

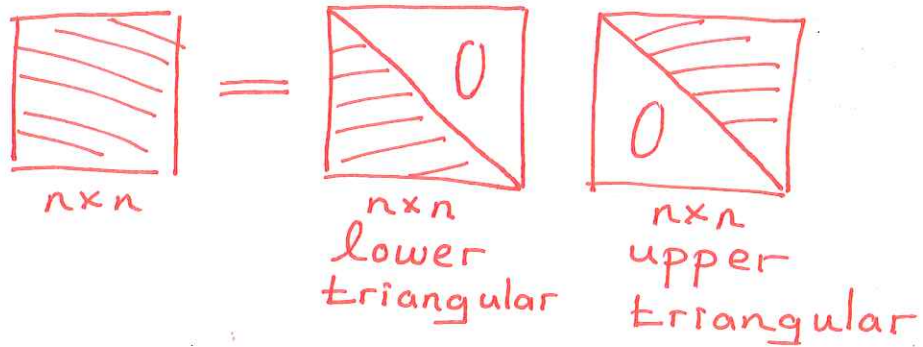


LU Factorization

Let

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(n)}] \in \mathbb{C}^{n \times n}$$

$$A = LU$$



Ex

$$\underbrace{\begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 2 \\ 3 & 2 & 1 \end{bmatrix}}_A \xrightarrow{\substack{\Gamma_2 := \Gamma_2 - 2\Gamma_1 \\ \Gamma_3 := \Gamma_3 - 3\Gamma_1}} \begin{bmatrix} 1 & 2 & 3 \\ 0 & -3 & -4 \\ 0 & -4 & -8 \end{bmatrix}$$

$$\xrightarrow{\Gamma_3 := \Gamma_3 - \frac{4}{3}\Gamma_2} \begin{bmatrix} 1 & 2 & 3 \\ 0 & -3 & -4 \\ 0 & 0 & -8/3 \end{bmatrix}$$

$$\underbrace{\begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 2 \\ 3 & 2 & 1 \end{bmatrix}}_A = \underbrace{\begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 3 & 4/3 & 1 \end{bmatrix}}_L \underbrace{\begin{bmatrix} 1 & 2 & 3 \\ 0 & -3 & -4 \\ 0 & 0 & -8/3 \end{bmatrix}}_U$$

Triangularization by lower triangular matrices

$$\underbrace{\begin{bmatrix} X & X & \dots & X \\ X & X & & X \\ X & X & & X \\ \vdots & \vdots & & \vdots \\ X & X & & X \end{bmatrix}}_A \xrightarrow{\text{lower triangular matrix from left}} \begin{bmatrix} X & X & X \\ 0 & \boxed{X} & X \\ 0 & X & X \\ \vdots & \vdots & \vdots \\ 0 & X & X \end{bmatrix} = A^{(2)}$$

$L^{(1)} A$

$$\xrightarrow{\quad} \begin{bmatrix} X & X & X \\ 0 & X & X \\ 0 & 0 & X \\ \vdots & \vdots & \vdots \\ 0 & 0 & X \end{bmatrix} = A^{(3)}$$

$L^{(2)} L^{(1)} A$

$$L^{(1)} = \begin{bmatrix} 1 & & & & 0 \\ -l_{21} & 1 & & & \\ -l_{31} & 0 & \ddots & & \\ \vdots & \vdots & & \ddots & \\ -l_{n1} & 0 & & & 1 \end{bmatrix},$$

$$\begin{aligned}
 l_{21} &= a_{21} / a_{11} \\
 l_{31} &= a_{31} / a_{11} \\
 &\vdots \\
 l_{n1} &= a_{n1} / a_{11}
 \end{aligned}$$

$$L^{(2)} = \begin{bmatrix} 1 & & & & 0 \\ 0 & 1 & & & \\ 0 & -l_{32} & 1 & & \\ \vdots & \vdots & & \ddots & \\ 0 & -l_{n2} & & & 1 \end{bmatrix},$$

$$\begin{aligned}
 l_{32} &= a_{32}^{(2)} / a_{22}^{(2)} \\
 &\vdots \\
 l_{n2} &= a_{n2}^{(2)} / a_{22}^{(2)}
 \end{aligned}$$

