

# 1. Linear systems: direct methods

Reference: Trefethen & Bau, Numerical Linear Algebra, SIAM, Lectures 20 & 21

Given  $n \in \mathbb{N} \setminus \{0, 1\}$ ,  $A \in \mathbb{R}^{n \times n}$ , and  $b \in \mathbb{R}^n$ , find  $x \in \mathbb{R}^n$  such that  $Ax = b$ .

$$A = \begin{bmatrix} a_{1,1} & \dots & a_{1,m} \\ \vdots & \ddots & \vdots \\ a_{m,1} & \dots & a_{m,m} \end{bmatrix} \quad b = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}$$

1. Triangular systems are easy to solve

$$A = \begin{bmatrix} a_{1,1} & & \\ & \ddots & \\ a_{m,1} & \dots & a_{m,m} \end{bmatrix} \quad \text{or} \quad A = \begin{bmatrix} a_{1,1} & \dots & a_{1,m} \\ & \ddots & \vdots \\ & & a_{m,m} \end{bmatrix}$$

lower triangular
upper triangular

→ forward substitution
→ back substitution

2. Gaussian elimination and LU factorization

Example:  $A = \begin{bmatrix} 2 & 1 & 1 & 0 \\ 4 & 3 & 3 & 1 \\ 8 & 7 & 9 & 5 \\ 6 & 7 & 9 & 8 \end{bmatrix}$  and  $b = \begin{bmatrix} 2 \\ 3 \\ 5 \\ 0 \end{bmatrix}$

Step 1:  $\begin{bmatrix} 2 & 1 & 1 & 0 & 2 \\ -2 & 1 & 1 & 1 & -1 \\ -4 & 3 & 5 & 5 & -3 \\ -3 & 4 & 6 & 8 & -6 \end{bmatrix} \quad \underbrace{\begin{bmatrix} 1 & & & \\ -2 & 1 & & \\ -4 & & 1 & \\ -3 & & & 1 \end{bmatrix}}_{=: L_1} \underbrace{\begin{bmatrix} 2 & 1 & 1 & 0 \\ 4 & 3 & 3 & 1 \\ 8 & 7 & 9 & 5 \\ 6 & 7 & 9 & 8 \end{bmatrix}}_{=: A} = \begin{bmatrix} 2 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 3 & 5 & 5 & 5 \\ 4 & 6 & 8 & 8 \end{bmatrix}$

Step 2:  $\begin{bmatrix} 2 & 1 & 1 & 0 & 2 \\ -2 & 1 & 1 & 1 & -1 \\ -4 & -3 & 2 & 2 & 0 \\ -3 & -4 & 2 & 4 & -2 \end{bmatrix} \quad \underbrace{\begin{bmatrix} 1 & & & \\ & 1 & & \\ -3 & & 1 & \\ -4 & & & 1 \end{bmatrix}}_{=: L_2} \begin{bmatrix} 2 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 3 & 5 & 5 & 5 \\ 4 & 6 & 8 & 8 \end{bmatrix} = \begin{bmatrix} 2 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ & 2 & 2 & 2 \\ & 2 & 4 & 2 \end{bmatrix}$

Step 3:  $\begin{bmatrix} 2 & 1 & 1 & 0 & 2 \\ +2 & 1 & 1 & 1 & -1 \\ +4 & -3 & 2 & 2 & 0 \\ +3 & -4 & 2 & 4 & -2 \end{bmatrix} \quad \underbrace{\begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & -1 & 1 \end{bmatrix}}_{=: L_3} \begin{bmatrix} 2 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ & 2 & 2 & 2 \\ & 2 & 4 & 2 \end{bmatrix} = \begin{bmatrix} 2 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ & 2 & 2 & 2 \\ & & & 2 \end{bmatrix} =: U$

to obtain L

Conclusion:  $x = \begin{bmatrix} 1 \\ -1 \\ 1 \\ -1 \end{bmatrix}$ ;  $L_3 L_2 L_1 A = U$  or  $A = LU$  with  $L := L_1^{-1} L_2^{-1} L_3^{-1}$

$L_1^{-1} = \begin{bmatrix} 1 & & & \\ 2 & 1 & & \\ 4 & 3 & 1 & \\ 6 & 7 & 9 & 1 \end{bmatrix}$ 
 $L_2^{-1} = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & 3 & & 1 \\ & 4 & & & 1 \end{bmatrix}$ 
 $L_3^{-1} = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & 1 & 1 \end{bmatrix}$ 
 $L = \begin{bmatrix} 1 & & & \\ 2 & 1 & & \\ 4 & 3 & 1 & \\ 6 & 7 & 9 & 1 \end{bmatrix}$

