

Tuesday 9-10-15

3) Gaussian elimination/direct methods in solving $Ax=b$

- ① • Solve $Ax=b$ directly, expensive $O(n^3)$ "exact" $O(10^{-16})$
• later you will see indirect or iterative methods $O(10^{-8})$ or set tolerance cheaper
• advantage of direct: once factor A can solve for any RHS \vec{b} .

Three cases arise when solving $Ax=b$

- 1) No solution: $b \notin R(A)$
- 2) infinitely many solutions $N(A) \neq \emptyset$ (non-trivial)
- ✓ 3) exactly one solution (A invertible)

① Elimination - forward and backward sub

- based on how to solve triangular systems.

- we use gaussian elimination to put A into upper triangular form U and then use back-substitution on U

• Recall a matrix is upper triangular $\Leftrightarrow u_{ij} = 0$ $i > j$ (below diagonal)

$$\begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{pmatrix} \sim \begin{pmatrix} u_{11} & u_{12} & \dots & u_{1n} \\ & u_{22} & & u_{2n} \\ & & \ddots & \vdots \\ & & & u_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}$$

algorithm) for $i = 1, \dots, n$
 $x_i = b_i$
 for $j = 1, \dots, i-1$
 $x_i = x_i - e_{ij} x_j$
 end
 $x_i = x_i / e_{ii}$
 end

Why do we care about triangular systems?
 - Use gaussian elimination (as we have been doing) to transform $Ax=b$ in upper tri system $Ux=y$ and use backsubstitution to solve

- The elementary row operations are applied to augmented system $[A|b]$ to introduce zeros below diagonal without changing solution

Elementary row ops

- ✓ 1) Multiply row by scalar
- 2) Add scalar multiples of 1 row to another
- 3) permute rows

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ | & & & | \\ a_{m-11} & & & a_{m-1n} \\ a_{m1} & \dots & & a_{mn} \end{bmatrix} \quad \begin{array}{l} 1) 0 \text{ all elements} \\ \text{in first column} \\ \text{below the diagonal} \end{array}$$

2) make sure pivot $\neq 0$ if it is swap in row

Three types of pivoting: (don't want pivot too small)

- ① partial: swap w/ largest in row (most comm)
- ② row: swap w/ largest in row or column
- ③ full: swap w/ largest in full submatrix (expensive)

$$\boxed{\text{row}(i') = \text{row}(i) - \frac{a_{i1}}{a_{11}} \text{row}(1)} \quad \forall i \ 2 \leq i \leq n$$

$$A' = \begin{bmatrix} a_{11} & a_{12} & \overset{\text{pivot}}{a_{13}} & \dots & a_{1n} \\ 0 & & & & \\ \vdots & & & & \\ 0 & & & & \end{bmatrix} \quad A' [2:m, 2:n]$$

• continue inductively on $(m-1) \times (n-1)$ submatrix

② LU decomposition

• So far, we have always written GU as equation on row operations.
we can write it in terms of matrix L. (premult by L) $LA = U$

• look at first ^{trans} transformation L_1
• 1st row of A should not be affected by L_1 : $\vec{e}_1 = [1, 0, 0, \dots, 0]$

$$L_1 = \begin{pmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{pmatrix}$$

$$\boxed{a'_{2j} = a_{2j} - \left(\frac{a_{21}}{a_{11}}\right) a_{1j}}$$

$$L_2 = \begin{bmatrix} -\frac{a_{21}}{a_{11}} & 1 & 0 & \dots & 0 \\ & \uparrow & & & \\ & \text{Column 1 of } L_2 \text{ multiplied by row 1 of } A & & & \\ & & \uparrow & & \\ & & \text{Column 2 of } L_1 \text{ row 2 of } A & & \end{bmatrix}$$

$$L_1 = \begin{bmatrix} 1 & & & & \\ -\frac{a_{21}}{a_{11}} & 1 & & & \\ -\frac{a_{31}}{a_{11}} & & 1 & & \\ \vdots & & & \ddots & \\ -\frac{a_{n1}}{a_{11}} & & & & 1 \end{bmatrix}$$

$$A' = L_1 A$$

$$2) A'' = L_2 A' = L_2 L_1 A$$

$$L_2 = \begin{bmatrix} 1 & & & \\ 0 & 1 & & \\ & -\frac{a_{32}}{a_{22}} & 1 & \\ & 0 & \frac{a_{12}}{a_{22}} & 1 \end{bmatrix}$$