kth step

$$\underbrace{\begin{bmatrix} \boxed{\begin{matrix} (k-1) \times (k-1) \\ X & \dots & X \\ 0 & \dots & X \end{matrix}} & X & \dots & X \\ & X & \dots & X \\ 0 & \boxed{\begin{matrix} X & X \\ X & X \\ \vdots & \vdots \\ X & X \end{matrix}} \end{bmatrix}}_{A^{(k)}} \longrightarrow \underbrace{\begin{bmatrix} \boxed{\begin{matrix} (k-1) \times (k-1) \\ X & \dots & X \\ 0 & \dots & X \end{matrix}} & X & \dots & X \\ & X & \dots & X \\ & X & & X \\ & 0 & & X \\ & \vdots & & \vdots \\ & 0 & & X \end{bmatrix}}_{L^{(k)} A^{(k)}}$$

$$\left(L^{(k-1)} \dots L^{(1)} A \right)$$

$$L^{(k)} = \begin{bmatrix} 1 & & & \\ & \dots & & \\ & & 1 & \\ & & -l_{(k+1)k} & \dots \\ & & \vdots & \\ & & -l_{nk} & \\ & & & 1 \end{bmatrix}$$

$$\begin{aligned}
 l_{(k+1)k} &= a_{(k+1)k}^{(k)} / a_{kk}^{(k)} \\
 &\vdots \\
 l_{nk} &= a_{nk}^{(k)} / a_{kk}^{(k)}
 \end{aligned}$$

$$= I_n - l^{(k)} [e^{(k)}]^T$$

where

$$l^{(k)} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ l_{(k+1)k} \\ \vdots \\ l_{nk} \end{bmatrix}, \quad e^{(k)} - k\text{th column of } I_n$$

After $(n-1)$ steps

$$L^{(n-1)} \dots L^{(2)} L^{(1)} A = U$$

$$A = \underbrace{[L^{(1)}]^{-1} [L^{(2)}]^{-1} \dots [L^{(n-1)}]^{-1}}_L \cdot U$$

Observe

$$L^{(k)} = I_n - l^{(k)} [e^{(k)}]^T$$

$$\Rightarrow [L^{(k)}]^{-1} = I_n + l^{(k)} [e^{(k)}]^T$$

since

$$\begin{aligned} L^{(k)} [L^{(k)}]^{-1} &= (I_n - l^{(k)} [e^{(k)}]^T) (I_n + l^{(k)} [e^{(k)}]^T) \\ &= I_n - l^{(k)} [e^{(k)}]^T + l^{(k)} [e^{(k)}]^T \\ &\quad - \underbrace{l^{(k)} [e^{(k)}]^T l^{(k)} [e^{(k)}]^T}_0 \\ &= I_n \end{aligned}$$

Hence,

$$L = \underbrace{(I_n + l^{(1)} [e^{(1)}]^T) (I_n + l^{(2)} [e^{(2)}]^T) \dots (I_n + l^{(n-1)} [e^{(n-1)}]^T)}_{I_n + l^{(1)} [e^{(1)}]^T + l^{(2)} [e^{(2)}]^T + l^{(1)} [e^{(1)}]^T l^{(2)} [e^{(2)}]^T}$$

$$\begin{aligned}
&= \underbrace{\left(I_n + l^{(1)} [e^{(1)}]^T + l^{(2)} [e^{(2)}]^T \right) \left(I_n + l^{(3)} [e^{(3)}]^T \right) \dots \left(I_n + l^{(n-1)} [e^{(n-1)}]^T \right)}_{I_n + l^{(1)} [e^{(1)}]^T + l^{(2)} [e^{(2)}]^T + l^{(3)} [e^{(3)}]^T} \\
&\quad + \underbrace{l^{(1)} [e^{(1)}]^T}_{0} l^{(3)} [e^{(3)}]^T + \underbrace{l^{(2)} [e^{(2)}]^T}_{0} l^{(3)} [e^{(3)}]^T \\
&= I_n + l^{(1)} [e^{(1)}]^T + \dots + l^{(n-1)} [e^{(n-1)}]^T
\end{aligned}$$

$$= \begin{bmatrix} 1 & & & & & \\ l_{21} & 1 & & & & \\ l_{31} & l_{32} & 1 & & & \\ \vdots & \vdots & \vdots & \ddots & \ddots & \\ l_{n1} & l_{n2} & & & & 1 \end{bmatrix}$$

due to $l^{(1)} [e^{(1)}]^T$ due to $l^{(2)} [e^{(2)}]^T$

Pseudocode (LU Factorization)

Input : $A \in \mathbb{C}^{n \times n}$

Output : A lower triangular $L \in \mathbb{C}^{n \times n}$
an upper triangular $U \in \mathbb{C}^{n \times n}$
such that $A = LU$.

