

SF2520 — Applied numerical methods

Lecture 15

Numerical linear algebra
Iterative methods, convergence

Anna Nissen
Numerical analysis
Department of Mathematics, KTH

2023-11-14

Today's lecture

- Iterative methods for linear systems of equations
 - Search methods (Krylov subspace methods)
 - Steepest descent
 - Conjugate gradient
 - GMRES
- Convergence theory
- Preconditioning

Iterative methods for linear systems of equations

We consider

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

- Solve this via iteration \Rightarrow

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

- Primarily used for sparse linear systems. Methods are mostly insensitive to the precise sparsity structure. No need to worry about "fill-in" as in direct methods.

- **Stationary methods:**

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

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

- Example 1: $M =$ diagonal of A gives "Jacobi method"
- Example 2: $M =$ lower triangular part of A gives "Gauss–Seidel method"

- In search methods, the iteration is of the form

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k, \quad \alpha_k \in \mathbb{R}, \quad \mathbf{p}_k \in \mathbb{R}^N,$$

where α_k and \mathbf{p}_k (the **search direction**) are computed based on the residual(s) $\mathbf{r}_\ell = \mathbf{b} - \mathbf{A}\mathbf{x}_\ell$, with $\ell \leq k$.

- **Examples:**

- Steepest descent method (SD)
 - Conjugate gradient method (CG)
 - GMRES ("generalized minimum residual" method)
- SD and CG require symmetric positive definite matrices. GMRES works for any nonsingular matrix.

Steepest descent method

Let A be a symmetric positive definite matrix. We note first that then

$$A\mathbf{x} = \mathbf{b} \quad \Leftrightarrow \quad \mathbf{x} \text{ minimizes } \frac{1}{2}\mathbf{x}^T A \mathbf{x} - \mathbf{x}^T \mathbf{b} =: \phi(\mathbf{x}),$$

since

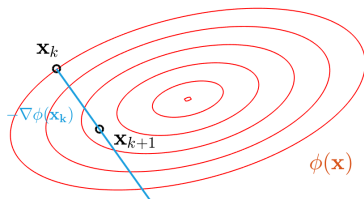
$$\nabla \phi(\mathbf{x}) = A\mathbf{x} - \mathbf{b}, \quad D^2 \phi(\mathbf{x}) = A > 0.$$

- In steepest descent we iterate

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla \phi(\mathbf{x}_k), \quad \mathbf{x}_0 = 0,$$

where α_k is chosen to minimize ϕ along the gradient direction,

$$\alpha_k = \operatorname{argmin}_{\alpha} \phi(\mathbf{x}_k - \alpha \nabla \phi(\mathbf{x}_k)).$$



Intuition: Move in direction of largest change ($\nabla \phi$) to lowest point in that direction.

Steepest descent method

We have

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k,$$

where

- The search direction is the residual,

$$\mathbf{p}_k = -\nabla \phi(\mathbf{x}_k) = \mathbf{b} - A\mathbf{x}_k = \mathbf{r}_k.$$

- α_k minimizes $\phi(\mathbf{x}_k + \alpha \mathbf{p}_k)$, i.e.

$$\begin{aligned} 0 &= \frac{d}{d\alpha} \phi(\mathbf{x}_k + \alpha \mathbf{p}_k) = \mathbf{p}_k^T \nabla \phi(\mathbf{x}_k + \alpha \mathbf{p}_k) = \mathbf{p}_k^T (A(\mathbf{x}_k + \alpha \mathbf{p}_k) - \mathbf{b}) \\ &= -\mathbf{p}_k^T \mathbf{r}_k + \alpha \mathbf{p}_k^T A \mathbf{p}_k. \end{aligned}$$

- Therefore,

$$\alpha_k = \frac{\mathbf{p}_k^T \mathbf{r}_k}{\mathbf{p}_k^T A \mathbf{p}_k} = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_k^T A \mathbf{r}_k}.$$

- Note also:

$$\mathbf{r}_{k+1} = \mathbf{b} - A\mathbf{x}_{k+1} = \mathbf{b} - A(\mathbf{x}_k + \alpha_k \mathbf{p}_k) = \mathbf{r}_k - \alpha_k A \mathbf{p}_k.$$