Algorithm:

```

for i = 1, ..., m-1 do
  for j = i+1, ..., m do
    aj,i ←  $\frac{a_{j,i}}{a_{i,i}}$ 
    for k = i+1, ..., m do
      aj,k ← aj,k - aj,i ai,k
    end
    bj ← bj - aj,i bi
  end
end
end
  
```

```

for i = m, ..., 1 do
  ri ←  $\frac{1}{a_{i,i}} \left( b_i - \sum_{j=i+1}^m a_{i,j} r_j \right)$ 
end
  
```

Number of arithmetic operations performed

$$\textcircled{*} \text{ Update of } b \quad \sum_{i=1}^{m-1} \sum_{j=i+1}^m 2 = 2 \sum_{i=1}^{m-1} (m-i) = 2 \sum_{i=1}^{m-1} i = 2 \frac{m(m-1)}{2} = m^2 - m$$

$$\begin{aligned} \textcircled{*} \text{ Update of } A, \text{ i.e., computation of } L \text{ and } U \\ \sum_{i=1}^{m-1} \sum_{j=i+1}^m \left( 1 + \sum_{k=i+1}^m 2 \right) &= \sum_{i=1}^{m-1} \sum_{j=i+1}^m 1 + 2 \sum_{i=1}^{m-1} \sum_{j=i+1}^m \sum_{k=i+1}^m 1 \\ &= 2 \sum_{i=1}^{m-1} (m-i) + \frac{1}{2} m^2 - \frac{1}{2} m \\ &= 2 \sum_{i=1}^{m-1} i^2 + \frac{1}{2} m^2 - \frac{1}{2} m \\ &= 2 \frac{(m-1)m(2m-1)}{6} + \frac{1}{2} m^2 - \frac{1}{2} m \\ &= \frac{2}{3} m^3 - \frac{1}{2} m^2 - \frac{1}{6} m \end{aligned}$$

$$\textcircled{*} \text{ Computation of } r \quad \sum_{i=1}^m (2 + 2(m-i)) = m + 2 \sum_{i=1}^m (m-i) = m + 2 \sum_{i=1}^{m-1} i = m^2$$

$$\begin{aligned} \text{Total: } (m^2 - m) + \left( \frac{2}{3} m^3 - \frac{1}{2} m^2 - \frac{1}{6} m \right) + m^2 &= \frac{2}{3} m^3 + \frac{3}{2} m^2 - \frac{7}{6} m \\ &\sim \frac{2}{3} m^3 \quad \left( \text{i.e., } \lim_{m \rightarrow \infty} \frac{\frac{2}{3} m^3 + \frac{3}{2} m^2 - \frac{7}{6} m}{\frac{2}{3} m^3} = 1 \right) \end{aligned}$$



⊛ Gaussian elimination

$$\begin{bmatrix} 8 & 7 & 9 & 5 & 5 \\ \frac{3}{4} & \frac{7}{4} & \frac{9}{4} & \frac{17}{4} & -\frac{15}{4} \\ \frac{1}{4} & -\frac{3}{4} & -\frac{2}{4} & \frac{4}{4} & -\frac{6}{4} \\ \frac{1}{2} & -\frac{2}{4} & -\frac{6}{4} & -\frac{2}{4} & -\frac{4}{4} \end{bmatrix}$$

$$L_2 := \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ & & & & 1 \end{bmatrix}$$

Step 3:

⊛ permutation of two rows

$$\begin{bmatrix} 8 & 7 & 9 & 5 & 5 \\ \frac{3}{4} & \frac{7}{4} & \frac{9}{4} & \frac{17}{4} & -\frac{15}{4} \\ \frac{1}{4} & -\frac{3}{4} & -\frac{2}{4} & \frac{4}{4} & -\frac{6}{4} \\ \frac{1}{2} & -\frac{2}{4} & -\frac{6}{4} & -\frac{2}{4} & -\frac{4}{4} \end{bmatrix}$$

$$P = \begin{bmatrix} 3 \\ 4 \\ 2 \\ 1 \end{bmatrix} \quad P_3 := \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ & & & & 1 \end{bmatrix}$$

⊛ Gaussian elimination

$$\begin{bmatrix} 8 & 7 & 9 & 5 & 5 \\ \frac{3}{4} & \frac{7}{4} & \frac{9}{4} & \frac{17}{4} & -\frac{15}{4} \\ \frac{1}{4} & -\frac{3}{4} & -\frac{2}{4} & \frac{4}{4} & -\frac{6}{4} \\ \frac{1}{2} & -\frac{2}{4} & -\frac{6}{4} & -\frac{2}{4} & -\frac{4}{4} \end{bmatrix}$$

$$L_3 := \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ & & & & 1 \end{bmatrix}$$

