



# **Discussion on Iterative Solvers**

TMA4280—Introduction to Supercomputing

NTNU, IMF

February 27. 2017

## Problem



- Solve

$$\mathbf{A}\mathbf{x} = \mathbf{b}, \quad \mathbf{b}, \mathbf{x} \in \mathbb{R}^N, \quad \mathbf{A} \in M_N(\mathbb{R})$$

where  $\mathbf{A}$  can be the system resulting from discretizing a Poisson problem using finite differences.

- We use standard notation for matrices and vectors, i.e.

$$\begin{pmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,N} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,N} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N,1} & a_{N,2} & \cdots & a_{N,N} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_N \end{pmatrix}$$

## Direct methods



- Computing the inverse of the matrix is unrealistic,
- Matrix inversion has the same time complexity as matrix multiplication (typically  $\mathcal{O}(n^3)$ ).
- Direct methods can theoretically compute exact solutions  $\mathbf{x} \in \mathbb{R}^N$  to linear systems in the form of:

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$

- Several methods were introduced based on factorizations of the type  $\mathbf{A} = \mathbf{P}\mathbf{Q}$
- $\mathbf{P}$  and  $\mathbf{Q}$  have a structure simplifying the resolution of the system: diagonal, banded, triangular.
- The structure and properties of the matrix  $\mathbf{A}$  determine which algorithms we can use.

Ex: LU, Cholevski involve triangular matrices, QR constructs an orthogonal basis.

## Iterative methods



All methods prove to be quite expensive, hard to parallelize due to the sequential nature of the algorithm and prone to error propagation.

Iterative methods have been developed for:

- solving very large linear systems with direct methods is in practice not possible due to the complexity in term of computational operations and data,
- taking advantage of sparse system for which the structure of the matrix can result in dramatic speed-up (this is the case for numerical schemes for PDEs),
- using the fact that some systems like PDEs discretizations are already formulated in an iterative fashion.

## General idea



Introduce a splitting of the form:

$$A = G - H$$

such the solution  $\mathbf{x}$  satisfies:

$$G\mathbf{x} = \mathbf{b} + H\mathbf{x}$$

Similarly to fixed-point methods we can define a sequence of approximate solutions  $(\mathbf{x}^k)$  satisfying relations of the form:

$$G\hat{\mathbf{x}}^{k+1} = \mathbf{b} + H\hat{\mathbf{x}}^k$$

with  $G$  invertible.

## Link to linear mappings



The matrix viewed as a linear mapping in  $\mathbb{R}^N$ :

- the counterpart embodied by the Brouwer Theorem in finite dimension,
- Continuous mapping  $f : \Omega \rightarrow \Omega$  with  $\Omega$  compact of  $\mathbb{R}^N$ 
  1. Admits a fixed-point  $\mathbf{x}^*$  satisfying  $f(\mathbf{x}^*) = \mathbf{x}^*$ ,
  2. and is contracting.

Following a sequence of approximate solutions  $(\mathbf{x}^k)$  in  $\mathbb{R}^N$

## Computational aspects

Methods introduced depend on the iteration defined by the splitting:

1. How can the convergence be ensured?
2. How fast is the convergence?
3. How expensive is each iteration?
4. How does the algorithm behave with respect to numerical error?

Estimate on error vectors in terms of iteration error  $\hat{\epsilon}^k = \hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}^k$  or global error:  $\epsilon^k = \hat{\mathbf{x}}^k - \mathbf{x}$ . With  $A = G - H$

$$\hat{\epsilon}^k = G^{-1}H \hat{\epsilon}^{k-1}$$

Existence of a contraction factor  $K < 1$  such that  $\|\hat{\epsilon}^k\|_\infty \leq K \|\hat{\epsilon}^{k-1}\|_\infty$

Spectral radius  $\rho(M)$  as  $\rho(M) < 1$  since in that case  $\lim_{k \rightarrow \infty} M^k \hat{\epsilon}^0 = \mathbf{0}_{\mathbb{R}^N}$ .  
The smaller the spectral radius, the faster the convergence.