Steepest descent method

In summary, the method can be written

Steepest descent method

- 0 $\mathbf{x}_0 = \mathbf{0}, \mathbf{r}_0 = \mathbf{b}.$
- 1 $\mathbf{p}_k = \mathbf{r}_k$
- 2 $\alpha_k = \frac{\mathbf{p}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}$
- 3 $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$
- 4 $\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k$
- 5 Goto 1

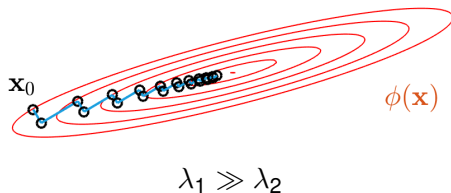
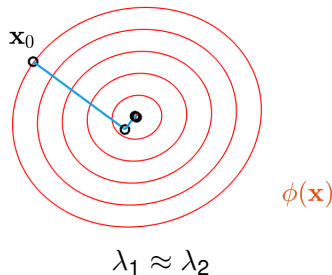
Note: The main cost of the iteration is the matrix-vector multiply $\mathbf{A} \mathbf{p}_k$. This should only be computed once per iteration!

(Of course in practice we would simplify and only keep \mathbf{r}_k , not \mathbf{p}_k in the algorithm. I kept \mathbf{p}_k here for comparison with the conjugate gradient method later on.)

Steepest descent method

Remark:

When the eigenvalues of A differ a lot, $\lambda_1 \gg \lambda_2$, the contour lines of ϕ become very elongated ellipses. This can lead to slow convergence. Compare:



Note that $\lambda_1 \gg \lambda_2$ means that the condition number $\kappa(A) \gg 1$. In general a large condition number of A leads to slow convergence.

Conjugate gradient method

The **conjugate gradient method** (CG) uses more optimally chosen search directions \mathbf{p}_k in

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k, \quad \mathbf{x}_0 = 0.$$

In the end, the rest of the algorithm will be the same as SD.

- **Observation:** Regardless of how \mathbf{p}_k are chosen we have

$$\mathbf{x}_{k+1} = \alpha_0 \mathbf{p}_0 + \alpha_1 \mathbf{p}_1 + \cdots + \alpha_k \mathbf{p}_k \in \text{span}\{\mathbf{p}_0, \dots, \mathbf{p}_k\}.$$

- Conjugate gradient idea:

- Let \mathbf{x}_{k+1} be the **best** value in all of $\text{span}\{\mathbf{p}_0, \dots, \mathbf{p}_k\}$,

$$\mathbf{x}_{k+1} = \underset{\mathbf{x} \in \text{span}\{\mathbf{p}_0, \dots, \mathbf{p}_k\}}{\text{argmin}} \quad \phi(\mathbf{x}).$$

- Choose \mathbf{p}_k so that:

- All $\{\mathbf{p}_k\}$ are linearly independent. (This implies full convergence in at most N iterations.)
- \mathbf{p}_{k+1} and \mathbf{x}_{k+1} can be computed easily from \mathbf{p}_k and \mathbf{x}_k .

Conjugate gradient method – sketch of derivation

- Let P be the matrix of search directions upto $k - 1$

$$P = \begin{pmatrix} | & | & \cdots & | \\ \mathbf{p}_0 & \mathbf{p}_1 & \cdots & \mathbf{p}_{k-1} \\ | & | & & | \end{pmatrix}$$

and, since \mathbf{x}_{k+1} should be in $\text{span}\{\mathbf{p}_j\}$, make the ansatz

$$\mathbf{x}_{k+1} = P\mathbf{y} + \alpha\mathbf{p}_k.$$

We need to determine \mathbf{y} , α and \mathbf{p}_k .

