

# SF2520 — Applied numerical methods

## Lecture 14

### Numerical linear algebra

Anna Nissen

Numerical analysis

Department of Mathematics, KTH

2023-11-13

# Today's lecture

- Direct methods for linear systems of equations (Gaussian elimination)
  - Full matrices, repetition
  - Sparse matrices
  - Banded matrices
- Stationary iterative methods

# Linear systems of equations

We consider

$$A\mathbf{x} = \mathbf{b}, \quad A \in \mathbb{R}^{N \times N}, \quad \mathbf{x}, \mathbf{b} \in \mathbb{R}^N.$$

- Finding  $\mathbf{x}$  given  $A$  and  $\mathbf{b}$  done using Gaussian elimination (GE).
- GE is called a "*direct method*", contains no approximations.
- ➊ Reduce the extended matrix  $A|\mathbf{b}$  to triangular form

The diagram shows the transformation of the augmented matrix  $[A|\mathbf{b}]$  into upper triangular form  $[U|\tilde{\mathbf{b}}]$  through four stages of Gaussian elimination. In each stage, the matrix is represented by a grid of 'x' marks. Brackets below the first three stages are labeled  $A$  and  $\mathbf{b}$ . The final stage is bracketed and labeled  $U$  and  $\tilde{\mathbf{b}}$ . The process involves using each row to eliminate elements below it, creating a staircase pattern of zeros.

Subtract multiple of current row from subsequent rows.

- ➋ Solve lower triangular system  $U\mathbf{x} = \tilde{\mathbf{b}}$  by backward substitution.

**Computational cost:** ➊ is  $O(N^3)$  and ➋ is  $O(N^2)$ .

Hard to improve if  $A$  is a full matrix.

# Linear systems of equations

## Connection to LU factorization

The reduction step ① can be seen as simultaneously

(a) Factorizing  $A = LU$  where

- $L$  is lower triangular with ones on diagonal
- $U$  is upper triangular

$$L = \begin{pmatrix} 1 & & & \\ \times & 1 & & \\ \vdots & \vdots & \ddots & \\ \times & \times & \times & 1 \end{pmatrix}, \quad U = \begin{pmatrix} \times & \times & \times & \times \\ & \ddots & \vdots & \vdots \\ & & \times & \times \\ & & & \times \end{pmatrix},$$

(b) Solving  $L\tilde{\mathbf{b}} = \mathbf{b}$  by forward substitution.

- Here  $U, \tilde{\mathbf{b}}$  are as in ① above and  $L$  is given by the GE multipliers.
- GE can then be seen as: LU factorization + forward substitution for  $L\tilde{\mathbf{b}} = \mathbf{b}$  + backward substitution for  $U\mathbf{x} = \tilde{\mathbf{b}}$ .

① LU factorize $A$	$\Rightarrow LU\mathbf{x} = \mathbf{b}$	Cost = $O(N^3)$
② Solve $L\tilde{\mathbf{b}} = \mathbf{b}$	$\Rightarrow U\mathbf{x} = \tilde{\mathbf{b}}$	Cost = $O(N^2)$
③ Solve $U\mathbf{x} = \tilde{\mathbf{b}}$	$\Rightarrow \mathbf{x}$	Cost = $O(N^2)$

- Note:  $LU$  factorization very useful when same linear system should be solved many times with different right hand sides, e.g. when time stepping PDEs with an implicit method.

# Pivoting

Pivoting often needed to stabilize GE.

$$\left( \begin{array}{cccccccc|c} \times & \times & \times & \times & \times & \times & \times & \times & \times \\ & \times & \times & \times & \times & \times & \times & \times & \times \\ & & \times & \times & \times & \times & \times & \times & \times \\ & & & \times & \times & \times & \times & \times & \times \\ & & & & \times_0 & x_1 & x_2 & x_3 & x_4 \\ & & & & y_0 & y_1 & y_2 & y_3 & y_4 \\ & & & & \times & \times & \times & \times & \times \\ & & & & \times & \times & \times & \times & \times \\ & & & & z_0 & z_1 & z_2 & z_3 & z_4 \end{array} \middle| \begin{array}{c} \times \\ \times \\ \times \\ \times \\ b_4 \\ b_5 \\ \times \\ \times \\ \times \end{array} \right) \Rightarrow \left( \begin{array}{cccccccc|c} \times & \times & \times & \times & \times & \times & \times & \times & \times \\ & \times & \times & \times & \times & \times & \times & \times & \times \\ & & \times & \times & \times & \times & \times & \times & \times \\ & & & \times & \times & \times & \times & \times & \times \\ & & & & z_0 & z_1 & z_2 & z_3 & z_4 \\ & & & & y_0 & y_1 & y_2 & y_3 & y_4 \\ & & & & \times & \times & \times & \times & \times \\ & & & & \times & \times & \times & \times & \times \\ & & & & x_0 & x_1 & x_2 & x_3 & x_4 \end{array} \middle| \begin{array}{c} \times \\ \times \\ \times \\ \times \\ b_8 \\ b_5 \\ \times \\ \times \\ b_4 \end{array} \right)$$