## Jacobi, methods of simultaneous displacements

$$(1) \quad \hat{\mathbf{x}}_i^{k+1} = \frac{1}{a_{ii}} \left( b_i - \sum_{i \neq j} a_{ij} \hat{\mathbf{x}}_j^k \right)$$

**Convergence:** the global error  $\epsilon^k$  is controlled by

$$\|\epsilon^{k+1}\| \leq \sum_{i \neq j} \left| \frac{a_{ij}}{a_{ii}} \right| \|\epsilon^k\| \leq K^k \|\epsilon^1\|$$

It is then enough if the matrix is strictly diagonally dominant. Expressing the iteration error gives  $M = G^{-1}H$  such that  $\rho(M) < 1$ .

1. Parallelization component by component is possible since there is only dependency on  $\hat{\mathbf{x}}^k$ .
2. Memory requirement for storing both  $\hat{\mathbf{x}}^{k+1}$  and  $\hat{\mathbf{x}}^k$  at each iteration.



## Gauss–Seidel, methods of successive displacements

In Jacobi iterations, notice that sequential ordered computation of terms

$$(2) \quad \hat{\mathbf{x}}_i^{k+1} = \frac{1}{a_{ii}} \left( b_i - \sum_{i \neq j} a_{ij} \hat{\mathbf{x}}_j^k \right)$$

involves components  $\hat{\mathbf{x}}_j^k$  which are also computed for  $\hat{\mathbf{x}}^{k+1}$  if  $j < i$ .

$$(3) \quad \hat{\mathbf{x}}_i^{k+1} = \frac{1}{a_{ii}} \left( b_i - \sum_{i < j} a_{ij} \hat{\mathbf{x}}_j^{k+1} - \sum_{i > j} a_{ij} \hat{\mathbf{x}}_j^k \right)$$

**Algorithm:** the splitting is

$$A = L - R_0$$

with  $L = D + L_0$  lower-triangular matrix and  $R_0$  strict upper-triangular matrix, thus

$$\hat{\mathbf{x}}^{k+1} = D^{-1} (\mathbf{b} - L_0 \hat{\mathbf{x}}^{k+1} + R_0 \hat{\mathbf{x}}^k)$$

## Gauss–Seidel, methods of successive displacements

Recast under the usual form:

$$\hat{\mathbf{x}}^{k+1} = \mathbf{L}^{-1}(\mathbf{b} + \mathbf{R}_0 \hat{\mathbf{x}}^k)$$


and the iteration matrix is  $\bar{\mathbf{M}} = \mathbf{L}^{-1} \mathbf{R}_0$ .

**Convergence:** the global error  $\epsilon^k$  is controlled by

$$\|\epsilon^{k+1}\| \leq \frac{\sum_{i>j} \left| \frac{a_{ij}}{a_{ii}} \right|}{1 - \sum_{i<j} \left| \frac{a_{ij}}{a_{ii}} \right|} \|\epsilon^k\| \leq \bar{K}^k \|\epsilon^1\|$$

If the Jacobi contraction factor  $K < 1$  then  $\bar{K} < 1$ . Expressing the iteration error gives directly that  $\bar{\mathbf{M}} = \mathbf{L}^{-1} \mathbf{R}_0$  such that  $\rho(\bar{\mathbf{M}}) < 1$ .

# Gauss–Seidel, methods of successive displacements



## Implementation:

1. Parallelization component by component is not possible easily since there is serialization for each row  $i$  due to the dependency on  $\hat{\mathbf{x}}_j^{k+1}$ ,  $j < i$ .
2. Memory requirement is only for storing one vector of  $\mathbb{R}^N$  at each iteration.

Parallelization possible if the matrix is sparse: components do not all possess connectivities with each other:

1. Component Dependency-Graph: generate a graph to reorder entries such that dependencies are avoided.
2. Red–Black coloring: special case for two-dimensional problems.

## Relaxation methods



Introduce the relaxation parameter  $\gamma \in (0, 1)$ : adding a linear combination of the approximate solution at the previous iteration to minimize the spectral radius for convergence.