Conclusion:  $\alpha = \begin{bmatrix} 1 \\ -1 \\ 1 \\ -1 \end{bmatrix}$  and  $PA = LU$  with  $P = \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ & & & & 1 \end{bmatrix}$ ,

$$L = \begin{bmatrix} 1 & & & & \\ \frac{3}{4} & & & & \\ \frac{1}{4} & & & & \\ \frac{1}{2} & & & & \\ \frac{1}{4} & -\frac{3}{4} & -\frac{2}{4} & \frac{4}{4} & -\frac{6}{4} \end{bmatrix}, \text{ and } U = \begin{bmatrix} 8 & 7 & 9 & 5 \\ & \frac{7}{4} & \frac{9}{4} & \frac{17}{4} \\ & & -\frac{6}{4} & -\frac{2}{4} \\ & & & \frac{2}{3} \end{bmatrix}$$

Justification:

$$L_3 P_3 L_2 P_2 L_1 P_1 A = U$$

$$L_3 P_3 L_2 P_2 L_1 P_1 = L_3 \underbrace{(P_3 L_2 P_3)}_{=: L'_2} \underbrace{(P_3 P_2 L_1 P_2 P_3)}_{=: L'_1} P_3 P_2 P_1$$

$$PA = LU \text{ with } P := P_3 P_2 P_1 \text{ and } L := L'_1{}^{-1} L'_2{}^{-1} L_3$$

$$\text{Here, } L'_2 = \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & 1 & \\ & & & & 1 \end{bmatrix}$$

Algorithm:

$$p \leftarrow \begin{bmatrix} 1 \\ \vdots \\ m \end{bmatrix}$$

for  $i = 1, \dots, m-1$  do

Choose  $l \in \operatorname{argmax}_{m \in \{1, \dots, m\}} |a_{m,i}|$

if  $a_{l,i} = 0$  then

stop:  $A$  is singular

else

$$a_{i,:} \leftrightarrow a_{l,:}$$

$$b_i \leftrightarrow b_l$$

$$p_i \leftrightarrow p_l$$

$$a_{i+1:m,i} \leftarrow \frac{a_{i+1:m,i}}{a_{i,i}}$$

for  $j = i+1, \dots, m$  do

$$a_{j,i+1:m} \leftarrow a_{j,i+1:m} - a_{j,i} a_{i,i+1:m}$$

$$b_j \leftarrow b_j - a_{j,i} b_i$$

end

end

end

Number of comparisons of real numbers performed for pivoting:

$$\sum_{i=1}^{m-1} (m-i) = \sum_{i=1}^{m-1} i = \frac{1}{2} m(m-1).$$

Conclusion: in exact arithmetic, linear systems can be solved exactly within a finite number of arithmetic operations.

#### 4. Effects of round-off errors

Computers represent real numbers in a floating-point format — typically the IEEE 754 double-precision binary floating-point format (binary64) — and use floating-point arithmetic (which is not exact arithmetic).

Backward error analysis shows that, given floating-point representations of  $A$  and  $b$ , Gaussian elimination in floating-point arithmetic returns an exact solution to a perturbed system. Therefore, the roadmap is as follows:

1. review of background material about matrix norms and condition numbers;
2. based on the reviewed material, state a perturbation theorem;
3. based on a backward error analysis of Gaussian elimination, analyze the effect of round-off errors by applying the perturbation theorem.

Additional references:

- J. W. Demmel, Applied Numerical Linear Algebra, SIAM, 1997;
- G. W. Stewart, Matrix Algorithms: Basic Decompositions, SIAM, 1998;
- N. J. Higham, Accuracy and Stability of Numerical Algorithms, second edition, SIAM, 2002;
- G. H. Golub & C. F. Van Loan, Matrix Computations, <sup>fourth edition</sup> Johns Hopkins University Press, 2013.

#### 4.1. Matrix norm and condition number

Proposition. If  $\|\cdot\|$  is a norm on  $\mathbb{R}^m$ , then the function

$$\mathbb{R}^{m \times m} \rightarrow \mathbb{R} : A \mapsto \sup_{x \in \mathbb{R}^m, \|x\|=1} \|Ax\|$$

is a norm on  $\mathbb{R}^{m \times m}$  called the induced norm and denoted also by  $\|\cdot\|$ .

Example.