$L \leftarrow I_n$

for $k = 1, 2, \dots, n-1$

for $j = k+1, \dots, n$

$l_{jk} \leftarrow a_{jk} / a_{kk}$

$A(j, k:n) \leftarrow A(j, k:n) - l_{jk} A(k, k:n)$

↓ ↓
② ①

end

$U \leftarrow A$

of flops

$$\textcircled{1} \quad \sim (n-k)$$

$$\textcircled{2} \quad \sim (n-k)$$

Total # of flops

$$\sum_{k=1}^{n-1} \sum_{j=k+1}^n 2(n-k)$$

=

$$\sum_{k=1}^{n-1} 2 \underbrace{(n-k)}_l^2$$

=

$$\sum_{l=1}^{n-1} 2l^2 = 2 \frac{(n-1)n(2n-1)}{6}$$

$$\sim 2n^3/3$$

Condition Numbers

$$f: V \rightarrow W$$

V, W - vector spaces with norms

$$\kappa = \lim_{\delta \rightarrow 0^+} \sup_{\substack{\delta x \in V \\ 0 < \|\delta x\| \leq \delta}} \frac{\|f(x + \delta x) - f(x)\|}{\|\delta x\|}$$

absolute condition number

$$\tilde{\kappa} = \lim_{\delta \rightarrow 0^+} \sup_{\substack{\delta x \in V \\ 0 < \|\delta x\| \leq \delta}} \frac{\|f(x + \delta x) - f(x)\| / \|f(x)\|}{\|\delta x\| / \|x\|}$$

relative condition number

$$= \frac{\|x\|}{\|f(x)\|} \kappa$$

Examples

① Matrix-vector product

$$f: \mathbb{C}^{m \times n} \rightarrow \mathbb{C}^m$$

$$f(A) = Ab \text{ for a given } b \in \mathbb{C}^n$$

Using 2-norms
on $\mathbb{C}^{m \times n}$ and \mathbb{C}^m

$$\kappa = \lim_{\delta \rightarrow 0^+} \sup_{0 < \|\delta A\|_2 \leq \delta} \frac{\|(A + \delta A)b - Ab\|_2}{\|\delta A\|_2}$$

$$= \lim_{\delta \rightarrow 0^+} \sup_{0 < \|\delta A\|_2 \leq \delta} \frac{\|\delta A b\|_2}{\|\delta A\|_2}$$

For all $\delta A \in \mathbb{C}^{m \times n}$

$$\frac{\|\delta A b\|_2}{\|\delta A\|_2} \leq \frac{\|\delta A\|_2 \|b\|_2}{\|\delta A\|_2} = \|b\|_2$$

Choosing $\delta A = \delta \frac{b b^*}{\|b\|^2}$

$$\frac{\left\| \left(\delta \frac{b b^*}{\|b\|^2} \right) b \right\|_2}{\underbrace{\left\| \delta \frac{b b^*}{\|b\|^2} \right\|}_\delta} = \|b\|_2$$

Hence,

$$\kappa = \|b\|_2$$

$$\tilde{\kappa} = \frac{\|b\|_2 \|A\|_2}{\|A b\|_2}$$

If $A \in \mathbb{C}^{n \times n}$ and invertible, $b = A^{-1} A b$

$$\|b\|_2 \leq \|A^{-1}\|_2 \|A b\|_2,$$

so

$$\tilde{\kappa} \leq \|A^{-1}\|_2 \|A\|_2$$

② Linear Systems

Solution x of

$$Ax = b$$

as a function of $b \in \mathbb{C}^n$
(and a given invertible $A \in \mathbb{C}^{n \times n}$)

$$x : \mathbb{C}^n \rightarrow \mathbb{C}^n \quad \left(\begin{array}{l} \text{Using 2-norms} \\ \text{on } \mathbb{C}^n \end{array} \right)$$
$$x(b) = A^{-1}b$$

$$\kappa = \lim_{\delta \rightarrow 0^+} \sup_{0 < \|\delta b\|_2 \leq \delta} \frac{\|A^{-1}(b + \delta b) - A^{-1}b\|_2}{\|\delta b\|_2}$$

$$= \lim_{\delta \rightarrow 0^+} \underbrace{\sup_{0 < \|\delta b\|_2 \leq \delta} \frac{\|A^{-1}\delta b\|_2}{\|\delta b\|_2}}_{\|A^{-1}\|_2} = \|A^{-1}\|_2$$

$$\tilde{\kappa} = \frac{\|b\|_2}{\|A^{-1}b\|_2} \quad \kappa = \frac{\|b\|_2 \|A^{-1}\|_2}{\|A^{-1}b\|_2}$$

