

Floating-point arithmetic

Higham 2002 (A. S. N. A)
2nd ed.

- ▶ Computers store number in base 2 with finite/fixed memory (bits)
- ▶ Irrational numbers are stored inexactly, e.g. $1/3 \approx 0.333\dots$
- ▶ Calculations are rounded to nearest floating-point number (rounding error)
- ▶ Thus the accuracy of the final error is nontrivial

(base 2)
 $\sqrt{3} = 1.011000\dots$
 $\times 2^{\lfloor \sqrt{3} \rfloor}$
 $-1000 \quad 1000$
 $10^{-300} \sim 10^{300}$

1. $\left(\frac{1}{3} \cdot 3\right) = 0.999\dots$

Two examples with MATLAB

- ▶ $((\text{sqrt}(2))^2 - 2) * 1e15 = 0.4441$ (should be 1..)
- ▶ $\sum_{n=1}^{\infty} \frac{1}{n} \approx 30$ (should be ∞ ..)

→ rounded off.

$+ 10^{16} \sum_{n=10^{16}}^{\infty} \frac{1}{n} = 30$

An important (but not main) part of numerical analysis/NLA is to study the effect of rounding errors

Canonical reference: Higham's book (2002)

Conditioning and stability



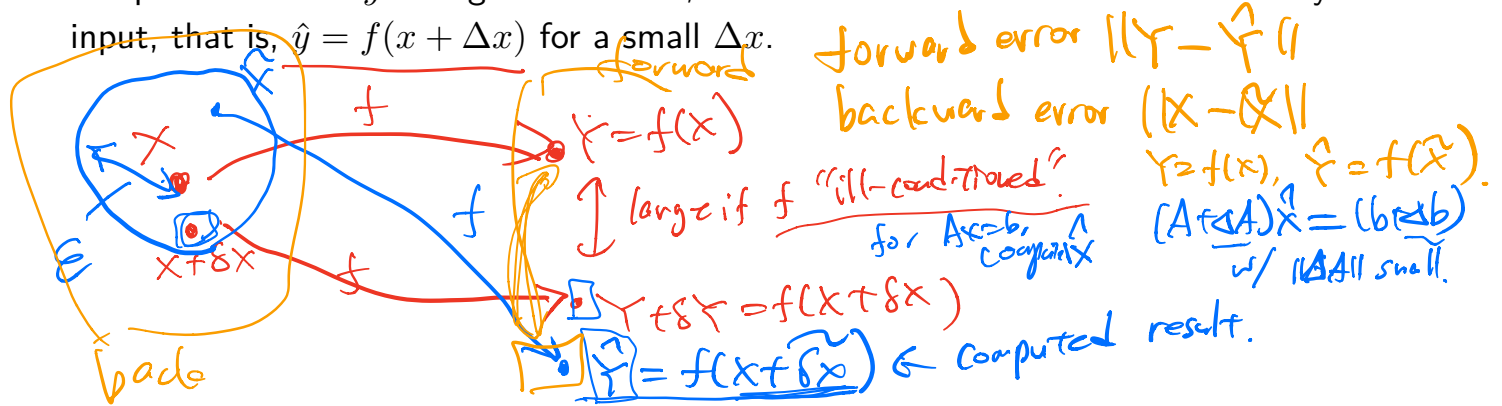
- ▶ **Conditioning** is the sensitivity of a problem (e.g. of finding $y = f(x)$ given x) to perturbation in inputs, i.e., how large $\kappa := \sup_{\delta x} \frac{\|f(x + \delta x) - f(x)\|}{\|\delta x\|}$ is in the limit $\delta x \rightarrow 0$.

(this is *absolute* condition number; equally important is *relative* condition number $\Rightarrow \kappa \gg 1$!)

$$\kappa_r := \sup_{\delta x} \frac{\|f(x + \delta x) - f(x)\|}{\|f(x)\|} \bigg/ \frac{\|\delta x\|}{\|x\|}$$

ill-conditioned $\Leftrightarrow \kappa (\gg 1)$ big
 well-cond. κ small $(\Leftrightarrow 0 < \kappa < 1)$

- ▶ **(Backward) Stability** is a property of an algorithm, which describes if the computed solution \hat{y} is a 'good' solution, in that it is an exact solution of a nearby input, that is, $\hat{y} = f(x + \Delta x)$ for a small Δx .