Norm on  $\mathbb{R}^m$

$$\|x\|_1 := \sum_{i=1}^m |x_i|$$

$$\|x\|_2 := \sqrt{\sum_{i=1}^m x_i^2}$$

$$\|x\|_\infty := \max_{i \in \{1, \dots, m\}} |x_i|$$

Induced norm on  $\mathbb{R}^{n \times n}$

$$\|A\|_1 = \max_{j \in \{1, \dots, m\}} \sum_{i=1}^m |a_{i,j}|$$

$$\|A\|_2 = \sigma_{\max}(A) = \sqrt{\lambda_{\max}(A^T A)}$$

$$\|A\|_\infty = \max_{i \in \{1, \dots, m\}} \sum_{j=1}^m |a_{i,j}|$$

Proposition. Every induced norm on  $\mathbb{R}^{n \times n}$  is submultiplicative: for all  $A, B \in \mathbb{R}^{n \times n}$ ,

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

Definition. The condition number of an invertible matrix  $A \in \mathbb{R}^{n \times n}$  is

$$\kappa(A) := \|A\| \|A^{-1}\|.$$

In this chapter, only induced norms are used.

#### 4.2. Sensitivity of linear systems

Theorem (Demmel 1997, equation (2.4))

Let  $\|\cdot\|$  denote both a norm on  $\mathbb{R}^n$  and the induced norm on  $\mathbb{R}^{n \times n}$ .

Let  $A \in \mathbb{R}^{n \times n}$  be invertible,  $b \in \mathbb{R}^n \setminus \{0\}$ ,  $\tilde{A} \in \mathbb{R}^{n \times n}$ , and  $\tilde{b} \in \mathbb{R}^n$ . If  $\frac{\|A - \tilde{A}\|}{\|A\|} < \frac{1}{\kappa(A)}$ , then  $\tilde{A}$  is invertible. If, moreover,  $x, \tilde{x} \in \mathbb{R}^n$  satisfy  $Ax = b$  and  $\tilde{A}\tilde{x} = \tilde{b}$ , then

$$\underbrace{\frac{\|x - \tilde{x}\|}{\|x\|}}_{\text{"relative error on } x"} \leq \frac{\kappa(A)}{1 - \kappa(A) \frac{\|A - \tilde{A}\|}{\|A\|}} \left( \underbrace{\frac{\|A - \tilde{A}\|}{\|A\|}}_{\text{"relative error on } A"} + \underbrace{\frac{\|b - \tilde{b}\|}{\|b\|}}_{\text{"relative error on } b"} \right)$$

### 4.3. Floating-point arithmetic

Computers represent real numbers in a floating-point format — typically the IEEE 754 double-precision binary floating-point format (binary64) — and use floating-point arithmetic (which is not exact arithmetic). The analysis of round-off errors is based on the existence of a positive real number  $\epsilon_{\text{machine}}$ , called machine epsilon, such that:

1. for every  $x \in \mathbb{R}$ , there exists  $\epsilon \in [-\epsilon_{\text{machine}}, \epsilon_{\text{machine}}]$  such that  $\text{fl}(x) = x(1 + \epsilon)$ , i.e.,  $\frac{\text{fl}(x) - x}{x} = \epsilon$ , where  $\text{fl}(x)$  is the floating-point number that is closest to  $x$ ;
2. for every arithmetic operation  $*$  (addition, subtraction, multiplication, or division) and every pair  $(x, y)$  of floating-point numbers, there exists  $\epsilon \in [-\epsilon_{\text{machine}}, \epsilon_{\text{machine}}]$  such that  $x \circledast y = (x * y)(1 + \epsilon)$ , where  $\circledast$  is  $*$  in floating-point arithmetic.

Property 1 is a property of the floating-point format, while Property 2 is a property of floating-point arithmetic. Property 1 actually holds only for  $x \in [-N_{\text{max}}, N_{\text{max}}] \setminus (-N_{\text{min}}, N_{\text{min}})$  but this has no practical impact. For example, with binary 64,  $N_{\text{max}} = (2 - 2^{-52}) 2^{1023} \approx 10^{308}$ ,  $N_{\text{min}} = 2^{-1022} \approx 10^{-308}$ , and  $\epsilon_{\text{machine}}$  is about  $2^{-52} \approx 2.22 \cdot 10^{-16}$  or even  $2^{-53} \approx 1.11 \cdot 10^{-16}$ .

### 4.4. Gaussian elimination in floating-point arithmetic

Two sorts of errors are involved.

The first sort comes from the representation of  $A$  and  $b$  in the chosen floating-point format:  $A$  becomes  $\text{fl}(A)$  and  $b$  becomes  $\text{fl}(b)$ , where  $\text{fl}$  is applied componentwise. Property 1 implies that  $|\text{fl}(A) - A| \leq \epsilon_{\text{machine}} |A|$  and  $|\text{fl}(b) - b| \leq \epsilon_{\text{machine}} |b|$ , where the absolute value is applied componentwise.