A series of $n-1$ elimination steps:

$$U = L_{n-1} L_{n-2} \dots L_2 L_1 A$$

$$\Rightarrow A = (L_{n-1} L_{n-2} \dots L_2 L_1)^{-1} U$$

$$= \underbrace{(L_1^{-1} L_2^{-1} \dots L_{n-1}^{-1})^{-1}}_L U = LU$$

• since inverse of lower triangular matrix is lower tri and product of lower tri is lower tri

Special properties of L_i

- easy to compute inverses
- each takes a row & subtracts constant times another row
- since the inverse should undo whatever L_i - it should add that constant

\Rightarrow only a sign difference $L_i^{-1} = \begin{bmatrix} 1 & & & \\ \frac{a_{21}}{a_{11}} & 1 & & \\ & \frac{a_{31}}{a_{11}} & 1 & \\ & \frac{a_{41}}{a_{11}} & & 1 \end{bmatrix}$

• Also simple to multiply - multiplication factors preserved at original location

$$L_1^{-1} L_2^{-1} \dots L_{n-1}^{-1} = \begin{bmatrix} 1 & & & \\ a_{21} & 1 & & \\ \vdots & & \ddots & \\ a_{n1} & a_{n2} & \dots & 1 \end{bmatrix} = L$$

- fills in each column

⇒ L always has ones on diagonal (unit lower triangular)

• U has pivots on diagonal
 $\det(A) = \det(LU) = \det(L) \det(U)$
 $= \det(U)$

② LU with pivoting: $PA = LU$

P, permutation matrix takes care of row swaps.

$$\det(P) = \pm 1 \quad \text{So}$$

$$\det(PA) = \det(P) \det(A) = \det(U)$$

$$\Rightarrow \det(A) = \pm \det(U)$$

Why do we need P?

Gaussian elimination from top to bottom

can fail (ignore 0 + small pivots)

numerically don't want to introduce large numbers

Ex: $A = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix}$ • swap rows 1 + 2 to get upper triangular.

- How to use to solve $Ax = b$
- Once have $A = LU$ can use to solve for any RHS b

$$Ax = L(Ux) = b. \text{ Let } Ux = y$$

- Solve $Ly = b$ for y (forward sub)
- Solve $Ux = y$ for x (back sub)

Further properties of matrix form.
Let us consider a vector $x \in \mathbb{R}^n$
and zeroing entries $i > j$

$$L_j \begin{pmatrix} x_1 \\ \vdots \\ x_j \\ x_{j+1} \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} x_1 \\ \vdots \\ x_j \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

\uparrow column of A

We can write it succinctly in outer product form:

$$\text{Recall } (e_i)_j = \begin{cases} 0 & j \neq i \\ 1 & j = i \end{cases} \quad (0 \text{ --- } | \text{ 0 --- } 0)$$

$$= I - \left(\frac{x_{j+1}}{x_j} \right) e_{j+1} e_j^T \quad (\text{To zero out } x_{j+1})$$

\uparrow look at $(j+1)$ th element

$$\left(I - \left(\frac{x_{j+1}}{x_j} \right) e_{j+1} e_j^T \right) x$$

$$= x - \frac{x_{j+1}}{x_j} e_{j+1} (e_j^T x)$$

$(0 \text{ --- } | \text{ 0 --- } 0) \begin{pmatrix} x_1 \\ \vdots \\ x_j \\ x_{j+1} \\ \vdots \\ x_n \end{pmatrix} = x_j$

$$= x - x_{j+1} e_{j+1} = \begin{pmatrix} x_1 \\ \vdots \\ x_{j+1} \\ \vdots \\ x_n \end{pmatrix} - \begin{pmatrix} 0 \\ \vdots \\ x_{j+1} \\ \vdots \\ 0 \end{pmatrix} = \begin{pmatrix} x_1 \\ \vdots \\ 0 \\ \vdots \\ x_n \end{pmatrix}$$