The relaxation parameter  $\gamma$  cannot be known *a priori* and is usually determined by heuristics.

## Relaxation methods

1. Jacobi Over-Relaxation (JOR):

$$(4) \quad \hat{\mathbf{x}}_i^{k+1} = (1 - \gamma) \hat{\mathbf{x}}_i^k + \gamma \frac{1}{a_{ii}} \left( b_i - \sum_{i \neq j} a_{ij} \hat{\mathbf{x}}_j^k \right)$$

which reads in matricial form

$$\hat{\mathbf{x}}^{k+1} = M_\gamma \hat{\mathbf{x}}^k + \gamma D^{-1} \mathbf{b}$$

with  $M_\gamma = (1 - \gamma)\mathbb{I} + \gamma D^{-1} H$

2. Successive Over-Relaxation (SOR):

$$(5) \quad \hat{\mathbf{x}}_i^{k+1} = (1 - \gamma) \hat{\mathbf{x}}_i^k + \gamma \frac{1}{a_{ii}} \left( b_i - \sum_{i < j} a_{ij} \hat{\mathbf{x}}_j^{k+1} - \sum_{i > j} a_{ij} \hat{\mathbf{x}}_j^k \right)$$

which reads in matricial form

$$\hat{\mathbf{x}}^{k+1} = M_\gamma \hat{\mathbf{x}}^k + \gamma C \mathbf{b}$$

with  $M_\gamma = (1 + \gamma D^{-1} L_0^{-1})^{-1} [(1 - \gamma)\mathbb{I} + \gamma D^{-1} R_0]$  and

$$C = (1 + \gamma D^{-1} L_0^{-1})^{-1} D^{-1}$$

# Krylov-subspace methods



Idea: decomposition on a sequence of orthogonal subspaces.

If  $A$  is symmetric definite positive it induces the corresponding scalar product:

$$\langle \mathbf{x}, \mathbf{y} \rangle = (A\mathbf{x}, \mathbf{y}) = \mathbf{y}^T A\mathbf{x}$$

with  $(A\cdot, \cdot)$  canonical scalar product in  $\mathbb{R}^N$ . The vectors  $(\mathbf{e}_1, \dots, \mathbf{e}_N)$  are said  $A$ -conjugate if  $\mathbf{e}_j^T A\mathbf{e}_i = 0$  for  $i \neq j$ : they are orthogonal for the scalar-product induced by  $A$ .

To bring the Conjugate Gradient method, first let us introduce the idea of descent method.

## Descent methods



Minimisation of the residual:

$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}} J(\mathbf{x}) = \frac{1}{2} \langle \mathbf{x}, \mathbf{x} \rangle - \langle \mathbf{b}, \mathbf{x} \rangle$$

Construct a sequence of solutions to approximate minimization problems, given  $\hat{\mathbf{x}}^k$ :

$$J(\hat{\mathbf{x}}^{k+1}) \leq J(\hat{\mathbf{x}}^k)$$

where  $\hat{\mathbf{x}}^{k+1} = \hat{\mathbf{x}}^k + \alpha_{k+1} \mathbf{e}^{k+1}$ , with  $\alpha_{k+1}$  a descent factor and  $\mathbf{e}_{k+1}$  a direction.

# Steepest Gradient



For the Steepest Gradient:

1. take the direction given by  $-\nabla J(\hat{\mathbf{x}}^k) = \mathbf{b} - A\hat{\mathbf{x}}^k$  which is the residual  $\mathbf{r}_k = \mathbf{b} - A\hat{\mathbf{x}}^k$ , thus  $\hat{\mathbf{x}}^{k+1} = \hat{\mathbf{x}}^k + \alpha_{k+1}\mathbf{r}_k$ .
2. choose the descent factor  $\alpha^{k+1}$  minimizing the functional  $J(\hat{\mathbf{x}}^k + \alpha_{k+1}\mathbf{r}_k)$ :

$$\alpha_{k+1} = \frac{\mathbf{r}_k^T \mathbf{b}}{\mathbf{r}_k^T A \mathbf{r}_k}$$

1. Speed of convergence is bounded by  $\mathcal{O}(1 - \mathcal{C}(A)^{-1})$  with  $\mathcal{C}(A)$  the conditioning of  $A$ .
2. Gradient direction may not be optimal: Conjugate Gradient methods improve the choice of  $(\mathbf{e}_k)$ .



