

NLA 2021-2022

Linear equation solving

Martin Sombra

22 September 2021

Linear equation solving

We want to solve

$$Ax = b$$

for a nonsingular $n \times n$ matrix A and an n -vector b

Gaussian elimination consists in computing a factorization

$$A = PLU$$

with P a permutation, L a unit lower triangular, and U an upper triangular

Allows to solve $Ax = b$ by solving the simpler equations

- 1 $Pz = b$ (permute the entries of b)
- 2 $Ly = z$ (forward substitution)
- 3 $Ux = y$ (backwards substitution)

Constructing the PLU factorization

Set

$$A = \begin{bmatrix} a_{1,1} & \dots & a_{1,n} \\ \vdots & & \vdots \\ a_{n,1} & \dots & a_{n,n} \end{bmatrix}.$$

For $n = 1$

$$P = L = [1] \quad \text{and} \quad U = [a_{1,1}]$$

Constructing the PLU factorization (cont.)

Let $n \geq 2$ and choose k such that $a_{k,1} \neq 0$

Gaussian elimination with partial pivoting (GEPP)

k such that $|a_{k,1}|$ is maximal

Swap rows 1 and k premultiplying by the permutation matrix

$$P_1 = \begin{bmatrix} 0 & 0 & \cdots & 0 & 1 & 0 & \cdots & 0 \\ 0 & 1 & & & 0 & & & \\ \vdots & & \ddots & & \vdots & & & \\ 0 & & & 1 & 0 & & & \\ 1 & 0 & \cdots & 0 & 0 & & & \\ 0 & & & & & 1 & & \\ \vdots & & & & & & \ddots & \\ 0 & & & & & & & 1 \end{bmatrix}.$$

Consider the 2×2 -block

$$P_1^T A = \begin{bmatrix} a_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}$$

Constructing the PLU factorization (cont.)

Then

$$\begin{bmatrix} a_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix} = \begin{bmatrix} 1 & \mathbb{0} \\ L_{2,1} & \mathbb{1}_{n-1} \end{bmatrix} \begin{bmatrix} u_{1,1} & U_{1,2} \\ \mathbb{0} & \tilde{A}_{2,2} \end{bmatrix}$$

with

$$u_1 = a_{1,1}, \quad L_{2,1} = a_{1,1}^{-1} A_{2,1}, \quad U_{1,2} = A_{1,2} \quad \text{and} \quad \tilde{A}_{2,2} = A_{2,2} - L_{2,1} U_{1,2}$$

The $(n-1) \times (n-1)$ matrix $\tilde{A}_{2,2}$ is the *Schur complement*

By the inductive hypothesis (case $n-1$):

$$\tilde{A}_{2,2} = \tilde{P} \tilde{L} \tilde{U}$$

Constructing the PLU factorization (cont.)

Then

$$\begin{aligned} P_1^T A &= \begin{bmatrix} 1 & 0 \\ L_{2,1} & \mathbb{1}_{n-1} \end{bmatrix} \begin{bmatrix} u_{1,1} & U_{1,2} \\ 0 & \tilde{P} \tilde{L} \tilde{U} \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 \\ 0 & \tilde{P} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \tilde{P}^T L_{2,1} & \tilde{L} \end{bmatrix} \begin{bmatrix} u_{1,1} & U_{1,2} \\ 0 & \tilde{U} \end{bmatrix} \end{aligned}$$

Hence $A = P L U$ with

$$P = P_1 \begin{bmatrix} 1 & 0 \\ 0 & \tilde{P} \end{bmatrix}, \quad L = \begin{bmatrix} 1 & 0 \\ \tilde{P}^T L_{2,1} & \tilde{L} \end{bmatrix} \quad \text{and} \quad U = \begin{bmatrix} u_{1,1} & U_{1,2} \\ 0 & \tilde{U} \end{bmatrix}.$$

GEPP \rightsquigarrow $L = [l_{i,j}]_{i,j}$ with $|l_{i,j}| \leq 1$ for all i, j

Complete pivoting

Gaussian elimination with complete pivoting (GECP): takes $|a_{k,l}|$ maximal among all the entries
 \rightsquigarrow swaps the rows 1 and k , and the columns 1 and l

Gives a factorization

$$A = P_1 L U P_2$$

with P_1 and P_2 permutations

Can be more *numerically stable* but is also *more expensive* (in terms of speed)

The GEPP algorithm

for $i = 1, \dots, n - 1$

 swap row k and row i of A and L for k such that
 $|a_{k,i}| = \max_{i \leq p \leq n} |a_{p,i}|$

 for $j = i + 1, \dots, n$ (compute column i of L)

$$l_{j,i} \leftarrow \frac{a_{j,i}}{a_{i,i}}$$

 for $j = 1, \dots, n$ (compute row i of U)

$$u_{i,j} \leftarrow a_{i,j}$$

 for $j, k = i + 1, \dots, n$ (update $A_{2,2}$)

$$u_{j,k} \leftarrow a_{j,k} - l_{j,i} u_{i,k}$$

Useful remark: the i -th column of A only used to compute the i -th column of L , and the i -th row of A only used to compute the i -th row of U

for $i = 1, \dots, n-1$

swap row k and row i of A and L for k such that

$$|a_{k,i}| = \max_{i \leq p \leq n} |a_{p,i}|$$

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

$$\cancel{l_{j,i}} \leftarrow \frac{a_{j,i}}{a_{i,i}}, \text{ (replace by } a_{j,i} \text{)}$$