- They should be selected so that
 - \mathbf{p}_k is not in $\text{span}\{\mathbf{p}_0, \dots, \mathbf{p}_{k-1}\}$.
 - We minimize

$$\phi(\mathbf{x}_{k+1}) = \dots = \phi(P\mathbf{y}) + \alpha\mathbf{y}^T P^T A \mathbf{p}_k + \frac{1}{2}\alpha^2 \mathbf{p}_k^T A \mathbf{p}_k - \alpha \mathbf{p}_k^T \mathbf{b}.$$

- We therefore pick \mathbf{p}_k such that it is "A conjugate" to \mathbf{p}_j , $j < k$,

$$\mathbf{p}_j^T A \mathbf{p}_k = 0, \quad j = 0, \dots, k - 1$$

$\Rightarrow \{\mathbf{p}_k\}$ are linearly independent and $\alpha\mathbf{y}^T P^T A \mathbf{p}_k = 0$.

Conjugate gradient method – sketch of derivation

- We are left to minimize, over \mathbf{y} and α

$$\phi(\mathbf{x}_{k+1}) = \underbrace{\phi(P\mathbf{y})}_{\text{only depends on } \mathbf{y}} + \underbrace{\frac{1}{2}\alpha^2 \mathbf{p}_k^T A \mathbf{p}_k - \alpha \mathbf{p}_k^T \mathbf{b}}_{\text{only depends on } \alpha}.$$

- We can minimize the two terms independently:
 - Minimum of $\phi(P\mathbf{y}) = \phi(\mathbf{x}_k)$, the solution in the previous iteration, by the definition of how \mathbf{x}_k are chosen. Hence $P\mathbf{y} = \mathbf{x}_k$.
 - Optimal α in second term obtained by putting α -derivative to zero,

$$\alpha_k = \frac{\mathbf{p}_k^T \mathbf{b}}{\mathbf{p}_k^T A \mathbf{p}_k} = \dots = \frac{\mathbf{p}_k^T \mathbf{r}_k}{\mathbf{p}_k^T A \mathbf{p}_k}.$$

- We get

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$$

- Finally, one can show that

$$\mathbf{p}_k = \mathbf{r}_k + \beta_k \mathbf{p}_{k-1}, \quad \beta_k = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_{k-1}^T \mathbf{r}_{k-1}}$$

is A -conjugate to \mathbf{p}_j for $j < k$ (and easy to compute).

Conjugate gradient method

In summary, the method can be written

Conjugate gradient method

0 $\mathbf{x}_0 = \mathbf{p}_{-1} = \mathbf{0}, \mathbf{r}_0 = \mathbf{b}, \beta_0 = 0.$

1 $\mathbf{p}_k = \mathbf{r}_k + \beta_k \mathbf{p}_{k-1}$ and $\beta_k = \frac{\mathbf{r}_k^T \mathbf{r}_k}{\mathbf{r}_{k-1}^T \mathbf{r}_{k-1}}$ (when $k > 0$)

2 $\alpha_k = \frac{\mathbf{p}_k^T \mathbf{r}_k}{\mathbf{p}_k^T \mathbf{A} \mathbf{p}_k}$

3 $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$

4 $\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k \mathbf{A} \mathbf{p}_k$

5 Goto 1

Note: As in the steepest descent method, the main cost of the iteration is the matrix-vector multiply $\mathbf{A} \mathbf{p}_k$. It should only be computed once per iteration!

Conjugate gradient method

Remarks:

- The span of \mathbf{p}_k satisfies

$$\text{span}\{\mathbf{p}_0, \mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_k\} = \text{span}\{\mathbf{b}, A\mathbf{b}, A^2\mathbf{b}, \dots, A^k\mathbf{b}\},$$

which is called a **Krylov space**, denoted by $\mathcal{K}(\mathbf{b}, A, k)$.

- \mathbf{x}_k in each iteration minimizes the "A-norm" of the error,

$$\|\mathbf{x} - \mathbf{x}^*\|_A^2 =: (\mathbf{x} - \mathbf{x}^*)^T A(\mathbf{x} - \mathbf{x}^*), \quad (A\mathbf{x}^* = \mathbf{b} \text{ here}),$$

over all vectors in $\mathcal{K}(\mathbf{b}, A, k)$. (Note $\|\cdot\|_A$ is a norm since A is symmetric positive definite.)