# Conjugate Gradient

The Conjugate Gradient (CG) is a Krylov-subspace algorithm for symmetric positive definite matrices.

Given  $\hat{\mathbf{x}}^0$ ,  $(\hat{\mathbf{x}}^k)$  is a sequence of solutions to approximate  $k$ -dimensional minimisation problems.

For the Conjugate Gradient:

1. take the direction  $\mathbf{e}_{k+1}$  such that  $(\mathbf{e}_1, \dots, \mathbf{e}_k, \mathbf{e}_{k+1})$  is  $A$ -conjugate, thus  $\hat{\mathbf{x}}^{k+1} = \hat{\mathbf{x}}^k + \alpha_{k+1} \mathbf{e}_{k+1}$ .
2. choose the descent factor  $\alpha^{k+1}$  minimizing the functional  $J(\hat{\mathbf{x}}^k + \alpha_{k+1} \mathbf{e}_k)$ , which is defined by

$$\alpha_j = \frac{\mathbf{e}_j^T \mathbf{b}}{\mathbf{e}_j^T \mathbf{A} \mathbf{e}_j}$$

and with  $\mathbf{e}_j^T \mathbf{b} \neq 0$  (unless the exact solution is reached).

# Conjugate Gradient



The construction of  $(\mathbf{e}_1, \dots, \mathbf{e}_{k+1})$  is done by orthogonalization of residuals by Gram–Schmidt:

$$\mathbf{e}_{k+1} = \mathbf{r}_k - \frac{\mathbf{e}_k^T \mathbf{A} \mathbf{r}_{k-1}}{\mathbf{e}_k^T \mathbf{A} \mathbf{e}_k}$$

so that  $\mathbf{r}_{k+1} = \mathbf{b} - \mathbf{A} \hat{\mathbf{x}}^{k+1} = \mathbf{r}^k - \alpha_{k+1} \mathbf{A} \mathbf{e}_{k+1}$

After  $N$  steps, the  $A$ -conjugate basis of  $\mathbb{R}^N$  is done and the exact solution is reached:

$$\mathbf{x} = \sum_{j=1}^N \alpha_j \hat{\mathbf{x}}^j$$

# Conjugate Gradient

For any  $k$ , the speed of convergence is bounded by

$$\mathcal{O} \left( \frac{1 - \sqrt{\mathcal{C}(A)}}{1 + \sqrt{\mathcal{C}(A)}} \right)^{2k}$$

in the norm induced by  $A$ , with  $\mathcal{C}(A)$  the conditioning of  $A$ .

The Conjugate Gradient can therefore be seen as a direct methods but in practice:

- the iterative computation of the  $A$ -conjugate basis suffers from the same issue of numerical error propagation as the QR factorization leading to a loss of orthogonality,
- the convergence is slow, which makes it unrealistic to compute the exact solution for large systems,

so it is used as an iterative method.



# Conjugate Gradient

First steps:

1. Given  $\hat{\mathbf{x}}^0 = \mathbf{0}$ , set  $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\hat{\mathbf{x}}^0$  and  $\mathbf{e}_1 = \mathbf{r}_0$ ,
2. Take  $\hat{\mathbf{x}}_1 = \alpha_1 \mathbf{e}_1$ , then  $\alpha_1 \mathbf{e}_1^T \mathbf{A} \mathbf{e}_1 = \mathbf{e}_1^T \mathbf{b}$ , thus

$$\alpha_1 = \frac{\mathbf{r}_0^T \mathbf{b}}{\mathbf{r}_0^T \mathbf{A} \mathbf{r}_0}$$