~~for $j = 1, \dots, n$~~

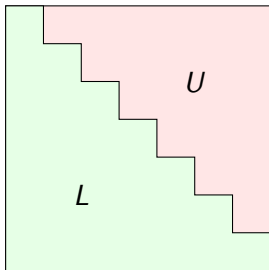
$$\cancel{u_{i,j}} \leftarrow \cancel{a_{i,j}}$$

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

$$\cancel{u_{j,k}} \leftarrow a_{j,k} - \cancel{l_{j,i}} \cancel{u_{i,k}} \text{ (replace by } a_{j,k}, a_{j,i} \text{ and } a_{i,k} \text{)}$$

Modified GEPP (cont.)

We need no extra space to store L and U !



Complexity

Recall the asymptotic formulae for $k \geq 1$:

$$\sum_{i=1}^n i^k = \frac{n^{k+1}}{k+1} + O(n^k),$$

The complexity (# flops) of GEPP on $n \times n$ matrices is

$$\sum_{i=1}^{n-1} \left(\sum_{j=i+1}^n 1 + \sum_{j=i+1}^n \sum_{k=i+1}^n 2 \right) = \sum_{i=1}^{n-1} ((n-i)+2(n-i)^2) = \frac{2}{3} n^3 + O(n^2)$$

Forward and backward substitution have each a complexity of

$$n^2 + O(n)$$

Hence GEPP solves the equation $Ax = b$ with complexity

$$\frac{2}{3} n^3 + O(n^2)$$

Floating point arithmetic

A *floating point number* is

$$f = \pm 0.d_1 d_2 \dots d_t \times \beta^e$$

$$\beta \geq 2, \quad 0 \leq d_i < \beta \text{ with } d_1 \neq 0 \quad \text{and} \quad L \leq e \leq U$$

$\beta \geq 2$ the *base*, L the *underflow* and U the *overflow*

The *range* is

$$\beta^{L-1} \leq |f| \leq \beta^U (1 - \beta^{-t}).$$

Floating point operations are defined by picking the closest floating point number:

$$a \odot b := \text{fl}(a \cdot b)$$

Machine epsilon: a bound for the relative error of roundoffs:

$$\frac{|a - \text{fl}(a)|}{|a|} \leq \varepsilon = \frac{\beta^{1-t}}{2}$$

Single precision IEEE standard

s	e	q
1	8	23
sign	exponent	coefficient/mantissa

The *coded number* is

$$f = (-1)^s (1 + q) 2^{e-127}$$

Corresponds to $t = 24$, $\beta = 2$, $L = -126$ and $U = 127$

The maximal *relative error* is

$$2^{-24} \approx 6 \cdot 10^{-8}$$

and the *range* is

$$2^{-127} \approx 6 \cdot 10^{-38} \leq f \leq 2^{129} (1 - 2^{-24}) \approx 7 \cdot 10^{39}$$

IEEE double precision standard

s	e	q
1	11	52
sign	exponent	coefficient/mantissa

The coded number is

$$f = (-1)^s (1 + q) 2^{e-1023}.$$

Corresponds to $t = 53$, $\beta = 2$, $L = -1022$ and $U = 1026$

The maximal error is

$$2^{-53} \approx 6 \cdot 10^{-16}$$

and the range is

$$2^{-1022} \approx 2 \cdot 10^{308} \leq f \leq 2^{1029} (1 - 2^{-24}) \approx 7 \cdot 10^{309}$$

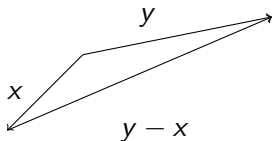
Vector and matrix norms

A *norm* on \mathbb{R}^n is a function

$$\| \cdot \|: \mathbb{R}^n \longrightarrow \mathbb{R}$$

such that

- 1 for every vector x we have that $\|x\| \geq 0$, and $\|x\| = 0$ if and only if $x = 0$
- 2 for every vector x and scalar α we have that $\|\alpha x\| = |\alpha| \|x\|$
- 3 for every vectors x, y we have that $\|x + y\| \leq \|x\| + \|y\|$
(*triangle inequality*).



Vector and matrix norms (cont.)