- CG converges in at most N iterations. (Although this is seldom of practical use.)
- Convergence for CG is much faster than for SD.
- Like SD, convergence is slower for matrices with large condition numbers $\kappa(A)$.
- $\phi(\mathbf{x}_{k+1}) < \phi(\mathbf{x}_k)$, for both SD and CG.
- Keep in mind: SD and CG only work for symmetric positive definite matrices.

- GMRES is a Krylov method for non-symmetric A .
- \mathbf{x}_k minimizes $\|\mathbf{A}\mathbf{x}_k - \mathbf{b}\|$ over all vectors in $\mathcal{K}(\mathbf{b}, \mathbf{A}, k)$.
- Implementation and algorithm much more complicated than CG.
E.g. \mathbf{p}_k depends on **all** previous $\mathbf{p}_j, j \leq k$.
- Algorithm builds an orthonormal basis for $\mathcal{K}(\mathbf{b}, \mathbf{A}, k)$ via modified Gram-Schmidt (Arnoldi) in each iteration.
- Matlab commands:

```
>> pcg(A,b);  
>> gmres(A,b);
```


+ many other arguments and options possible

- Multigrid is an iterative method tailored for discretizations of elliptic PDEs (primarily).
- Solves PDE on different grid sizes N , $N/2$, $N/4$, \dots and uses differences between numerical solutions to eliminate coarse scale errors.
- Uses simple stationary methods (e.g. Jacobi) to eliminate fine scale errors.
- Highly efficient.

Computational costs

To achieve a certain given accuracy, the costs can be divided as

$$\text{computational cost} = \text{cost/iteration} \times \# \text{ iterations needed.}$$

Cost/iteration

We suppose $A \in \mathbb{R}^{N \times N}$ is sparse with $O(N)$ non-zero elements. (Typical for PDE discretizations.) Then

- Stationary methods

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

Since A sparse, then so is M and T , and the cost of $T\mathbf{x}_k$ is $O(N)$. If also $M\mathbf{x} = \mathbf{f}$ can be solved in $O(N)$ time (as it can for Jacobi and Gauss–Seidel), then cost/iteration is $O(N)$.

- Search methods

Main cost is the matrix-vector multiply $A\mathbf{p}_k$ which costs $O(N)$.

Conclusion: Cost/iteration is $O(N)$. This is true for any sparsity pattern with $O(N)$ non-zero elements, not just banded matrices.

- For the methods discussed one can derive convergence estimates of the type, that if $A\mathbf{x}^* = \mathbf{b}$, then

$$\|\mathbf{x}_k - \mathbf{x}^*\| \leq C\beta^k \|\mathbf{x}_0 - \mathbf{x}^*\|,$$

for some norm, constant C and $0 < \beta < 1$ called the **convergence rate**. It depends on A and on the method.

- Smaller β gives faster convergence. Often $\beta = 1 - \delta$ with $\delta \ll 1$ however.

Stationary methods

- We can write the stationary methods in concise form as

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

- Since $\mathbf{x}^* = R\mathbf{x}^* + \mathbf{c}$ we get

$$\mathbf{x}_k - \mathbf{x}^* = R\mathbf{x}_{k-1} + \mathbf{c} - R\mathbf{x}^* - \mathbf{c} = R(\mathbf{x}_{k-1} - \mathbf{x}^*) = R^k(\mathbf{x}_0 - \mathbf{x}^*)$$

- And therefore,

$$\|\mathbf{x}_k - \mathbf{x}^*\| \leq \|R^k\| \|\mathbf{x}_0 - \mathbf{x}^*\|.$$

- Finally, if R can be diagonalized, then

$$\|R^k\| \leq \rho(R)^k$$

where $\rho(R)$ is the spectral radius of R .

- We conclude that one can take $\beta = \rho(R)$, if R can be diagonalized.
- In stationary methods, the convergence rate thus depends on the structure of M and T in a complicated way.