Conditioning and stability

- ▶ **Conditioning** is the sensitivity of a problem (e.g. of finding $y = f(x)$ given x) to perturbation in inputs, i.e., how large $\kappa := \sup_{\delta x} \|f(x + \delta x) - f(x)\| / \|\delta x\|$ is in the limit $\delta x \rightarrow 0$.

(this is *absolute* condition number; equally important is *relative* condition number

$$\kappa_r := \sup_{\delta x} \frac{\|f(x+\delta x) - f(x)\|}{\|f(x)\|} / \frac{\|\delta x\|}{\|x\|}$$

$A + \Delta A \quad x = b$
 $Ax = b$
 not exact anyway!

- ▶ **(Backward) Stability** is a property of an algorithm, which describes if the computed solution \hat{y} is a 'good' solution, in that it is an exact solution of a nearby input, that is, $\hat{y} = f(x + \Delta x)$ for a small Δx .

$$\|x\|_{\infty}^p \leq \|x\|_2^p \leq \|x\|_{\infty}^p \sqrt{n}$$

If problem is **ill-conditioned** $\kappa \gg 1$, then blame the problem not the algorithm

Notation/convention: \hat{x} denotes a computed approximation to x (e.g. of $x = A^{-1}b$)

ϵ denotes a small term $O(u)$, on the order of unit roundoff/working precision; so we write e.g. $u, 10u, (m+n)u, mnu$ all as ϵ

$$K \begin{pmatrix} 10^{10} & 1 \\ & 10^{10} \end{pmatrix} > 10^{10}$$

- ▶ Consequently (in this lecture/discussion) norm choice does not matter

Numerical stability: backward stability

For computational task $Y = f(X)$ and computed approximant \hat{Y} ,

▶ Ideally, error $\|Y - \hat{Y}\|/\|Y\| = \epsilon$: seldom true

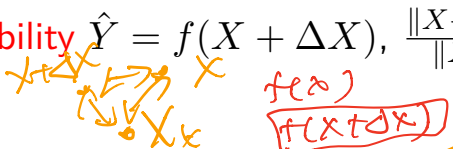
(u : unit roundoff, $\approx 10^{-16}$ in standard double precision)

▶ Good alg. has **Backward stability** $\hat{Y} = f(X + \Delta X)$, $\frac{\|X - \hat{X}\|}{\|X\|} = \epsilon$ “exact solution of slightly wrong input”

Numerical stability: backward stability

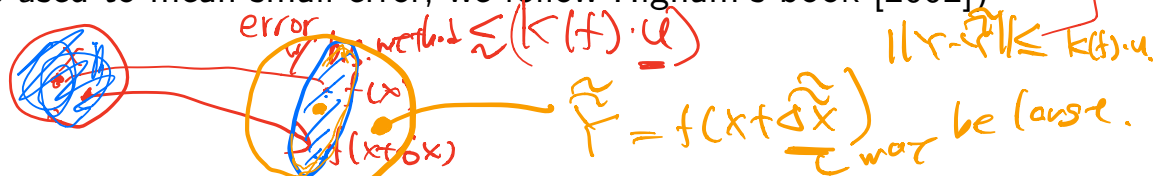
For computational task $Y = f(X)$ and computed approximant \hat{Y} ,

- ▶ Ideally, error $\|Y - \hat{Y}\|/\|Y\| = \epsilon$: seldom true *“forward stability”*
 (u : unit roundoff, $\approx 10^{-16}$ in standard double precision)
- ▶ Good alg. has **Backward stability** $\hat{Y} = f(X + \Delta X)$, $\frac{\|X - \hat{X}\|}{\|X\|} = \epsilon$ “exact solution of slightly wrong input”



