New topic :

Direct and Iterative Methods for Solving Linear Systems

Goal : Solve $\quad Ax = b$

$\qquad$ where $\quad A \in \mathbb{R}^{n \times n}, \quad \text{rank}(A) = n$

---

Easiest case : $\quad A$ is upper/lower triangular matrix

$$A \qquad x \qquad = \qquad b$$

$$\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ & a_{12} & \cdots & a_{2n} \\ & & \ddots & \vdots \\ & & & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$$

- Backward substitution : $\quad x_n = b_n / a_{nn}$

$$x_{n-1} = (b_{n-1} - a_{n-1,n} \, x_n) / a_{n-1,n-1}$$

$$\vdots$$

$$x_i = (b_1 - \sum_{j=i+1}^{n} a_{ij} x_j) / a_{ii}$$

- Operation count :

$$\sum_{k=1}^{n} [2(n-k)+1] = n^2 \text{ FLOPs}$$

- Backward stability :

Backward substitution is very stable, analyzing

its backward error is very similar to that of inner product.

<span style="background-color:magenta">**Lemma,**</span> If $\quad y = (c - \sum_{i=1}^{k-1} a_i b_i) / b_k$ is evaluated in floating

point arithmetic, no matter what the order of evaluation,

the computed $\hat{y}$ satisfies

$$b_k \, \hat{y} \, (1 + \theta_k^{(0)}) = c - \sum_{i=1}^{k-1} a_i b_i (1 + \theta_k^{(i)})$$

where $|\theta_k^{(i)}| \le \gamma_k := \dfrac{k \, \varepsilon_{mach}}{1 - k \varepsilon_{mach}}$

Proof: exercise (Hint: Use binary tree) 🔲

Applying this lemma to triangular system, we have

**Thm** (Backward substitution is backward stable)

Let the triangular system $Ax = b$, where $A \in \mathbb{R}^{n \times n}$ is

nonsingular, be solved by substitution, with any ordering.

then the computed solution $\hat{x}$ satisfies

$$(A + \Delta A)\hat{x} = b, \qquad |\Delta A| \le \gamma_n |A|$$

$$\nwarrow |\Delta a_{ij}| \le \gamma_n |a_{ij}|, \forall i,j$$
$$(\text{implies } \|\Delta A\| \le C_n \|A\|)$$

---

Triangular system is easy to solve and stable.

For general $A$, use different strategies to factor it

into triangular matrices

**Method 1:** **QR factorization**

$$A = QR, \qquad Q^T Q = I, \qquad R \text{ upper triangular}$$

$$Ax = b \iff QRx = b \iff Rx = Q^T b \iff x = R^{-1} Q^T b$$

- Implementation:

Operation count:

Step 1: Householder QR $A = QR \longrightarrow \approx 2n^3 - \frac{2}{3}n^3$
Step 2: Compute $d = Q^T b \longrightarrow \approx O(n^2)$
Step 3: Solve $Rx = d \longrightarrow = n^2$

$$\approx \frac{4}{3}n^3 \text{ FLOPs}$$

- **Stability :** Since 1) Householder $QR$ is backward stable,

   2) computed $\hat{Q}$ is very close to orthogonal matrix

   3) solving a triangular system is stable.

   solving linear system with Householder $QR$

   is also backward stable.

**Thm** Let $A \in \mathbb{R}^{n \times n}$ be nonsingular. Suppose we solve

the system $Ax = b$ with the aid of $QR$ factorization

computed by the Householder algorithm. The computed

$\hat{x}$ satisfies

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

where $\dfrac{\|\Delta A\|_2}{\|A\|_2} \leq C_n \, \varepsilon_{mach}$ , $\dfrac{\|\Delta b\|_2}{\|b\|_2} \leq C_n' \, \varepsilon_{mach}$

**Pf :** From the backward stability of back substitution,

(∗)  $(\hat{R} + \Delta \hat{R})\, \hat{x} = \hat{Q}^T b$ , with $\|\Delta \hat{R}\|_2 \leq C_n \varepsilon_{mach} \|\hat{R}\|_2 \quad \cdots (1)$

From the stability of Householder $QR$, $\exists\, Q$, $Q^T Q = I$ and

$\| Q - \hat{Q} \|_2 \leq C_n' \varepsilon_{mach} \quad \cdots (2)$,   $\| Q\hat{R} - A \|_2 \leq C_n'' \varepsilon_{mach} \|A\| \quad \cdots (3)$