3x3 case

Example

$$L_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} - \begin{bmatrix} 0 & & \\ x_{21}/x_{11} & & \\ x_{31}/x_{11} & & \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & 0 \\ -x_{21}/x_{11} & 1 & 0 \\ -x_{31}/x_{11} & 0 & 1 \end{bmatrix}$$

To zero all elements below we define $L_j = I - \tau e_j^T$ put in j^{th} column

$$\tau_k = \begin{cases} 0 & k \leq j \\ \frac{x_k}{x_j} & k > j \end{cases}$$

Gaussian transform

We can use closed form to prove inverse property:

$$L_j^{-1} = I + \tau e_j^T$$

Proof: $(I - \tau e_j^T)(I + \tau e_j^T)$
 $= I - \tau e_j^T + \tau e_j^T - \tau(e_j^T \tau e_j^T)$

But $e_j^T \tau = 0$, since $\tau_j = 0$ by definition.

③ we have special factorizations to exploit special properties in matrices:

- 1) symmetric $A = A^T$ (less storage)
- 2) positive definite: $\forall x \neq 0 \in \mathbb{R}^n$ $x^T A x > 0$

- all eigenvalues are positive
- positive diagonals $a_{ii} > 0$
- invertible

Ex: $A^T A$ is positive definite

$$x^T A^T A x = (Ax)^T Ax = \|Ax\|_2^2 \geq 0$$

$$= 0 \Leftrightarrow x = 0$$

- negative definite: $z^T M z < 0$, $z \neq 0$
- positive semi-definite: $x^T A x \geq 0$ (could have 0 eigenvalues) non-invertible
- indefinite: $x^T A x > 0$ and $y^T A y < 0$, $\exists x, y$
 ————— positive & negative eigenvalues

- A_r submatrix positive definite

Special Factorizations

A matrix is symmetric positive definite (spd)
 $\Leftrightarrow \exists$ lower triangular matrix L w/ real positive diagonal entries s.t.

$$\boxed{A = LL^T = R^T R} \quad L, \text{ unique}$$

Cholesky Factorization

- no need for pivoting - positive diagonal so can take sqrt

Symmetric Case

- In GE, we use row operations to introduce 0's below diagonal in each column
- To maintain symmetry we can apply the same operations to columns of matrix to introduce 0's in first row

$\boxed{LDL^T}$ Since A is symmetric, it is clear that strictly upper triangular portion of U must be symmetric to strictly lower triangular part of L

- We can't write $U = L^T$ since diagonals of L & U aren't the same

- L is unit lower w/ all 1s on diagonal
- + U has pivots on diagonal
- Store pivots on diagonal D
- Let $U = DL^T \Rightarrow A = LU = \boxed{LDL^T}$

Special case) Cholesky - positive definite

all diagonal elements $a_{ii} > 0$, so we can take the sqrt

$$A = L D^{\frac{1}{2}} D^{\frac{1}{2}} L^T$$

$$= \boxed{(L D^{\frac{1}{2}}) (L D^{\frac{1}{2}})^T}$$

since $D^T = D$

$$D^{-1} = \left(\frac{1}{d_{11}}, \dots, \frac{1}{d_{nn}} \right)$$

$$= \boxed{L_1 L_1^T}$$

since nonzero no pivoting necessary for Cholesky

- Exercise: Symmetric and triangular \Rightarrow diagonal

④ Conditioning and perturbation studies for $Ax = b$

• ill-conditioning

① Will small perturbations in the data + round-off error lead to small or large errors in solution vector x ?