- ▶ Justification: Input (matrix) is usually inexact anyway! $f(X + \Delta X)$ is just as good at $f(X)$ at approximating $f(X_*)$ where $\|\Delta X\| = O(\|X - X_*\|)$
 We shall 'settle with' such solution, though it may not mean $\hat{Y} - Y$ is small

- ▶ **Forward stability** $\|Y - \hat{Y}\|/\|Y\| = O(\kappa(f)u)$ “error is as small as backward stable alg.” (sometimes used to mean small error; we follow Higham’s book [2002])



Matrix condition number

$$\kappa_2(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)} (\geq 1)$$

e.g. for linear systems. A backward stable soln for $Ax = b$, s.t. $(A + \Delta A)\hat{x} = b$ satisfies, assuming backward stability $\|\Delta A\| \leq \epsilon \|A\|$ and $\kappa_2(A) \ll \epsilon^{-1}$ (so $\|A^{-1}\Delta A\| \ll 1$),

$$\frac{\|\hat{x} - x\|}{\|x\|} \lesssim \epsilon \kappa_2(A)$$

in practice $\kappa_2(A) < 10^k$
2-norm cond num. κ_2

Matrix condition number

$$\kappa_2(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)} (\geq 1)$$

e.g. for linear systems. A backward stable soln for $Ax = b$, s.t. $(A + \Delta A)\hat{x} = b$ satisfies, assuming backward stability $\|\Delta A\| \leq \epsilon \|A\|$ and $\kappa_2(A) \ll \epsilon^{-1}$ (so $\|A^{-1}\Delta A\| \ll 1$),

$$\frac{\|\hat{x} - x\|}{\|x\|} \lesssim \epsilon \kappa_2(A)$$

$$\hat{x} = \frac{(A + \Delta A)^{-1} b}{(b \neq 0)}$$

wird durch $\|x\| < 1$

'proof': By Neumann series

$$(I - X)^{-1} = 1 + X + X^2 + \dots = \sum_{i=0}^{\infty} X^i \quad (\|X\| < 1)$$

(any norm.)

$$(A + \Delta A)^{-1} = (A(I + A^{-1}\Delta A))^{-1} = (I - A^{-1}\Delta A + O(\|A^{-1}\Delta A\|^2))A^{-1}$$

So $\hat{x} = (A + \Delta A)^{-1}b = A^{-1}b - A^{-1}\Delta A A^{-1}b + O(\|A^{-1}\Delta A\|^2) =$
 $x - A^{-1}\Delta Ax + O(\|A^{-1}\Delta A\|^2)$, Hence

$$\|x - \hat{x}\| \lesssim \|A^{-1}\Delta Ax\| \leq \|A^{-1}\| \|\Delta A\| \|x\| = \kappa_2(A) \|x\|$$

ρ

Backward stable + well conditioned = accurate solution

Suppose

▶ $Y = f(X)$ computed backward stably i.e., $\hat{Y} = f(X + \Delta X)$, $\|\Delta X\| = \epsilon$.

▶ Conditioning $\|f(X) - f(X + \Delta X)\| \leq \kappa \|\Delta X\|$ $(+ \|\epsilon X^2\|)$

Then (relative version possible)

$$\|\hat{Y} - Y\| \leq \kappa \epsilon$$

Backward stable+well conditioned=accurate solution

Suppose

- ▶ $Y = f(X)$ computed backward stably i.e., $\hat{Y} = f(X + \Delta X)$, $\|\Delta X\| = \epsilon$.
- ▶ Conditioning $\|f(X) - f(X + \Delta X)\| \lesssim \kappa \|\Delta X\|$

Then (relative version possible)

$$\|\hat{Y} - Y\| \lesssim \kappa \epsilon$$

'proof':

$$\|\hat{Y} - Y\| = \|f(X + \Delta X) - f(X)\| \lesssim \kappa \underbrace{\|\Delta X\|}_{\epsilon} \cancel{\|f(X)\|} = \kappa \epsilon$$