Search methods

- In search methods β depends on the condition number κ of the matrix A .
- For steepest descent,

$$\beta = 1 - \frac{1}{\kappa}.$$

- For conjugate gradient,

$$\beta = 1 - \frac{1}{\sqrt{\kappa}} \quad \left(\text{which is } < 1 - \frac{1}{\kappa} \text{ for } \kappa \text{ large} \right)$$

- Well-conditioned matrices A give faster convergence for search methods.

Convergence rate vs # iterations needed

Suppose we want an error $\leq \varepsilon$ and $\beta = 1 - \delta$ with $0 < \delta \ll 1$. How many iterations are then needed?

- We need to find the number of iterations m such that

$$\|\mathbf{x}_m - \mathbf{x}^*\| \leq C \underbrace{\beta^m}_{e_0} \|\mathbf{x}_0 - \mathbf{x}^*\| \leq \varepsilon.$$

- This will be satisfied when

$$m \ln(\beta) \leq \ln\left(\frac{\varepsilon}{C e_0}\right) \Rightarrow m |\ln(\beta)| \geq \left| \ln\left(\frac{\varepsilon}{C e_0}\right) \right| \Rightarrow m \geq \frac{D}{|\ln(\beta)|},$$

where $D = |\ln(\varepsilon/(C e_0))|$ only depends on e_0 and ε .

- Since $\delta \ll 1$,

$$\ln(\beta) = \ln(1 - \delta) \approx -\delta.$$

- Consequently, we need, roughly,

$$m \approx \frac{D}{\delta} \text{ iterations.}$$

- Small δ means more iterations m needed. Precise number also depends on initial error e_0 and tolerance ε .

Computational costs

In summary, to achieve a certain given accuracy, we have

computational cost = cost/iteration \times # iterations needed.

- Assuming that the number of non-zero entries in A is $O(N)$, then

$$\text{Cost} = O(Nm) = O(N/\delta),$$

where m is the number iterations needed and $\beta = 1 - \delta$.

- For stationary methods $\beta = \rho(R)$, so $\delta = 1 - \rho(R)$. This gives

$$\text{Cost} = O\left(\frac{N}{1 - \rho(R)}\right).$$

- For steepest descent, $\delta = 1/\kappa$. This gives

$$\text{Cost} = O(N\kappa).$$

- For conjugate gradient, $\delta = 1/\sqrt{\kappa}$. This gives

$$\text{Cost} = O(N\sqrt{\kappa}).$$

(Also true for GMRES when A close to normal.)

Costs of 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. For **direct methods** we saw before that

Dimension	Unknowns (N)	Bandwidth (p)	Cost ($\sim Np^2$)
1	n	1	$O(n)$
2	n^2	n	$O(n^4)$
3	n^3	n^2	$O(n^7)$

Remark: Direct methods can sometimes be improved if problem is highly structured (square, constant coefficients, ...). E.g. $d = 2$ can be solved in $O(n^3)$ or even $O(n^2 \log n)$ time.

Costs of solving Poisson equation

For **iterative methods** we have when $A\mathbf{x} = \mathbf{b}$ comes from a finite difference discretization of Poisson:

- For Jacobi, Gauss–Seidel: $\rho(R) \sim 1 - 1/n^2$
(Can be derived from discretization.)
- Condition number $\kappa \sim h^{-2}$, where $h = O(1/n)$, for all d .
(We derived this in one dimension when we talked about elliptic equations.)
- Multigrid (MG): $\beta \leq \beta_0 < 1$ independent of N .

	δ	$m \sim 1/\delta$	cost ($Nm = n^d m$)
Jacobi	$1/n^2$	n^2	n^{d+2}
Gauss–Seidel	$1/n^2$	n^2	n^{d+2}
SD	$1/\kappa \sim 1/n^2$	n^2	n^{d+2}
CG	$1/\sqrt{\kappa} \sim 1/n$	n	n^{d+1}
GMRES	$1/\sqrt{\kappa} \sim 1/n$	n	n^{d+1}
Multigrid	$O(1)$	1	n^d

Costs of solving Poisson equation

In summary, we have the following costs for solving the Poisson equation with finite differences in 1D, 2D and 3D.