- Element  $\mathbf{x}_0$  called the "*pivot element*". In the reduction step, elements in the current row are divided by it.
- Small pivot elements introduces instabilities. Risk of large amplification of rounding errors.
- **Pivoting:** Changing around rows to get a large pivot element. Above:  $|z_0| > |x_0|$  is assumed. (Sometimes also columns are changed.)
- Gives modified  $LU$  factorization,

$$PA = LU,$$

where  $P$  is a permutation matrix. (An identity matrix with shuffled rows.)

If  $A$  is a non-singular matrix, then

- $LU$  and GE may not always work. e.g.  $A = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$
- $LU$  and GE with pivoting always works.
- If  $A$  is also symmetric positive definite<sup>1</sup>  $LU$  and GE always works without pivoting...
- ... and GE gives an even simpler "Cholesky" factorization  $A = LL^T$  (with  $L$  lower triangular).

**Remark:** Do not compute  $A^{-1}$  explicitly and multiply  $\mathbf{x} = A^{-1}\mathbf{b}$ . This is more expensive and less stable than solving  $A\mathbf{x} = \mathbf{b}$  with GE, in particular for sparse matrices.

---

<sup>1</sup>Or diagonally dominant,  $|a_{kk}| > \sum_{j \neq k} |a_{kj}|$  for all  $j$ .

# Sparse linear systems

- The system

$$A\mathbf{x} = \mathbf{b}, \quad A \in \mathbb{R}^{N \times N}, \quad \mathbf{x}, \mathbf{b} \in \mathbb{R}^N.$$

is **sparse** if the number of non-zero elements in  $A$  is  $\ll N^2$ .

- A typical situation is that the number of non-zero elements are  $O(N)$ , e.g. if each row or column has a fixed number of them.
- Cost of matrix-vector multiply is then  $O(N)$ .
- A special type is banded matrices...

$$A = \left( \begin{array}{ccc} \text{shaded band} & & 0 \\ & \neq 0 & \\ 0 & & \text{shaded band} \end{array} \right) \left. \begin{array}{l} \\ \\ \end{array} \right\} N, \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} p$$

Here  $p \ll N$ .

$p$  called the *bandwidth*.

Ex:  $p = 2$  for tridiagonal matrices.

# PDE discretizations

- Discretizations of PDEs lead to sparse systems, e.g. in elliptic equations or implicit time stepping methods for parabolic PDEs.
- Finite difference methods on simple geometries (with standard ordering) give banded systems:

$$A = \frac{1}{h^2} \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix} \in \mathbb{R}^{N \times N}$$

$$\begin{aligned} -u_{xx} &= f \text{ in 1D} \\ h &\sim 1/N, \\ \text{bandwidth} &= 2 \end{aligned}$$

$$\underbrace{A = \frac{1}{h^2} \begin{pmatrix} T & -I & & \\ -I & T & -I & \\ & \ddots & \ddots & \ddots \\ & & -I & T \end{pmatrix}}_{\in \mathbb{R}^{N \times N}}, \quad \underbrace{T = \begin{pmatrix} 4 & -1 & & \\ -1 & 4 & -1 & \\ & \ddots & \ddots & \ddots \\ & & -1 & 4 \end{pmatrix}}_{\in \mathbb{R}^{n \times n}}$$

$$\begin{aligned} -\Delta u &= f \text{ in 2D} \\ h &\sim 1/n, \\ \text{bandwidth} &= n \sim \sqrt{N} \end{aligned}$$

- Finite element methods also give banded systems with  $\approx$  same bandwidth if a good ordering of the elements is used. (Good ordering needed for FD too in general.)