If  $A$  is ill-conditioned, i.e.,  $\kappa(A)$  is large, then the perturbation of  $f$  can be significant. If  $\epsilon_{\text{machine}} \kappa(A) = \frac{1}{2}$ , then the upper bound given by the perturbation theorem is merely upper bounded by 2 since  $\frac{\|A - f(A)\|_{\infty}}{\|A\|_{\infty}} \leq \epsilon_{\text{machine}}$  and  $\frac{\|b - f(b)\|_{\infty}}{\|b\|_{\infty}} \leq \epsilon_{\text{machine}}$ .

Example. Let  $A = \begin{bmatrix} 1 & 2^{-54} \\ 1 & 0 \end{bmatrix}$  and  $b = \begin{bmatrix} 1+2^{-54} \\ 1 \end{bmatrix}$ . Then,  $A^{-1}b = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ . With binary 64,  $f(A) = \begin{bmatrix} 1 & 2^{-54} \\ 1 & 0 \end{bmatrix}$  and  $f(b) = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ , and thus  $f(A)^{-1}f(b) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ . Since  $A^{-1} = \begin{bmatrix} 0 & 1 \\ 2^{54} & -2^{54} \end{bmatrix}$ ,  $\|A\|_{\infty} = 1+2^{-54}$ , and  $\|A^{-1}\|_{\infty} = 2^{55}$ , it holds that  $\kappa_{\infty}(A) = 2^{55+2}$  and the upper bound given by the perturbation theorem equals 2, twice the actual relative error.

In conclusion, if  $A$  is ill-conditioned, then algorithms cannot be expected to solve the system with high accuracy.

The second sort of errors comes from floating-point arithmetic. Backward error analysis shows that the round-off errors made in the course of Gaussian elimination can be projected back on the original matrix.

Theorem (Higham 2002, Theorems 9.3 and 9.4)

Given  $A$  and  $b$  in floating-point format, Gaussian elimination yields  $(L, U)$  such that  $LU - A \leq \frac{m \epsilon_{\text{machine}}}{1 - m \epsilon_{\text{machine}}} \|L\| \|U\|$  (provided that  $m \epsilon_{\text{machine}} < 1$ ).

Moreover, the computed <sup>1-m  $\epsilon_{\text{machine}}$  solution</sup>  $\tilde{x}$  satisfies  $\tilde{A} \tilde{x} = b$  with  $\|A - \tilde{A}\| \leq \frac{3m \epsilon_{\text{machine}}}{1 - 3m \epsilon_{\text{machine}}} \|L\| \|U\|$  (provided that  $3m \epsilon_{\text{machine}} < 1$ ). Thus,  $\frac{\|A - \tilde{A}\|_{\infty}}{\|A\|_{\infty}} \leq \frac{3m^3 \epsilon_{\text{machine}}}{1 - 3m \epsilon_{\text{machine}}} \cdot \underbrace{\frac{\|U\|_{\text{max}}}{\|A\|_{\text{max}}}}_{\text{"growth factor"}}$ .

Proof of the "Thus" statement. It holds that

$$\|A - \tilde{A}\|_{\infty} \leq \frac{3m \epsilon_{\text{machine}}}{1 - 3m \epsilon_{\text{machine}}} \|L\| \|U\|_{\infty}.$$

"growth factor"

$$\|A\|_{\text{max}} := \max_{i,j \in \{1, \dots, m\}} |a_{ij}|$$

Moreover,

$$\|L\| \|U\|_{\infty} \leq \|L\|_{\infty} \|U\|_{\infty} = \|L\|_{\infty} \|U\|_{\infty} \leq m \|U\|_{\infty} \leq m^2 \|U\|_{\text{max}}.$$

The result follows from the inequality  $\|A\|_{\text{max}} \leq \|A\|_{\infty}$ . ■



This yields the so-called normal equations:

$$A^T A x = A^T b.$$

Since  $\text{rank } A^T A = \text{rank } A = n$ , this system has a unique solution.

Example. Let  $A = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix}$  and  $b = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}$ . Then,  $b \notin \text{im } A$ . We compute

$$A^T A = \begin{bmatrix} 3 & 6 \\ 6 & 14 \end{bmatrix} \text{ and } A^T b = \begin{bmatrix} 4 \\ 8 \end{bmatrix}.$$

thus, the least-squares solution is  $\begin{bmatrix} 4/3 \\ 0 \end{bmatrix}$ .