Backward stable + well conditioned = accurate solution

Suppose

- ▶ $Y = f(X)$ computed backward stably i.e., $\hat{Y} = f(X + \Delta X)$, $\|\Delta X\| = \epsilon$.
- ▶ Conditioning $\|f(X) - f(X + \Delta X)\| \lesssim \kappa \|\Delta X\|$

Then (relative version possible)

$$\|\hat{Y} - Y\| \lesssim \kappa \epsilon$$

'proof':

$$\|\hat{Y} - Y\| = \|f(X + \Delta X) - f(X)\| \lesssim \kappa \|\Delta X\| \|f(X)\| = \kappa \epsilon$$

If well-conditioned $\kappa = O(1)$, good accuracy! Important examples:

- ▶ Well-conditioned linear system $Ax = b$, $\kappa_2(A) \approx 1$ *incl. $Qx = b$*
- ▶ Eigenvalues of symmetric matrices (via Weyl's bound $\lambda_i(A + E) \in \lambda_i(A) + [-\|E\|_2, \|E\|_2]$) *orth. $a_i > b.s. \text{ alg}$*
- ▶ Singular values of any matrix $\sigma_i(A + E) \in \sigma_i(A) + [-\|E\|_2, \|E\|_2]$

Note: eigvecs/singvecs can be highly ill-conditioned

$$\kappa(A) = \frac{\sigma_{\max}}{\sigma_{\min}}$$

orth. $Qx = b$

$$\begin{aligned} |\hat{\lambda}_i(A) - \lambda_i(A)| &\leq \|E\|_2 \\ |\hat{\sigma}_i(A) - \sigma_i(A)| &\leq \|E\|_2 \end{aligned}$$

Backward stability of triangular systems

Recall $Ax = b$ via $Ly = b$, $Ux = y$ (triangular systems).

The computed solution \hat{x} for a (upper/lower) triangular linear system $Rx = b$ solved via back/forward substitution is backward stable, i.e., it satisfies

$$(R + \Delta R)\hat{x} = b, \quad \|\Delta R\| = O(\epsilon \|R\|).$$

Proof: Trefethen-Bau or Higham (nonexaminable but interesting)

- ▶ (backward error can be bounded componentwise)
- ▶ this means $\|\hat{x} - x\|/\|x\| \leq \epsilon \kappa_2(R)$
 - ▶ (unavoidably) poor worst-case (and attainable) bound when ill-conditioned
 - ▶ often better with triangular systems

pessimistic but attainable

(In)stability of $Ax = b$ via LU with pivots

Fact (proof nonexaminable): Computed $\hat{L}\hat{U}$ satisfies $\frac{\|\hat{L}\hat{U} - A\|}{\|L\|\|U\|} = \epsilon$

(note: not $\frac{\|\hat{L}\hat{U} - A\|}{\|A\|} = \epsilon$)

- ▶ If $\|L\|\|U\| = O(\|A\|)$, then $(L + \Delta L)(U + \Delta U)\hat{x} = b$
 $\Rightarrow \hat{x}$ backward stable solution (exercise)

$$\frac{\|\hat{L}\hat{U} - A\|}{\|L\|\|U\|} = \epsilon$$

b.s. if

$$L\hat{U} = A + \Delta A$$

$$\|\Delta A\| \leq \epsilon \|A\|$$

$$\|\Delta A\| \leq \epsilon \|L\| \|U\|$$

possible

$$\|L\| \|U\| \gg \|A\|$$

$$A \sim L^T U$$

$$\|A\| \leq \|L\| \|U\|$$

may or
 may not
 \Rightarrow b.s. of LU.

$$(A + \Delta A)\hat{x} = b$$

(In)stability of $Ax = b$ via LU with pivots

Fact (proof nonexaminable): Computed $\hat{L}\hat{U}$ satisfies $\frac{\|\hat{L}\hat{U} - A\|}{\|\hat{L}\|\|\hat{U}\|} = \epsilon$

(note: not $\frac{\|\hat{L}\hat{U} - A\|}{\|A\|} = \epsilon$)

- ▶ If $\|\hat{L}\|\|\hat{U}\| = O(\|A\|)$, then $(L + \Delta L)(U + \Delta U)\hat{x} = b$
 $\Rightarrow \hat{x}$ backward stable solution (exercise)

$$\|\hat{L}\|\|\hat{U}\| = O(\|A\|)$$

Question: Does $\underline{LU = A + \Delta A}$ or $\underline{LU = PA + \Delta A}$ with $\|\Delta A\| = \epsilon\|A\|$ hold?