## Remarks:

- Discretization of second derivatives ( $-\partial_{xx}$ ,  $-\Delta$ ) typically leads to symmetric positive definite matrices (with the right BC). These have real positive eigenvalues.
- Discretization of first derivatives ( $\partial_x$ ,  $\nabla$ ) typically leads to skew symmetric matrices  $A = -A^T$ , (with the right BC). These have purely imaginary eigenvalues.

# Solving sparse linear systems

Suppose  $A$  is sparse, when can we solve  $A\mathbf{x} = \mathbf{b}$  faster than  $O(N^3)$ ?

- If  $A$  is triangular and sparse, backward/forward substitution cost is  $O(\# \text{ non zero elements in } A)$ . (Most operations involve only zeros and can be skipped.)
- In general, if  $A = LU$ , both  $L$  and  $U$  must also be sparse, in order to be able to solve  $A\mathbf{x} = \mathbf{b}$  fast.

(Then there are few entries in  $L$  and  $U$  to compute, and both  $L\tilde{\mathbf{b}} = \mathbf{b}$  and  $U\mathbf{x} = \tilde{\mathbf{b}}$  can be solved fast by the previous point.)

- Unfortunately,  $A$  sparse does not imply that  $L$  and  $U$  are sparse.

**Example:**

$$\underbrace{\begin{pmatrix} \times & \times & \times & \times & \times & \times & \times & \times \\ \times & \times & & & & & & \\ \times & & \times & & & & & \\ \times & & & \times & & & & \\ \times & & & & \times & & & \\ \times & & & & & \times & & \\ \times & & & & & & \times & \\ \times & & & & & & & \times \end{pmatrix}}_A = \underbrace{\begin{pmatrix} \times & & & & & & & \\ \times & \times & & & & & & \\ \times & \times & \times & & & & & \\ \times & \times & \times & \times & & & & \\ \times & \times & \times & \times & \times & & & \\ \times & \times & \times & \times & \times & \times & & \\ \times & \times & \times & \times & \times & \times & \times & \\ \times & \times & \times & \times & \times & \times & \times & \times \end{pmatrix}}_L \underbrace{\begin{pmatrix} \times & \times & \times & \times & \times & \times & \times & \times \\ & \times & \times & \times & \times & \times & \times & \times \\ & & \times & \times & \times & \times & \times & \times \\ & & & \times & \times & \times & \times & \times \\ & & & & \times & \times & \times & \times \\ & & & & & \times & \times & \times \\ & & & & & & \times & \times \\ & & & & & & & \times \end{pmatrix}}_U$$

$L$  and  $U$  full. New non-zero entries in  $L$  and  $U$  called "fill-in".

# Solving sparse linear systems

- Appearance of fill-in depends a lot on column and row ordering. If we reverse columns and rows in the previous example we get:

$$\underbrace{\begin{pmatrix} \times & & & & & & & \times \\ & \times & & & & & & \times \\ & & \times & & & & & \times \\ & & & \times & & & & \times \\ & & & & \times & & & \times \\ & & & & & \times & & \times \\ & & & & & & \times & \times \\ \times & \times & \times & \times & \times & \times & \times & \times \end{pmatrix}}_{PAQ} = \begin{pmatrix} \times & & & & & & & \\ & \times & & & & & & \\ & & \times & & & & & \\ & & & \times & & & & \\ & & & & \times & & & \\ & & & & & \times & & \\ & & & & & & \times & \\ \times & \times & \times & \times & \times & \times & \times & \times \end{pmatrix} \begin{pmatrix} \times & & & & & & & \times \\ & \times & & & & & & \times \\ & & \times & & & & & \times \\ & & & \times & & & & \times \\ & & & & \times & & & \times \\ & & & & & \times & & \times \\ & & & & & & \times & \times \\ \times & \times & \times & \times & \times & \times & \times & \times \end{pmatrix}$$

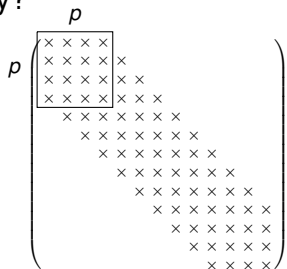
where  $P$ ,  $Q$  are permutation matrices.  $L$  and  $U$  sparse!

- When solving a linear system we are free to use any row and column ordering. Finding the optimal reordering (or permutations  $P$ ,  $Q$ ) to **minimize** fill-in is an NP-complete problem, however. Too expensive!
- Note also: Pivoting may restrict possible reorderings, or cause fill-in by itself.



