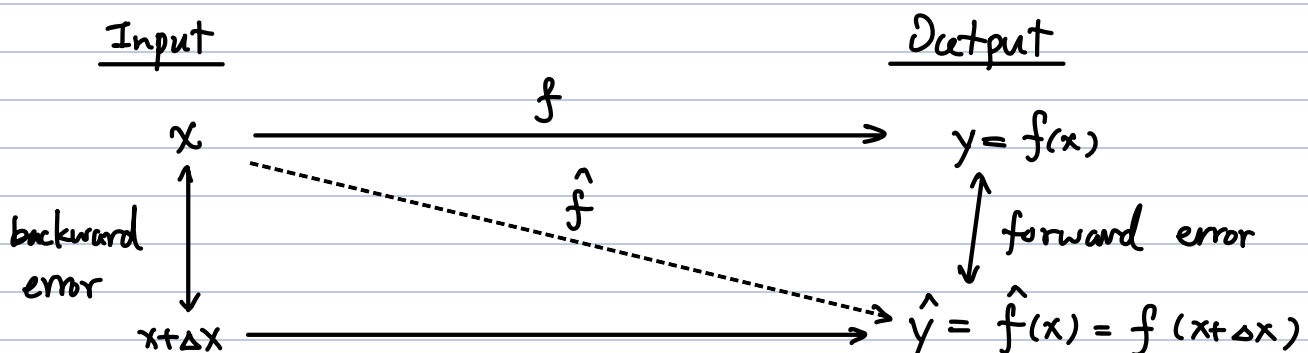
Method 2: Rearrange calculation

$$\text{First compute } \hat{y} = e^x(1+\delta_1)$$

$$\text{then } \text{Output}_2 = \frac{\hat{y} - 1}{\log \hat{y}} (1+\delta_2)$$

exercise: while the relative errors of numerator and denominator are  $O(1)$  for  $x \sim O(\epsilon_{\text{mach}})$ ,  
Output<sub>2</sub> has  $O(\epsilon_{\text{mach}})$  relative error and is accurate.

Last time:  $y = f(x)$ ,  $x \in \mathbb{R}^n$ ,  $y \in \mathbb{R}^m$ , approximated by  $\hat{f}(x)$



$$\text{relative forward error} = \frac{\|\hat{f}(x) - f(x)\|}{\|f(x)\|}$$

$$\text{relative backward error} = \frac{\|\Delta x\|}{\|x\|} \quad \text{s.t. } f(x+\Delta x) = \hat{f}(x)$$

Definition: An algorithm is not unique. can often be chosen to be arbitrary large

- backward stable if  $\exists \Delta x$ , s.t.  $\hat{f}(x) = f(x+\Delta x)$ ,

$$\frac{\|\Delta x\|}{\|x\|} = O(\epsilon_{\text{mach}})$$

- (numerically) stable if  $\exists \Delta x, \Delta y$  s.t.  $\hat{f}(x) + \Delta y = f(x+\Delta x)$

$$\|\Delta y\|/\|y\|, \|\Delta x\|/\|x\| = O(\epsilon_{\text{mach}})$$

- accurate (forward stable) if  $\frac{\|\hat{f}(x) - f(x)\|}{\|f(x)\|}$  is small ( $\sim O(\epsilon_{\text{mach}})$ )

ex 1. Inner product is backward stable

inner product using floating point numbers  $\rightarrow f_l(x^T y) = (x+\Delta x)^T y$  with  $\|\Delta x\| \leq \gamma_n \|x\|$ ,  $\gamma_n = O(n \epsilon_{\text{mach}})$  (\*)

$|x^T y - f_l(x^T y)| = |\Delta x^T y| \leq \gamma_n \|x\|^T \|y\|$

ex 2. Outer product is not backward stable

but satisfies  $f_l(x y^T) = x y^T + E$ ,  $\|E\| \leq \epsilon_{\text{mach}} \|x y^T\|$

hence numerically stable

↑ exercise

Remark: backward stability implies numerical stability.