3. Compute the residual:

$$\mathbf{r}_1 = \mathbf{b} - \mathbf{A}\hat{\mathbf{x}}^1$$

4. Compute the direction:

$$\mathbf{e}_2 = \mathbf{r}_1 - \frac{\mathbf{e}_1^T \mathbf{A} \mathbf{r}_1}{\mathbf{e}_1^T \mathbf{A} \mathbf{e}_1} \mathbf{e}_1$$

5. Compute the factor:

$$\alpha_2 = \frac{\mathbf{e}_2^T \mathbf{b}}{\mathbf{e}_2^T \mathbf{A} \mathbf{e}_2}$$

6. Update the solution:

$$\hat{\mathbf{x}}^2 = \hat{\mathbf{x}}^1 + \alpha_2 \mathbf{e}_2$$

7. ...



# Conjugate Gradient

The algorithm iteration reads:

1. Compute the residual:

$$\mathbf{r}_k = \mathbf{b} - \mathbf{A}\hat{\mathbf{x}}^k$$

2. Compute the direction:

$$\mathbf{e}_{k+1} = \mathbf{r}_k - \frac{\mathbf{e}_k^T \mathbf{A} \mathbf{r}_{k-1}}{\mathbf{e}_k^T \mathbf{A} \mathbf{e}_k}$$

3. Compute the factor:

$$\alpha_{k+1} = \frac{\mathbf{e}_{k+1}^T \mathbf{b}}{\mathbf{e}_{k+1}^T \mathbf{A} \mathbf{e}_{k+1}}$$

4. Update the solution:

$$\hat{\mathbf{x}}^{k+1} = \hat{\mathbf{x}}^k + \alpha_{k+1} \mathbf{e}_{k+1}$$

which requires two matrix-vector multiplications per loop,  $\mathbf{A}\hat{\mathbf{x}}^k$  then  $\mathbf{A}\mathbf{e}_{k+1}$   
Using  $\mathbf{r}_{k+1} = \mathbf{r}^k - \alpha_{k+1} \mathbf{A}\mathbf{e}_{k+1}$  saves one matrix-vector multiplication.



# Conjugate Gradient

While the residual norm  $\rho_k = \|\mathbf{r}_k\|_2$  is big:

1. Compute the projection:

$$\beta_k = -\frac{\rho_k}{\rho_{k-1}}$$

2. Compute the direction:

$$\mathbf{e}_{k+1} = \mathbf{r}_k + \beta_k \mathbf{e}_k$$

3. Compute the factor:

$$\mathbf{w} = \mathbf{A} \mathbf{e}_{k+1}; \alpha_{k+1} = \frac{\rho_k}{\mathbf{e}_{k+1}^T \mathbf{w}}$$

4. Update the solution:

$$\hat{\mathbf{x}}^{k+1} = \hat{\mathbf{x}}^k + \alpha_{k+1} \mathbf{e}_{k+1}$$

5. Update the residual:

$$\mathbf{r}^{k+1} = \mathbf{r}^k - \alpha_{k+1} \mathbf{w}$$

Cost is one matrix-multiplication and five vector operations.



## Preconditioners

The convergence is still slow as soon as the condition number of the matrix is bad.

Preconditioning the system consists in finding a non-singular symmetric matrix  $C$  such that  $\tilde{A} = C^{-1}AC^{-1}$  and the conjugate gradient is applied to

$$\tilde{A}\tilde{\mathbf{x}} = \tilde{\mathbf{b}}$$

with  $\tilde{\mathbf{x}} = C^{-1}\mathbf{x}$  and  $\tilde{\mathbf{b}} = C^{-1}\mathbf{b}$ .

With:

- $M = C^2$
- $\mathbf{e}_k = C^{-1}\tilde{\mathbf{e}}_k$
- $\hat{\mathbf{x}}_k = C^{-1}\tilde{\mathbf{x}}_k$
- $\mathbf{z}_k = C^{-1}\tilde{\mathbf{r}}_k$
- $\mathbf{r}_k = C\tilde{\mathbf{r}}_k = \mathbf{b} - A\hat{\mathbf{x}}_k$

and  $M$  is a symmetric positive definite matrix called the preconditioner.