Without pivot ($P = I$): $\|\hat{L}\|\|\hat{U}\| \gg \|A\|$ unboundedly (e.g. $\begin{bmatrix} \epsilon & 1 \\ 1 & 1 \end{bmatrix}$) unstable

(In)stability of $Ax = b$ via LU with pivots

Fact (proof nonexaminable): Computed $\hat{L}\hat{U}$ satisfies $\frac{\|\hat{L}\hat{U} - A\|}{\|L\|\|U\|} = \epsilon$

(note: not $\frac{\|\hat{L}\hat{U} - A\|}{\|A\|} = \epsilon$)

- ▶ If $\|L\|\|U\| = O(\|A\|)$, then $(L + \Delta L)(U + \Delta U)\hat{x} = b$
 $\Rightarrow \hat{x}$ backward stable solution (exercise)

Question: Does $LU = A + \Delta A$ or $LU = PA + \Delta A$ with $\|\Delta A\| = \epsilon\|A\|$ hold?

Without pivot ($P = I$): $\|L\|\|U\| \gg \|A\|$ unboundedly (e.g. $\begin{bmatrix} \epsilon & 1 \\ 1 & 1 \end{bmatrix}$) unstable

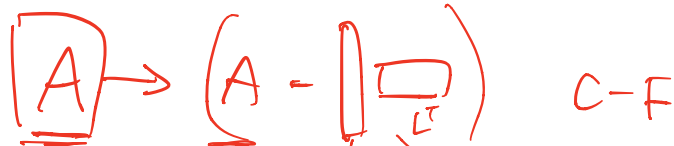
With pivots:

- $\frac{\|L\|\|U\|}{\|A\|} \sim 2^n$
- ▶ Worst-case: $\|L\|\|U\| \gg \|A\|$ grows exponentially with n , unstable
 - ▶ growth governed by that of $\frac{\|L\|\|U\|}{\|A\|} \Rightarrow \frac{\|U\|}{\|A\|}$
 - ▶ **In practice (average case): perfectly stable**
 - ▶ Hence this is how $Ax = b$ is solved, despite alternatives with guaranteed stability exist (but slower; e.g. via SVD, or QR (next))

Resolution/explanation: among biggest open problems in numerical linear algebra!

Stability of Cholesky for $A \succ 0$

Cholesky $A = R^T R$ for $A \succ 0$



▶ succeeds without pivot (active matrix is always positive definite)

▶ R never contains entries $> \|A\|_2$

⇒ backward stable! Hence positive definite linear system $Ax = b$ stable via Cholesky

$$\hat{R}^T \hat{R} = A + \epsilon \|A\|_2$$

$\epsilon = \mathcal{O}(u)$

(caveat;

$$\text{if } \kappa_2(A) \approx \frac{1}{u}$$

then Cholesky may break down;



fails to be positive due to num. errors.

(In)stability of Gram-Schmidt

$$A = QR$$

▶ Gram-Schmidt is subtle

▶ plain (classical) version: $\|Q^T Q - I\| \leq \epsilon \kappa_2(A)^2$ 10^{-5}

▶ modified Gram-Schmidt (orthogonalise 'one vector at a time'): $\|\hat{Q}^T \hat{Q} - I\| \leq \epsilon \kappa_2(A)$ 1 row

▶ Gram-Schmidt twice (G-S again on computed \hat{Q}): $\|\hat{Q}^T \hat{Q} - I\| \leq \epsilon$

$$GS(\hat{Q}) = \begin{bmatrix} Q \\ R_2 \end{bmatrix}$$

Stability of Householder QR

With Householder QR, the computed \hat{Q}, \hat{R} satisfy

$$\|\hat{Q}^T \hat{Q} - I\| = O(\epsilon), \quad \|A - \hat{Q} \hat{R}\| = O(\epsilon \|A\|),$$

and (of course) \hat{R} upper triangular.

best possible!