Example: for $1 \leq p \leq +\infty$ the p -norm is defined as

$$\|x\|_p = \begin{cases} (\sum_{i=1}^n |x_i|^p)^{1/p} & \text{if } 1 \leq p < +\infty \\ \max_i |x_i| & \text{if } p = +\infty \end{cases}$$

Any two norms can be compared: there are $c_1, c_2 > 0$ such that

$$\|\cdot\|_1 \leq c_1 \|\cdot\|_2 \quad \text{and} \quad \|\cdot\|_2 \leq c_2 \|\cdot\|_1$$

Example:

$$\|\cdot\|_2 \leq \|\cdot\|_1 \leq n^{1/2} \|\cdot\|_2 \quad \text{and} \quad \|\cdot\|_{\text{inf}} \leq \|\cdot\|_1 \leq n \|\cdot\|_{\infty}$$

Vector and matrix norms (cont.)

A *matrix norm* is a norm on $\mathbb{R}^{n \times n}$ s.t. for all A, B we have that

$$\|AB\| \leq \|A\| \|B\|$$

Example: the *Frobenius norm*

$$\|A\|_F = \left(\sum_{i,j} |a_{ij}|^2 \right)^{1/2}$$

Vector and matrix norms (cont.)

Definition: Given a norm on \mathbb{R}^n , the associated *operator norm* is

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|}$$

It is a matrix norm

Properties:

- ① $\|A\|_\infty = \max_i \sum_j |a_{i,j}|$ (*maximum row sum*)
- ② $\|A\|_1 = \max_j \sum_i |a_{i,j}|$ (*maximum column sum*)
- ③ $\|A\|_2 = \lambda_{\max}(A^*A)^{1/2}$ with $A^* = \overline{A}^T$ and λ_{\max} the maximal eigenvalue
- ④ if Q and Q' are *orthogonal* then

$$\|Q A Q'\|_2 = \|A\|_2 \quad \text{and} \quad \|Q A Q'\|_F = \|A\|_F$$

In particular $\|Q\|_2 = 1$ and $\|Q\|_F = n^{1/2}$

Perturbation theory

A matrix A is *well/badly (or ill) conditioned* if small changes in A can cause small/large changes in the solution of

$$Ax = b$$

Let x and $\hat{x} = x + \delta x$ solutions to

$$Ax = b \quad \text{and} \quad (A + \delta A)\hat{x} = b + \delta b$$

We have that

$$\begin{array}{rcl} (A + \delta A)(x + \delta x) & = & b + \delta b \\ - \quad Ax & = & b \\ \hline \delta Ax + (A + \delta A)\delta x & = & \delta b \end{array}$$

Then

$$\delta x = A^{-1}(-\delta A\hat{x} + \delta b)$$

The condition number

Fix a norm $\|\cdot\|$. Then

$$\|\delta x\| \leq \|A^{-1}\| (\|\delta A\| \|\hat{x}\| + \|\delta b\|)$$

or equivalently

$$\frac{\|\delta x\|}{\|\hat{x}\|} \leq \kappa_{\|\cdot\|} \left(\frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|A\| \|\hat{x}\|} \right) \quad (1)$$

with

$$\kappa_{\|\cdot\|} := \|A\| \|A^{-1}\|$$

the *condition number* of A with respect to $\|\cdot\|$

Precision of approximations

Let $\lambda \in \mathbb{R}$ and $\hat{\lambda} = \lambda + \delta\lambda$ an approximation of λ with k correct digits in base $\beta \geq 2$. Then

$$\lambda = \beta^e \times d_1 \cdots d_k d_{k+1} \cdots \quad \text{and} \quad \hat{\lambda} = \beta^e \times d_1 \cdots d_k \tilde{d}_{k+1} \cdots$$

and so

$$\frac{|\delta\lambda|}{|\lambda|} \leq \beta^{-k}$$

or equivalently

$$-\log_{\beta} \left(\frac{|\delta\lambda|}{|\lambda|} \right) \geq k.$$

Roundoff with IEEE single precision and double precision give approximations with 24 and 53 correct bits, respectively:

$$-\log_2 \left(\frac{|\delta_{\text{single}}\lambda|}{|\lambda|} \right) \geq 24 \quad \text{and} \quad -\log_2 \left(\frac{|\delta_{\text{double}}\lambda|}{|\lambda|} \right) \geq 53.$$

The inequality (1) translates into

$$-\log_{\beta} \frac{\|\delta x\|}{\|x\|} \geq -\log_{\beta} \kappa(A) - \log_{\beta} \left(\frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|A\| \|\hat{x}\|} \right)$$

Warning: ill-conditioned matrices destroy the quality of your approximations!

For instance, for *exact* data truncated with IEEE single or double precision, the computed result of $Ax = b$ will be meaningless as soon

$$\kappa(A) > 2^{24} \approx 6 \cdot 10^8 \text{ (single precision)}$$

and

$$\kappa(A) > 2^{53} \approx 10^{16} \text{ (double precision)}$$

Distance to the ill-posed problems

The condition number of the 2-norm has a geometrical interpretation as the inverse of its distance to the set of ill-posed problems:

$$\kappa_2(A) = \frac{1}{\text{distance}(A, \Sigma)} \quad (2)$$

with $\Sigma = \{A \mid \det(A) = 0\}$