## Preconditioned CG

While the residual norm  $\rho_k = \|\mathbf{r}_k\|_2$  is big:

1. Solve:

$$\mathbf{M}\mathbf{z}_k = \mathbf{r}_k$$

2. Compute the projection:

$$\beta_k = -\frac{\mathbf{z}_k^T \mathbf{r}_k}{\mathbf{z}_{k-1}^T \mathbf{r}_{k-1}}$$

3. Compute the direction:

$$\mathbf{e}_{k+1} = \mathbf{z}_k + \beta_k \mathbf{e}_k$$

4. Compute the factor:

$$\alpha_{k+1} = \frac{\mathbf{z}_k^T \mathbf{r}_k}{\mathbf{e}_{k+1}^T \mathbf{A} \mathbf{e}_{k+1}}$$

5. Update the solution:

$$\hat{\mathbf{x}}^{k+1} = \hat{\mathbf{x}}^k + \alpha_{k+1} \mathbf{e}_{k+1}$$

6. Update the residual:

$$\mathbf{r}^{k+1} = \mathbf{r}^k - \alpha_{k+1} \mathbf{w}$$





## Preconditioned CG



The linear system  $M\mathbf{z}_k = \mathbf{r}_k$  should be easy to solve and can lead to fast convergence, typically  $\mathcal{O}(\sqrt{N})$ . Since

$$M\mathbf{z}_k = \mathbf{b} - A\hat{\mathbf{x}}_k$$

Then an iterative relation appears:

$$\hat{\mathbf{x}}_{k+1} = M^{-1}(\mathbf{b} - A\hat{\mathbf{x}}_k)$$

therefore iterative methods like Jacobi, Gauss-Seidel and relaxation methods can be used.

## Power method

Find the dominant eigenvalues of a matrix  $A \in M_N(\mathbb{R})$  of  $N$  eigenvectors ( $\mathbf{v}_i$ ) with associated eigenvalues ( $\lambda_i$ ) ordered in decreasing module. The eigenvalues are either real or conjugate complex pairs.

Given a random vector  $\mathbf{x}^0$ , construct a sequence of vectors ( $\hat{\mathbf{x}}^k$ ) such that

$$\hat{\mathbf{x}}^{k+1} = A\hat{\mathbf{x}}^k$$

then  $\forall k \geq 0$

$$\hat{\mathbf{x}}^k = \sum_{i=0}^{N-1} \lambda_i^k \xi_i \mathbf{v}_i$$

for some coefficients ( $\mathbf{v}_i$ ).

## Power method



Assume that  $\lambda_0$  is a dominant real eigenvalue and  $\xi_0 \neq 0$ , then

$$\hat{\mathbf{x}}^k = \lambda_0^k (\xi_0 \mathbf{v}_0 + \mathbf{r}_k)$$

with the residual  $\mathbf{r}_k$  defined as

$$\mathbf{r}_k = \lambda_0^{-k} \sum_{i=1}^{N-1} \lambda_i^k \xi_i \mathbf{v}_i$$

and  $\lim_{k \rightarrow \infty} \mathbf{r}_k = \mathbf{O}_{\mathbb{R}^N}$ . To the limit  $\hat{\mathbf{x}}_{k+1} \approx \lambda_0 \hat{\mathbf{x}}^k \approx \lambda_0 \xi_0 \mathbf{v}_0$  almost parallel to the first eigenvector.

## Power method



- This method is fast to compute the spectral radius for the Jacobi method and relaxation parameters.
- The convergence is geometric and the speed depends on the ratio  $|\lambda_1/\lambda_0|$ .
- If the matrix is symmetric, the convergence speed can be doubled.
- If  $\lambda_0$  is very large or very small then taking high powers lead to numerical issues, the algorithm requires a normalization.

# Software Packages



Libraries like PETSc and Trilinos offer interfaces to:

- a wide-range of iterative solvers based on Krylov-spaces,
- preconditioned by block Jacobi. ILU, AMG, ...
- so better design your software packages in consequence.