Exploiting $\|b\|_2 \leq \|A\|_2 \|A^{-1}b\|_2$

$$\tilde{\kappa} \leq \|A\|_2 \|A^{-1}\|_2$$

Backward Stability

 $f: V \rightarrow W$ exact problem $\hat{f}: V \rightarrow W$ computed solution
by the algorithm

The algorithm is called backward stable if

$$\hat{f}(x) = f(x + \delta x)$$

$$\exists \delta x \in V \text{ s.t. } \frac{\|\delta x\|}{\|x\|} = O(\epsilon_{\text{mach}})$$

for all $x \in V, x \neq 0$.Reminder

$$g(\epsilon_{\text{mach}}) = O(h(\epsilon_{\text{mach}}))$$

means there exists a constant c s.t.

$$|g(\epsilon_{\text{mach}})| \leq c |h(\epsilon_{\text{mach}})|$$

for all ϵ_{mach} close to 0.

e.g. $2\epsilon_{\text{mach}} = O(\epsilon_{\text{mach}})$

$2\epsilon_{\text{mach}} = O(\sqrt{\epsilon_{\text{mach}}})$

$$g(\epsilon_{\text{mach}}) = o(h(\epsilon_{\text{mach}}))$$

e.g.

means $\lim_{\epsilon_{\text{mach}} \rightarrow 0^+} \frac{g(\epsilon_{\text{mach}})}{h(\epsilon_{\text{mach}})} = 0$

$2\epsilon_{\text{mach}} = o(\sqrt{\epsilon_{\text{mach}}})$

①

Examples

① Addition in IEEE floating point arithmetic

Suppose $x, y \in \mathbb{R}$
is representable
in IEEE arith.

$$\begin{aligned}x \oplus y &= \text{fl}(x+y) \\ &= (x+y)(1+\epsilon)\end{aligned}$$

$$\exists \epsilon \text{ s.t. } |\epsilon| \leq \epsilon_{\text{mach}}$$

$$\left(\text{as } \frac{\{\text{fl}(x+y) - (x+y)\}}{(x+y)} \leq \epsilon_{\text{mach}}\right)$$

Hence,

$$x \oplus y = (x + \delta x) + (y + \delta y)$$

$$\text{where } \delta x = x \epsilon_{\text{mach}} \quad \delta y = y \epsilon_{\text{mach}}$$

$$\left\| \begin{bmatrix} \delta x \\ \delta y \end{bmatrix} \right\|_2 = |\epsilon| \left\| \begin{bmatrix} x \\ y \end{bmatrix} \right\|_2$$

$$\leq \epsilon_{\text{mach}} \left\| \begin{bmatrix} x \\ y \end{bmatrix} \right\|_2$$

Addition in IEEE arith.
is backward stable.

② Inner Product

$$f(x) = x^T y$$

$$= x_1 y_1 + x_2 y_2$$

for a given $y \in \mathbb{R}^2$, $y \neq 0$.

$$\hat{f}(x) = (x_1 \otimes y_1) \oplus (x_2 \otimes y_2)$$

$$\hat{f}(x) = \left\{ (x_1 y_1) (1 + \frac{\epsilon_1}{\epsilon_{mach}}) \right\} \oplus \left\{ (x_2 y_2) (1 + \frac{\epsilon_2}{\epsilon_{mach}}) \right\}$$

$$= \left[\left\{ (x_1 y_1) (1 + \epsilon_1) \right\} + \left\{ (x_2 y_2) (1 + \epsilon_2) \right\} \right] (1 + \epsilon_3)$$

$$= (x_1 y_1) (1 + \epsilon_1) (1 + \epsilon_3) + x_2 y_2 (1 + \epsilon_2) (1 + \epsilon_3)$$

$$= (x_1 + \delta x_1) y_1 + (x_2 + \delta x_2) y_2$$

where $\epsilon_1, \epsilon_2, \epsilon_3$ are s.t. $|\epsilon_1|, |\epsilon_2|, |\epsilon_3| \leq \epsilon_{mach}$