# Solving sparse linear systems

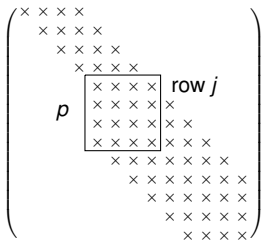
Why?



## First step

Row 1 is subtracted only from rows  $2, \dots, p$ . (Rows  $p, \dots, N$  already have a leading zero.)

$\Rightarrow$  Only the squared part of matrix is affected. Cost =  $O(p^2)$ .



## Step $j$

Row  $j$  is subtracted only from rows  $j+1, \dots, p+j-1$ .

$\Rightarrow$  Only the squared part of matrix is affected. Cost =  $O(p^2)$ .

Same cost  $O(p^2)$  in every step.

# Solving sparse linear systems

- Same cost  $O(p^2)$  in every step,  $A \in \mathbb{R}^{N \times N}$
- $\Rightarrow$  total cost, **GE for banded systems** =  $O(Np^2)$ .
- $U$  becomes upper triangular with bandwidth  $p$
- $L$  becomes lower triangular with bandwidth  $p$
- Backward/forward substitution has smaller cost than the reduction step:  $O(Np)$ .
- For Cholesky factorization  $A = LL^T$  cost will be same  $O(Np^2)$  and  $L$  as bandwidth  $p$ .
- When pivoting is needed,  $PA = LU$ , the  $L$  and  $U$  factors are *still sparse* but less so:
  - $U$  has bandwidth  $2p$
  - $L$  not banded, but lower triangular with  $p + 1$  non-zero entries per column. (I.e. forward substitution still fast, at  $O(Np)$ .)

# Costs for solving Poisson equation

Consider a finite difference discretization of

$$-\Delta u = f,$$

in  $d$  dimensions. This leads to a matrix equation

$$A\mathbf{u} = \mathbf{f}, \quad A \in \mathbb{R}^{N \times N}.$$

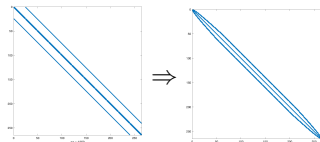
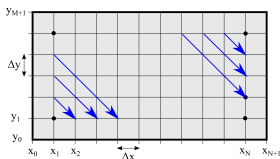
Suppose discretization has  $n$  grid points in each coordinate direction. Then

- Total number of points is  $N = n^d$ .
- Bandwidth is  $p = O(n^{d-1})$ .
- Since cost is  $O(Np^2)$  we get
  - Poisson in 1D. Then  $N = n$ ,  $p = 1$ , gives cost =  $O(N)$ .
  - Poisson in 2D. Then  $N = n^2$ ,  $p = n$ , gives cost =  $O(n^4) = O(N^2)$ .
  - Poisson in 3D. Then  $N = n^3$ ,  $p = n^2$ , gives cost =  $O(n^7) = O(N^{2\frac{1}{3}})$ .
- Higher dimension means not only more unknowns  $N$  but also higher cost to solve matrix equation for the same  $N$ .

# Sparse but not banded systems

## Remarks on sparse but not banded systems

- There are many heuristic methods (greedy algorithms) to reorder rows and columns to a better form that reduces fill-in. Works well in practice (e.g. for FEM) but no guarantee of success.
- Cost of these methods are  $O(\# \text{ non-zero elements in } A)$ , i.e. smaller than cost of LU-factorization and GE.
- Matlab-commands, examples:
  - "Reverse Cuthill–McKee" (RCM)  
`>> p = symrcm(A);`  $\Rightarrow A(p,p)$  reordered to low bandwidth.
  - Minimum degree, for  $A$  symmetric positive definite  
`>> p = symamd(A);`  $\Rightarrow A(p,p)$  has sparse Cholesky factors.
- **Ex.:** 2D Poisson with RCM gives ordering along *diagonals*. Can be solved in  $O(n^3)$  instead of  $O(n^4)$ !





# Iterative methods

We consider

$$A\mathbf{x} = \mathbf{b}, \quad A \in \mathbb{R}^{N \times N}, \quad \mathbf{x}, \mathbf{b} \in \mathbb{R}^N.$$

- Solve this via iteration  $\Rightarrow$

$$\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_2, \dots$$

- Stop e.g. when  $\|A\mathbf{x}_k - \mathbf{b}\| \leq \varepsilon$  or when  $\|\mathbf{x}_k - \mathbf{x}_{k-1}\| \leq \varepsilon$ , for some given tolerance  $\varepsilon$ .
- Primarily used for large sparse linear systems.
- Main cost of an iteration is typically a multiplication of  $\mathbf{x}_k$  by (a part of) the matrix  $A$ . Hence, cost of one iteration  $= O(\# \text{ non-zero elements of } A)$ .
- Methods are mostly insensitive to the precise sparsity structure. No need to worry about "fill-in" as in direct methods.
- Easier to exploit sparseness with iterative methods than with direct methods.