• Q: When is a backward/numerically stable algorithm accurate?

$$\frac{\|\hat{f}(x) - f(x)\|}{\|f(x)\|} = \frac{\|f(x + \Delta x) + \Delta y - f(x)\|}{\|f(x)\|}$$

forward error

$$\begin{aligned} & \stackrel{\text{triangle ineq}}{\leq} \frac{\|f(x + \Delta x) - f(x)\|}{\|f(x)\|} + \frac{\|\Delta y\|}{\|y\|} \\ & \qquad \qquad \qquad \underbrace{\qquad \qquad \qquad}_{O(\epsilon_{\text{mach}})} \\ & = \frac{\|f(x + \Delta x) - f(x)\| / \|f(x)\|}{\|\Delta x\| / \|x\|} \underbrace{\frac{\|\Delta x\|}{\|x\|}}_{\text{backward error}} + O(\epsilon_{\text{mach}}) \end{aligned}$$

(Relative) condition number

$$K(x) := \sup_{\|\Delta x\| \leq \epsilon_{\text{mach}} \|x\|} \frac{\|f(x + \Delta x) - f(x)\| / \|f(x)\|}{\|\Delta x\| / \|x\|}$$

Thm: For a numerically stable algorithm.

$$\frac{\|\hat{f}(x) - f(x)\|}{\|f(x)\|} = O(K(x) \epsilon_{\text{mach}})$$

Remark: Forward error  $\leq$  condition number  $\times$  backward error

The condition #  $K$  measures the sensitivity of  $f$  to perturbed inputs, which is independent of the algorithm used.

---

Detour: Vector and matrix norm

To quantify errors for vectors/matrices, we use norms

$$\|\cdot\| : \mathbb{C}^n \text{ (or } \mathbb{C}^{m \times n}) \rightarrow \mathbb{R}$$

satisfying

- 1)  $\|x\| \geq 0$ ,  $\|x\| = 0$  iff  $x = 0$
- 2)  $\|\alpha x\| = |\alpha| \|x\|$ ,  $\forall \alpha \in \mathbb{C}$ ,  $x \in \mathbb{C}^n$
- 3)  $\|x+y\| \leq \|x\| + \|y\|$

example:

Vector norm:

$x^*$  conjugate  
+ transpose

- 1)  $\|x\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}$ ,  $1 \leq p < +\infty$ ,  $p$ -norm
- 2)  $\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$   $\infty$ -norm
- 3)  $p=2 \Rightarrow \|x\|_2 = (x^* x)^{1/2}$  Euclidean norm

Matrix norm:

- 1)  $\|A\|_F = \left( \sum_{i,j} |a_{ij}|^2 \right)^{1/2} = [\text{tr}(A^* A)]^{1/2}$  Frobenius norm
- 2)  $\|A\|_\infty = \max_{i,j} |a_{ij}|$  max norm
- 3)  $\|A\|_{\alpha, \beta} = \max_{x \neq 0} \frac{\|Ax\|_\beta}{\|x\|_\alpha}$  subordinate norm

The subordinate matrix norm measures the size of the output relative to the size of the input.

• example of subordinate norm:

$$1) \|x\|_1 = \sum_{i=1}^n |x_i| \quad \text{is } 1\text{-norm}$$

$$Ax = \begin{bmatrix} | & & | \\ a_1 & \dots & a_n \\ | & & | \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = x_1 \begin{bmatrix} | \\ a_1 \\ | \end{bmatrix} + \dots + x_n \begin{bmatrix} | \\ a_n \\ | \end{bmatrix}$$

$$\|Ax\|_1 = \left\| x_1 \begin{bmatrix} | \\ a_1 \\ | \end{bmatrix} + \dots + x_n \begin{bmatrix} | \\ a_n \\ | \end{bmatrix} \right\|_1$$

$$\leq \sum_{i=1}^n |x_i| \|a_i\|_1$$

$$\leq \left[ \max_{1 \leq i \leq n} \|a_i\|_1 \right] \|x\|_1$$

$$= \|A\|_1$$

( $\leq$  holds for  $x = e_i = (0, \dots, 0, \underset{\uparrow}{1}, 0, \dots, 0)^T$   
picks out max  $\|a_i\|_1$  column of  $A$ )

$$2) \|A\|_2 = \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2} = \max_{\|x\|_2=1} \sqrt{x^* A^* A x}$$

$$= \max_{\|x\|_2=1} \sum_{i=1}^n |x^* u_i|^2 \lambda_i$$

$$= \lambda_{\max}(A^* A) \quad (\text{"take } x = u_{i_{\max}})$$

$(u_i, \lambda_i), i=1, \dots, n$   
eigenvector,  
eigenvalue  
of matrix  
 $A^* A$ .  
 $\lambda_i \in \mathbb{R}$