$$\delta x_1 := x_1 (1 + \epsilon_1) (1 + \epsilon_3) - x_1$$

$$= \cancel{x_1} (\epsilon_1 + \epsilon_3) x_1 + O(\epsilon_{mach}^2) x_1$$

$$\delta x_2 := \cancel{x_2} (1 + \epsilon_2) (1 + \epsilon_3) - \cancel{x_2}$$

$$= (\epsilon_2 + \epsilon_3) \cancel{x_2} + O(\epsilon_{mach}^2) x_2$$

Hence,

$$\hat{f}(x) = f(x + \delta x)$$

for some $\delta x = \begin{bmatrix} \delta x_1 \\ \delta x_2 \end{bmatrix}$ such that

$$\|\delta x\|_2 = \left\| \begin{bmatrix} (\epsilon_1 + \epsilon_3) x_1 + O(\epsilon_{mach}^2) x_1 \\ (\epsilon_2 + \epsilon_3) x_2 + O(\epsilon_{mach}^2) x_2 \end{bmatrix} \right\|_2 \leq 2 \epsilon_{mach} \|x\|_2 + O(\epsilon_{mach}^2)$$

$$\Rightarrow \|\delta x\| / \|x\| = O(\epsilon_{mach}).$$

Backward Error Analysis

Relative error

$$\tilde{E} := \frac{\|\hat{f}(x) - f(x)\|}{\|f(x)\|}$$

Now suppose the algorithm is backward stable, i.e.

$$\hat{f}(x) = f(x + \tilde{\delta}x)$$

$$\exists \tilde{\delta}x \quad \|\tilde{\delta}x\| \leq c \epsilon_{\text{mach}} \|x\|$$

where c is a constant independent of ϵ_{mach} .

$$\tilde{E} \leq \sup_{0 < \|\delta x\| \leq \delta} \frac{\|f(x + \delta x) - f(x)\| / \|f(x)\|}{\|\delta x\| / \|x\|} \cdot \frac{\|\delta x\|}{\|x\|}$$

for $\delta := c \epsilon_{\text{mach}} \|x\|$ $\tilde{\kappa}_\delta$

$$\tilde{E} \leq \tilde{\kappa}_\delta \frac{\|\tilde{\delta}x\|}{\|x\|} \Rightarrow \boxed{\tilde{E} = \tilde{\kappa}_\delta \cdot O(\epsilon_{\text{mach}})}$$

Example

Inner product in \mathbb{R}^n

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$f(x) = x^T y$$

$$= x_1 y_1 + \dots + x_n y_n$$

for a given $y \in \mathbb{R}^n$.

$$\tilde{K}_\delta = \sup_{0 < \|\delta x\| \leq \delta} \frac{|(x + \delta x)^T y - x^T y| / |x^T y|}{\|\delta x\| / \|x\|}$$

$$= \sup_{0 < \|\delta x\| \leq \delta} \frac{|\delta x^T y| \|x\|}{\|\delta x\| |x^T y|}$$

Observe

$$|\delta x^T y| \leq \|\delta x\| \|y\| \quad \forall \delta x$$

Specifically, for $\delta x_* := (\delta) \cdot y / \|y\|$

$$|\delta x_*^T y| = \delta \cdot \|y\| = \|\delta x_*\| \|y\|.$$

Hence

$$\tilde{K}_\delta = \frac{\|y\| \|x\|}{|x^T y|} = \frac{1}{|\cos \theta|}$$

where θ is the angle between x and y .

Computed in IEEE arith.

$$\hat{f}(x) = (x_1 \otimes y_1) \oplus \dots \oplus (x_n \otimes y_n)$$

It can be shown that (as we have done for the \mathbb{R}^2 case above)

$$\hat{f}(x) = f(x + \tilde{\delta}x)$$

$$\exists \delta x \text{ s.t. } \frac{\|\tilde{\delta}x\|}{\|x\|} = n \epsilon_{\text{mach}} + O(\epsilon_{\text{mach}}^2)$$

It follows that

$$\begin{aligned} \frac{\|\hat{f}(x) - f(x)\|}{|f(x)|} &\leq \frac{1}{|\cos \theta|} \{n \epsilon_{\text{mach}} + O(\epsilon_{\text{mach}}^2)\} \\ &= \frac{1}{|\cos \theta|} O(\epsilon_{\text{mach}}) \end{aligned}$$