Rough proof

- ▶ Each reflector satisfies $fl(H_i A) = H_i A + \epsilon_i \|A\|$
- ▶ Hence $(\hat{R} =) fl(H_n \cdots H_1 A) = H_n \cdots H_1 A + \epsilon \|A\|$
- ▶ $fl(H_n \cdots H_1) =: \hat{Q}^T = H_n \cdots H_1 + \epsilon \|A\|$
- ▶ Thus $\hat{Q} \hat{R} = A + \epsilon \|A\|$

$\sum \epsilon_i$
15-kt

Stability of Householder QR

With Householder QR, the computed \hat{Q}, \hat{R} satisfy

$$\hat{Q}^T \hat{Q} - I = O(\epsilon), \quad \|A - \hat{Q}\hat{R}\| = O(\epsilon\|A\|),$$

and (of course) R upper triangular.

Rough proof

- ▶ Each reflector satisfies $fl(H_i A) = H_i A + \epsilon_i \|A\|$
- ▶ Hence $(\hat{R} =) fl(H_n \cdots H_1 A) = H_n \cdots H_1 A + \epsilon \|A\|$
- ▶ $fl(H_n \cdots H_1) =: \hat{Q}^T = H_n \cdots H_1 + \epsilon \|A\|$,
- ▶ Thus $\hat{Q}\hat{R} = A + \epsilon \|A\|$

$$A = QR$$

Notes:

- ▶ This doesn't mean $\|\hat{Q} - Q\|, \|\hat{R} - R\|$ are small at all! Indeed \hat{Q}, R are as ill-conditioned as A $\approx O(\kappa_2(A)\epsilon)$

▶ $Ax = b$ via QR, least-squares stable
(please see arXiv:2009.11392 for application in low-rank approximation)

Orthogonality matters for stability

With orthogonal matrices Q ,

$$\frac{\|fl(QA) - QA\|}{\|QA\|} \leq \epsilon, \quad \frac{\|fl(AQ) - AQ\|}{\|AQ\|} \leq \epsilon$$

Orthogonality matters for stability

With orthogonal matrices Q ,

$$\frac{\|fl(QA) - QA\|}{\|QA\|} \leq \epsilon, \quad \frac{\|fl(AQ) - AQ\|}{\|AQ\|} \leq \epsilon$$

whereas in general, $\|fl(AB) - AB\| \leq \epsilon \|A\| \|B\|$, so

$$\|fl(AB) - AB\| / \|AB\| \leq \epsilon \min(\kappa_2(A), \kappa_2(B)) \quad (\text{why? exercise})$$

Orthogonality matters for stability

With orthogonal matrices Q ,

$$\frac{\|fl(QA) - QA\|}{\|QA\|} \leq \epsilon, \quad \frac{\|fl(AQ) - AQ\|}{\|AQ\|} \leq \epsilon$$

whereas in general, $\|fl(AB) - AB\| \leq \epsilon \|A\| \|B\|$, so

$$\|fl(AB) - AB\| / \|AB\| \leq \epsilon \min(\kappa_2(A), \kappa_2(B)) \text{ (why? exercise)}$$

Hence algorithms involving ill-conditioned matrices are **unstable** (e.g. eigenvalue decomposition of non-normal matrices, Jordan form, etc), whereas those based on orthogonal matrices are **stable**, e.g.

- ▶ Householder QR factorisation
- ▶ **QR algorithm** for $Ax = \lambda x$
- ▶ **Golub-Kahan** algorithm for $A = U\Sigma V^T$
- ▶ **QZ algorithm** for $Ax = \lambda Bx$

We next turn to the algorithms in boldface