# Stationary methods

We consider

$$A\mathbf{x} = \mathbf{b}, \quad A \in \mathbb{R}^{N \times N}, \quad \mathbf{x}, \mathbf{b} \in \mathbb{R}^N.$$

Construct iteration as follows:

- Split  $A = M - T$ , so that  $M\mathbf{x} = T\mathbf{x} + \mathbf{b}$ .
- Iterate

$$M\mathbf{x}_{k+1} = T\mathbf{x}_k + \mathbf{b}.$$

- Choose  $M$  invertible, and so that  $M\mathbf{x} = \mathbf{f}$  is easy to solve.
- **Example 1:**  $M =$  diagonal of  $A$  gives "Jacobi method"
- **Example 2:**  $M =$  lower triangular part of  $A$  gives "Gauss–Seidel method"

# Stationary methods

## Jacobi and Gauss–Seidel methods

Let  $A = \{a_{j,\ell}\}$  and  $\mathbf{x}_k = \{x_j^k\}$ . Consider the methods elementwise, for a  $3 \times 3$  matrix.

Jacobi

$$\begin{aligned} a_{11}x_1^{k+1} + a_{12}x_2^k + a_{13}x_3^k &= b_1, \\ a_{21}x_1^k + a_{22}x_2^{k+1} + a_{23}x_3^k &= b_2, \\ a_{31}x_1^k + a_{32}x_2^k + a_{33}x_3^{k+1} &= b_3 \end{aligned} \quad \Rightarrow \quad \begin{aligned} x_1^{k+1} &= \frac{1}{a_{11}}(b_1 - a_{12}x_2^k - a_{13}x_3^k), \\ x_2^{k+1} &= \frac{1}{a_{22}}(b_2 - a_{21}x_1^k - a_{23}x_3^k), \\ x_3^{k+1} &= \frac{1}{a_{33}}(b_3 - a_{31}x_1^k - a_{32}x_2^k). \end{aligned}$$

Gauss–Seidel

$$\begin{aligned} a_{11}x_1^{k+1} + a_{12}x_2^k + a_{13}x_3^k &= b_1, \\ a_{21}x_1^{k+1} + a_{22}x_2^{k+1} + a_{23}x_3^k &= b_2, \\ a_{31}x_1^{k+1} + a_{32}x_2^{k+1} + a_{33}x_3^{k+1} &= b_3 \end{aligned} \quad \Rightarrow \quad \begin{aligned} x_1^{k+1} &= \frac{1}{a_{11}}(b_1 - a_{12}x_2^k - a_{13}x_3^k), \\ x_2^{k+1} &= \frac{1}{a_{22}}(b_2 - a_{21}x_1^{k+1} - a_{23}x_3^k), \\ x_3^{k+1} &= \frac{1}{a_{33}}(b_3 - a_{31}x_1^{k+1} - a_{32}x_2^{k+1}). \end{aligned}$$

Hence, Gauss–Seidel not much harder to implement than Jacobi. Just use the previously computed  $x_j^{k+1}$  values in the update, instead of the old ones  $x_j^k$ .

# Convergence

Consider the stationary method

$$M\mathbf{x}_{k+1} = T\mathbf{x}_k + \mathbf{b}.$$

- Can be written as fixed point iteration

$$\mathbf{x}_{k+1} = \phi(\mathbf{x}_k), \quad \phi(\mathbf{x}) = M^{-1}(T\mathbf{x} + \mathbf{b}).$$

- This iteration converges if spectral radius of jacobian  $< 1$ . Here,

$$\rho(M^{-1}T) < 1, \quad (\rho(\cdot) \text{ is the spectral radius.})$$

- For the Jacobi method this is fulfilled when  $A \in \mathbb{R}^{N \times N}$  is *diagonally dominant*,

$$|a_{kk}| > \sum_{\ell \neq k} |a_{k,\ell}|, \quad k = 1, \dots, N.$$

- For Gauss–Seidel it is fulfilled if  $A$  is diagonally dominant or symmetric positive definite.
- Gauss–Seidel usually converges a bit faster than Jacobi.