Hence, apply $Q$ to (∗). we get

$$\left( A + \underbrace{(Q\hat{R} - A) + Q \Delta \hat{R}}_{=: \Delta A} \right) \hat{x} = b + \underbrace{Q(Q^T - \hat{Q}^T)\, b}_{=: \Delta b}$$

$\|\Delta A\|_2 \overset{(3)}{\leq} C_n'' \varepsilon_{mach} \|A\|_2 + \|\Delta \hat{R}\|_2$

$\overset{(1)}{\leq} C_n'' \varepsilon_{mach} \|A\|_2 + C_n \varepsilon_{mach} \left( \| Q\hat{R} - A\|_2 + \|A\|_2 \right)$

$\overset{(3)}{\leq} (2 C_n'' + C_n)\, \varepsilon_{mach} \|A\|_2$

$$\|\Delta b\|_2 \leq \|Q^T - \hat{Q}^T\|_2 \|b\|_2 \overset{(2)}{\leq} C_n' \, \varepsilon_{mach} \|b\|_2$$

See Higham Thm 19.5     for a refined bound     ▨

---

## Method 2:   LU factorization

$$A = L \ U$$

↑ (unit) lower triangular matrix

↖ upper triangular matrix

- Idea: Gaussian elimination

$$[A \, ; \, b] \longrightarrow [U \, ; \, d] \longrightarrow [I \, ; \, x]$$

↑ $\sum_{i=1}^{n-1} (n-i) + 2(n-i)^2 \approx \frac{2}{3} n^3$

↑ $O(n^2)$

\* : Reduce Linear system to a triangular system via elementary row operation

\* : At the $k^{th}$ stage, zero the element below the diagonal in the $k^{th}$ column     ( $A^{(k)} x = b^{(k)}$ )

$$a_{ij}^{(k+1)} = a_{ij}^{(k)} - m_{ik} a_{kj}^{(k)} , \quad i,j = k+1, \dots n$$

$$b_i^{(k+1)} = b_i^{(k)} - m_{ik} b_k^{(k)} , \quad i = k+1, \dots n$$

where $m_{ik} = a_{ik}^{(k)} / a_{kk}^{(k)}$

\* : operation count $\approx 2n^3/3$ FLOPs for each $b$

- GE in matrix form

ex.  $A = \begin{bmatrix} 2 & 1 & 1 \\ 4 & 3 & 3 \\ 8 & 7 & 9 \end{bmatrix}$

$$E_1 \qquad\qquad A \qquad\qquad E_1 A$$

$$\begin{bmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ -4 & 0 & 1 \end{bmatrix} \quad \begin{bmatrix} 2 & 1 & 1 \\ 4 & 3 & 3 \\ 8 & 7 & 9 \end{bmatrix} = \begin{bmatrix} 2 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 3 & 5 \end{bmatrix}$$

$$E_2 \qquad\qquad E_1 A \qquad\qquad E_2 E_1 A$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -3 & 1 \end{bmatrix} \quad \begin{bmatrix} 2 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 3 & 5 \end{bmatrix} = \begin{bmatrix} 2 & 1 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 2 \end{bmatrix} = U$$

$$\longleftarrow \text{upper triangular matrix}$$

$$E_2 E_1 A = U$$

$$\Rightarrow \quad A = \underbrace{E_1^{-1} E_2^{-1}}_{} U$$

$$=: L \longleftarrow \text{unit lower triangular matrix}$$

$$L = \overset{E_1^{-1}}{\begin{bmatrix} 1 & & \\ 2 & 1 & \\ 4 & & 1 \end{bmatrix}} \overset{E_2^{-1}}{\begin{bmatrix} 1 & & \\ & 1 & \\ & 3 & 1 \end{bmatrix}} = \begin{bmatrix} 1 & & \\ 2 & 1 & \\ 4 & 3 & 1 \end{bmatrix}$$

- Computing an LU factorization $A = LU$ is equivalent to solving the equation

$$a_{ij} = \sum_{r=1}^{\min(i,j)} l_{ir} u_{rj} \qquad i,j = 1,\cdots,n$$

These nonlinear equation are easily solved if examined in the right order.

