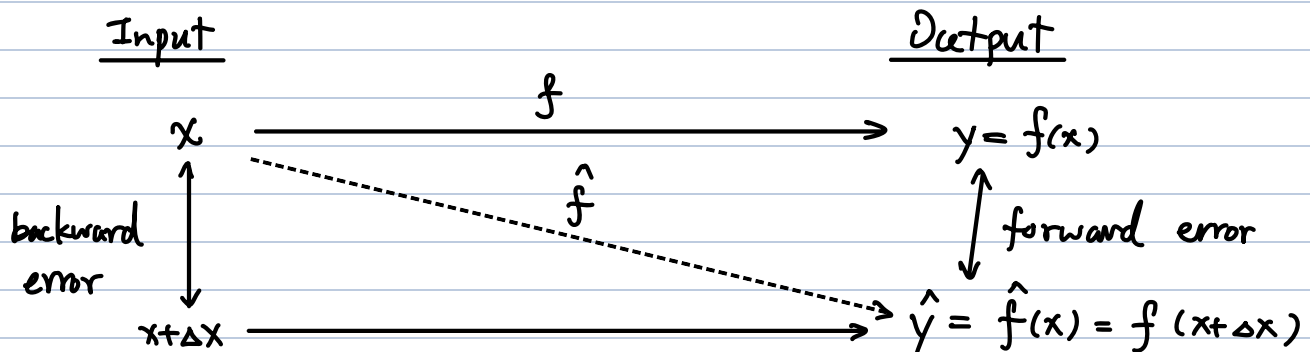


$y = f(x)$, $x \in \mathbb{R}^n$, $y \in \mathbb{R}^m$ approximated by $\hat{f}(x)$

- How should we measure the "quality" of \hat{y} ?



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

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

- Why "backward" analysis?
 - Uncertainty in data \Leftrightarrow rounding error in algorithm
 - Reduce bounding forward error to perturbation theory which is well understood.

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

ex 1. Inner product is backward stable

inner product using floating point numbers

$$f_l(x^T y) = (x + \Delta x)^T y \quad \text{with } \|\Delta x\| \leq \gamma_n \|x\|, \quad \gamma_n = O(n \epsilon_{\text{mach}})$$

$\|x\| = (\sum_{i=1}^n |x_i|^2)^{1/2}$

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

ex 2. Outer product is not backward stable

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

hence numerically stable ↖ exercise

Remark: backward stability implies numerical stability.

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

$$\underbrace{\frac{\|\hat{f}(x) - f(x)\|}{\|f(x)\|}}_{\text{forward error}} = \frac{\|f(x + \Delta x) + \Delta y - f(x)\|}{\|f(x)\|}$$

forward error

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

(Relative) condition number

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

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) $\|\alpha x\| = |\alpha| \|x\|$, $\forall \alpha \in \mathbb{C}$, $x \in \mathbb{C}^n \text{ (or } \mathbb{C}^{m \times n})$
- 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) $\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$ ∞ -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

$$3) \|A\|_\alpha = \max_{x \neq 0} \frac{\|Ax\|_\alpha}{\|x\|_\alpha} \quad (\text{take } m=n \text{ for simplicity})$$

subordinate norm

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

• example of subordinate norm:

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

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

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

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

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

$$= \|A\|_1$$

('=' holds for $x = e_i = (0, \dots, 0, 1, 0, \dots, 0)^T$
picks out max $\| \cdot \|_1$ 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:

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

2) Any subordinate norm is consistent with the

$$\text{vector norm that induce it : } \|Ax\|_\alpha \leq \|A\|_\alpha \|x\|_\alpha$$

Any subordinate norm is submultiplicative: $\|AB\|_\alpha \leq \|A\|_\alpha \cdot \|B\|_\alpha$

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

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

3) The Frobenius norm is consistent with the Euclidean norm

$$\|Ax\|_2 \leq \|A\|_F \|x\|_2, \text{ and submultiplicative: } \|AB\|_F \leq \|A\|_F \|B\|_F.$$

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,

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

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

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

Recall the definition.

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

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

consistency
of matrix norm

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

condition number
for differentiable system

usually negligible
or comparable to the previous
term.

example 1: Summation function

$$f(x) = \sum_{i=1}^n x_i \quad (\text{a special case of inner product with } y = \mathbb{1}, \text{ hence backward stable})$$

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

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

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

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

The forward error

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

Remarks:

1) Estimating the backward error $\frac{\|\Delta x\|}{\|x\|}$ is call
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. 2).

example 2: Solving linear equations

$$f(b) = A^{-1}b$$

$$K = \frac{\|b\| \|A^{-1}\|}{\|A^{-1}b\|} \leq \|A\| \|A^{-1}\|$$