First, we must define matrix and vector norms

• norm on a vector space V is a function

$\|\cdot\| : V \rightarrow \mathbb{R}^+$ u.s.o.s that satisfies:

1) absolute homogeneity: $\|\alpha v\| = |\alpha| \|v\| \quad \forall \alpha \in \mathbb{R}, v \in V$

2) subadditivity (triangle inequality)
 $\|u+v\| \leq \|u\| + \|v\|$

3) nondegeneracy $\|v\| = 0 \Leftrightarrow v = 0$

It gives the length of a vector + is a convex function

• L^p norms $\|x\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}$

1) Standard Euclidean 2 norm

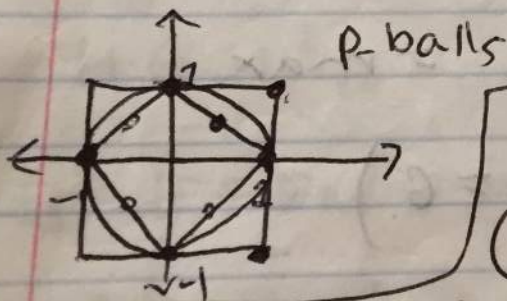
$$\|x\|_2 = \left(\sum_{i=1}^n |x_i|^2 \right)^{1/2} = \sqrt{x^T x}$$

↑ dot product

2) Infinity (maximum) norm

$$\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$$

$$3) \|x\|_1 = \sum_{i=1}^n |x_i|$$



Matrix norms

① Frobenius norm is extension

of Euclidean norm: $\|A\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{1/2}$

sqrt of sum of all elements squared

$$\boxed{\text{Tr}(A^T A)} = \sqrt{\text{Tr}(A A^T)}$$

$$\text{tr}(A^T A) = \sum_{i=1}^n \sum_{j=1}^n a_{ij}^T a_{ji} = \sum_{i=1}^n \sum_{j=1}^n a_{ji}^2 \checkmark$$

Submultiplicative $\|AB\|_F \leq \|A\|_F \|B\|_F$

2-norm $\|A\|_2 = \max_{x \neq 0} \frac{\|Ax\|_2}{\|x\|_2}$

$= \sup_{\|x\|_2=1} \|Ax\|_2$

- max amplitude of length of vector x when subjected to A
- difficult to calculate - will see later that it is the maximum singular value SVD

Other p -norms which are simpler to calculate:

• key idea is to express in terms of vector norms

$\|A\|_\infty = \sup_{\|x\|_\infty=1} \|Ax\|_\infty = \max \text{ row sum}$

$= \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|$ | Ex: $A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$
 $\|A\|_\infty = 7$

$\|A\|_1 = \sup_{\|x\|_1=1} \|Ax\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}| = \max$

col sum (Eg: $\|A\|_1 = 6$)

Proof: ~~state~~ we want $\max \|Ax\|_1$

s.t. $\|x\|_1 = \sum_{j=1}^n |x_j| = 1$

$\|Ax\|_1 = \left\| \sum_{j=1}^n x_j \vec{a}_j \right\|_1 = \sum_{j=1}^n |x_j| \|\vec{a}_j\|_1$
 triangle inequality
 scalar

$$\leq \sum_{j=1}^n |x_j| \left(\max_{1 \leq j \leq n} \|a_j\|_2 \right) \rightarrow \max_{1 \leq j \leq n} \|a_j\|_2$$

$$= \sqrt{\max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}|^2} \quad \checkmark$$

This is achieved when $x = e_j$

Important Properties

1) All induced matrix norms are submultiplicative: $\|AB\| \leq \|A\| \|B\|$
 $\|A \times\| \leq \|A\| \|x\|$

2) 2 norm and Frobenius are invariant under orthogonal transformations

$$\|QA\|_2 = \|A\|_2 \quad \text{for orthogonal matrix } Q^T Q = I = QQ^T$$

$$\|QA\|_F = \|A\|_F$$

Proof: $\|QA\|_2^2 = x^T A^T Q^T Q A x = (Ax)^T A x = \|Ax\|_2^2$

$$\|QA\|_F^2 = \text{tr}(A^T Q^T Q A) = \text{tr}(A^T A) = \|A\|_F^2 \quad \checkmark$$

① Effect on x^{\rightarrow} from perturbation in b^{\rightarrow}
 $Ax^{\rightarrow} = b^{\rightarrow}$

$A\vec{y} = \vec{b} + \delta\vec{b}$ ← perturbation in b^{\rightarrow}

$$\vec{y} = \vec{x} + \delta\vec{x}$$

$$\rightarrow A(\vec{x} + \delta\vec{x}) = \vec{b} + \delta\vec{b}$$

* How big is relative change in \vec{x} as a result of relative change in \vec{b} ?