Some properties:

$$1) \|A\|_{\alpha, \beta} = \max_{x \neq 0} \left\| A \frac{x}{\|x\|_\alpha} \right\|_\beta = \max_{\|x\|_\alpha=1} \|Ax\|_\beta$$

2) Any subordinate norm is consistent with the vector norm that induce it :  $\|Ax\|_\beta \leq \|A\|_{\alpha, \beta} \|x\|_\alpha$

Any subordinate norm is submultiplicative :  $\|AB\|_{\alpha, \gamma} \leq \|A\|_{\beta, \gamma} \cdot \|B\|_{\alpha, \beta}$

$$\text{Pf: } \|ABx\|_\gamma \leq \|A\|_{\beta, \gamma} \|Bx\|_\beta \leq \|A\|_{\beta, \gamma} \|B\|_{\alpha, \beta} \|x\|_\alpha$$

Divide both sides by  $\|x\|_\alpha$  and take supreme  $x \neq 0$   $\square$

3) The Frobenius norm is consistent with the Euclidean norm

$$\|Ax\|_2 \leq \|A\|_F \|x\|_2, \text{ and submultiplicative. (exercise)}$$

max norm is not submultiplicative :  $\|AB\|_\infty \leq n \|A\|_\infty \|B\|_\infty$  (exercise)

4) (Equivalence of norms)

For any two vector/matrix norm,  $\|\cdot\|_\alpha, \|\cdot\|_\beta$ ,

$$\text{we have } r \|A\|_\alpha \leq \|A\|_\beta \leq s \|A\|_\alpha$$

for some  $r, s > 0$ , for all  $A \in \mathbb{C}^{m \times n}$   
( $r, s$  only depend on how the norm  $\|\cdot\|_\alpha, \|\cdot\|_\beta$  are defined and the dimension  $m, n$ )

$$\text{ex. } \frac{1}{\sqrt{n}} \|x\|_2 \leq \|x\|_1 \leq \sqrt{n} \|x\|_2, \quad \frac{1}{\sqrt{n}} \|A\|_2 \leq \|A\|_1 \leq \sqrt{n} \|A\|_2$$

Now we are ready to handle condition #'s

If  $f(x) = (f_1(x), \dots, f_m(x)) : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is differentiable then,

$$f_j(x + \Delta x) = f_j(x) + \sum_{i=1}^n \frac{\partial f_j}{\partial x_i}(x) \Delta x_i + O(\|\Delta x\|^2)$$

Jacobian  $Df(x) = \left( \frac{\partial f_j}{\partial x_i}(x) \right)_{\substack{1 \leq i \leq n \\ 1 \leq j \leq m}}$

Then  $f(x + \Delta x) = f(x) + Df(x) \Delta x + O(\|\Delta x\|^2)$

Recall the definition.

$$K(x) := \sup_{\frac{\|\Delta x\|}{\|x\|} \leq \epsilon_{\text{mach}}} \frac{\|f(x + \Delta x) - f(x)\| / \|\Delta x\|}{\|f(x)\|}$$

$$= \sup_{\Delta x} \frac{\|Df(x) \Delta x + O(\|\Delta x\|^2)\|}{\|\Delta x\|} \frac{\|x\|}{\|f(x)\|}$$

consistency  
of matrix norm

$$\leq \frac{\|Df(x)\| \|x\|}{\|f(x)\|} + O(\epsilon_{\text{mach}} \|x\|^2 / \|f(x)\|)$$

condition number  
for differentiable system

usually negligible  
or comparable to the previous  
term.

example: Summation function

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

(a special case of inner product  
with  $y = \mathbf{1}$ . hence backward  
stable)

$$Df(x) = [1, \dots, 1]$$

Take  $\|\cdot\|_1$  in the following

$$\|Df(x)\|_1 = 1$$

$$K(x) = \frac{\|Df(x)\|_1 \|x\|_1}{|f(x)|} = \frac{\sum_{i=1}^n |x_i|}{\left| \sum_{i=1}^n x_i \right|}$$

The forward error

$$\frac{|\hat{f}(x) - f(x)|}{|f(x)|} = O\left(\frac{\sum_{i=1}^n |x_i|}{\left| \sum_{i=1}^n x_i \right|} \epsilon_{mach}\right)$$

Remarks:

1) Estimating the backward error  $\frac{\|\Delta x\|}{\|x\|}$  is called backward error analysis. Combining backward error (of an algorithm) and condition # yields forward error (of a problem).

2) Forward error bound can also be obtained directly here by using the error bound (\*) (on pp.3).