Suppose now row $1,\cdots, k-1$ of $U$ are computed

$\qquad\qquad$ column $1,\cdots, k-1$ of $L$ are computed

Setting $l_{kk} = 1$,

$$a_{kj} = l_{k1} u_{1j} + l_{k2} u_{2j} + \cdots + l_{k,k-1} u_{k-1,j} + \boxed{u_{kj}}, \qquad j = k, \cdots, n$$

$$a_{ik} = l_{i1} u_{1k} + l_{i2} u_{2k} + \cdots + l_{i,k-1} u_{k-1,k} + \boxed{l_{ik}} u_{kk}, \quad i = k+1, \cdots, n$$

- Implementation ( Doolittle's method )

<span style="color:red">same flops as</span>
<span style="color:red">↓ GE</span>

for $k = 1, \ldots, n$

     for $j = k, \cdots, n$

(△)       $u_{kj} = a_{kj} - \sum\limits_{r=1}^{k-1} l_{kr} u_{rj}$ ⟶ $2(k-1)(n-k+1)$ FLOPS

     end
     for $i = k+1, \cdots, n$

(△△)     $l_{ik} = \left( a_{ik} - \sum\limits_{r=1}^{k-1} l_{ir} u_{rk} \right) \Big/ u_{kk}$ ⟶ $[2(k-1)+1](n-k)$ FLOPS

     end
end

Operation count $= \sum\limits_{k=1}^{n} 2(k-1)(n-k+1)$
$+ [2(k-1)+1](n-k)$
$\approx \frac{2}{3} n^3$

- Solve $Ax = b$ via LU factorization

Step 1: Compute LU factorization $A = LU$

Step 2: Solve $Ly = b$ and then $Ux = y$

Remark:

$[b_1, \cdots, b_m]$ , $b_i \in \mathbb{R}^n$
↓

1) We can solve $Ax = b$ or $AX = B$ by

applying GE to augmented matrix $[A \vdots b]$ or $[A \vdots B]$.

This is faster and save memory.

2) If $b_1, \ldots, b_m$ are not available at the same time,

applying LU first then solve $Ly_i = b_i$, $Ux_i = y_i$ is faster.

($O(n^2)$ for each new $b_i$)

3) Doolittle's method is mathematically equivalent to GE

if we identify

$$a_{kj} - l_{k1} u_{1j} - l_{k2} u_{2j} - \cdots - l_{k,s} u_{s,j} = a_{kj}^{(s+1)}, \quad j = k+1, \ldots, n$$

- Error Analysis

Note that GE & Doolittle simply correspond to different order of evaluating (△) (△△), which are of the form

$$y = \left( c - \sum_{i=1}^{k-1} a_i b_i \right) / b_k.$$

By the lemma at the beginning of the notes, the rounding error in evaluating such formula is independent of the ordering. Hence, the worst case errors of GE is the same as Doolittle.

we focus on Doolittle in the following.

By the lemma at the beginning of the notes, the computed $\hat{L}$, $\hat{U}$ satisfy

$$\left| a_{kj} - \sum_{i=1}^{k-1} \hat{l}_{ki} \hat{u}_{ij} - \hat{u}_{kj} \right| \leq \gamma_k \sum_{i=1}^{k} |\hat{l}_{ki}| |\hat{u}_{ij}|, \quad j \geq k$$

$$\left| a_{ik} - \sum_{j=1}^{k} \hat{l}_{ij} \hat{u}_{jk} \right| \leq \gamma_k \sum_{j=1}^{k} |\hat{l}_{ij}| |\hat{u}_{jk}|, \quad i > k$$

Thm If Doolittle's algorithm runs to completion, then the computed $\hat{L}$, $\hat{U}$ satisfy

$$\hat{L} \hat{U} = A + \Delta A, \quad |\Delta A| \leq \gamma_n |\hat{L}| |\hat{U}|$$

## Thm (Backward error of GE)

Let $A \in \mathbb{R}^{n \times n}$ and suppose GE produces computed LU factor $\hat{L}$, $\hat{U}$ and a computed solution $\hat{x}$ to $Ax = b$

then $\qquad (A + \Delta A) \hat{x} = b$, with $|\Delta A| \leq \gamma_{3n} |\hat{L}| |\hat{U}|$ $\qquad (\bigstar)$