$$\frac{\|\delta \vec{x}\|}{\|\vec{x}\|} \text{ related } \frac{\|\delta \vec{b}\|}{\|\vec{b}\|} ?$$

Subtract 2 equations: $A(x + \delta x) = b + \delta b$
 $Ax = b$

$$\Rightarrow A \delta x = \delta b \Rightarrow \delta x = A^{-1} \delta b$$

use
Submultiplicativity

$$\|\delta x\| = \|A^{-1} \delta b\| \leq \|A^{-1}\| \|\delta b\|$$

$$Ax = b \Rightarrow \|Ax\| \leq \|A\| \|x\| = \|b\|$$

$$\Rightarrow \frac{1}{\|x\|} \leq \frac{\|A\|}{\|b\|}$$

$$\Rightarrow \frac{\|\delta x\|}{\|x\|} \leq \|A^{-1}\| \|A\| \frac{\|\delta b\|}{\|b\|}$$

At worst the relative change is $\|A\| \|A^{-1}\| > \kappa(A)$ - condition number

ill-conditioned matrices: $\kappa(A) = \infty$

$$\text{Ex: } A = \begin{bmatrix} 1 & 1 \\ 1 & 1.001 \end{bmatrix}$$

↑
close to being singular

$$A^{-1} = \begin{bmatrix} 10001 & -10000 \\ -10000 & 10000 \end{bmatrix}$$

$\Rightarrow \kappa(A) \approx 40,000$ in Frobenius norm

- ill-conditioning is a property of the system, not algorithm. Check

$K(A)$ ~~is~~ ^{when} solving matrix vector eqn

② Round-off error accumulation (LU)

Ex: from small pivot

Solve $A\vec{x} = \vec{b}$ $A = \begin{bmatrix} 0.01 & 1 \\ 1 & -1 \end{bmatrix}$ $b = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

$K(A) \approx 3$ well-conditioned

$\sim \begin{bmatrix} 0.01 & 1 \\ 0 & -101 \end{bmatrix}$ $c = \begin{bmatrix} 1 \\ -100 \end{bmatrix}$ solve $Ux = c$

$x = \begin{bmatrix} 0.99099099 \\ 0.99099099 \end{bmatrix}$

Now suppose we can only work w/ 2pt floating pt arithmetic

-101 truncated to -100

$\Rightarrow x_2 = 1$ and $\frac{1}{100} x_1 = 1 - x_2 = 1 - 1 = 0$

$\begin{bmatrix} 0 \\ 1 \end{bmatrix} \Rightarrow$ huge change in answer!

How to write in floating point (2 digits)

$0.01 = 10 \times 10^{-3}$

$-101 = -10 \times 10^{-3}$

can't store the additional 7

$1 = 10 \times 10^{-1}$

Imagine this on a real example of $\epsilon = 10^{-16}$ machine precision.

* Problem is caused by small pivot 0.01. The large multiplication (100), $\ln 8$ in

$a_{22} = -1$ is overwhelmed +
truncated by large -100 subtraction
 \Rightarrow catastrophic cancellation (CC)

To avoid CC: swap rows for larger

pivot: $\frac{1}{100} \begin{bmatrix} 1 & -1 \\ 0.01 & 1 \end{bmatrix} x = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$

$$\sim \begin{bmatrix} 1 & -1 & | & 0 \\ 0 & 1.01 & | & 1 \end{bmatrix}$$

\uparrow
truncated
rounding error

$$\begin{bmatrix} 1 & -1 & | & 0 \\ 0 & 1 & | & 1 \end{bmatrix} \quad \text{w/ rounding error}$$

$$x_2 = 1 \quad x_1 = x_2 = 1$$

$$\begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

\Rightarrow very small perturbation of exact

* need to avoid small pivots

- partial pivoting is most common strategy

even well-conditioned matrices $(|A| \gg 1)$ can suffer from round off error (algorithmic design)

$PA = LU$ for more stable LU