Dimension	Direct method	Jacobi/GS/SD	CG/GMRES	MG
1	n	n^3	n^2	n
2	n^4	n^4	n^3	n^2
3	n^7	n^5	n^4	n^3

- Little or no gain from iterative methods in 1D, compared to direct methods.
- Large gains in 3D.
- Unclear gains in 2D.
- Multigrid superior to the other methods.
- Jacobi/GS/SD not competitive. Mostly used as components in Multigrid or as "preconditioners".

Preconditioning

- For CG and GMRES the convergence rate is in general

$$\beta = 1 - \frac{1}{\sqrt{\kappa}} \quad \Rightarrow \quad m \sim \sqrt{\kappa}.$$

Slow convergence therefore for ill-conditioned matrices.

- Preconditioning** can be used to lower κ .

- Idea:

Replace $A\mathbf{x} = \mathbf{b}$ by equivalent system $M^{-1}A\mathbf{x} = M^{-1}\mathbf{b}$, where

(a) $\kappa(M^{-1}A)$ smaller than $\kappa(A)$

(b) $M\mathbf{x} = \mathbf{y}$ is easy to solve

Then iterate with CG/GMRES on $M^{-1}A\mathbf{x} = M^{-1}\mathbf{b}$ instead of $A\mathbf{x} = \mathbf{b}$.

\Rightarrow Faster convergence.

- Note that (a) and (b) are conflicting goals. Best for (a) is to take $M = A$, which gives $\kappa = 1$, but then (b) is not satisfied. Best for (b) is to take $M = I$, but then κ does not change. Find good trade-off!

The **preconditioner** M can be quite simple and still help. Some examples:

- $M =$ diagonal part of A (similar to one step with Jacobi, with $\mathbf{x}_0 = \mathbf{0}$).
- Multiple steps with Jacobi or Gauss–Seidel.
- Incomplete LU- or Cholesky factorizations.
- Good choice highly problem dependent.

Preconditioned Conjugate Gradient (PCG)

- Need extra trick to precondition conjugate gradient since $M^{-1}A$ is in general not symmetric positive definite (SPD) even if both M and A are.
- Let M be SPD and set $C = M^{1/2}$, i.e. an SPD matrix such that $C^2 = M$. (Such C exists.)
- Then $C^{-1}AC^{-1}$ is also SPD and since

$$A\mathbf{x} = \mathbf{b} \quad \Rightarrow \quad C^{-1}AC^{-1}C\mathbf{x} = C^{-1}\mathbf{b}$$

we can apply conjugate gradient to the SPD system

$$C^{-1}AC^{-1}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}, \quad C\mathbf{x} = \tilde{\mathbf{x}}, \quad C\tilde{\mathbf{b}} = \mathbf{b},$$

instead.

- This only changes CG very little in the end. (Main change: In each iteration an additional system $M\mathbf{z}_k = \mathbf{r}_k$ must be solved.)

Incomplete LU/Cholesky factorization

- Perform standard LU/Cholesky (via gaussian elimination) but let elements of L and U be zero whenever corresponding element of A is zero ($\ell_{j,k} = u_{j,k} = 0$ if $a_{j,k} = 0$).
- Gives approximate factors \tilde{L} and \tilde{U} with $\tilde{L}\tilde{U} \approx A$ for LU and $\tilde{L}\tilde{L}^T \approx A$ for Cholesky.
- In Matlab: `[L,U]=ilu(A)` and `L=ichol(A)`
- Cost = $O(N)$ if there is a fixed number of non-zeros in each row and column.
- Use $M = \tilde{L}\tilde{U}$ and $M = \tilde{L}\tilde{L}^T$ as preconditioners. Then $M\mathbf{x} = \mathbf{y}$ can easily be solved in $O(N)$ time.
- In Matlab, call `pcg` and `gmres` as
 - `>> pcg(A,b,tol,maxiter,M);` (if M is preconditioner)
 - `>> pcg(A,b,tol,maxiter,L,L');` (if LL^T is preconditioner)
 - `>> gmres(A,b,tol,maxiter,M);` (if M is preconditioner)
 - `>> gmres(A,b,tol,maxiter,L,U);` (if LU is preconditioner)