

## Solving Linear Systems with Iterative Methods

# 5 Solving Linear Systems with Iterative Methods

## 5.1 Splitting Methods

**Assumption** in this section:  $A \in \mathbb{R}^{n \times n}$  is invertible, so that  $x^* = A^{-1}b$  is the unique solution.

### 5.1.1 Motivation and Overview

**Problem:** *The matrix  $A$  can be very large ( $n \geq 10^5$ )!*

- If  $A$  can be stored, direct methods (such as LU, QR) are very slow or even not feasible due to large byproducts.
- Often  $A$  cannot be stored, so that direct methods are not an option.

**Example:**

- A float with double precision (standard) needs 8 bytes of memory.
- Considering 3d measurements with just 100 measurements in each dimension.  
 $\Rightarrow 100 \cdot 100 \cdot 100 = 10^6$  measurements
- Discretization methods (e.g., finite element method) for physical models (e.g., heat diffusion) interrelate these measurements.  
 $\Rightarrow$  gives a matrix  $A \in \mathbb{R}^{n \times n}$  with  $n = 10^6$   
 $\Rightarrow n \cdot n = 10^{12}$  numbers have to be stored  
 $\Rightarrow 10^{12}$  bytes = 1 **Terabyte** of memory has to be allocated in RAM ⚡
- A standard PC nowadays has 8–32 Gbytes of RAM (fast memory).

What if the matrix is *sparse* (= many 0 entries = redundancy)?

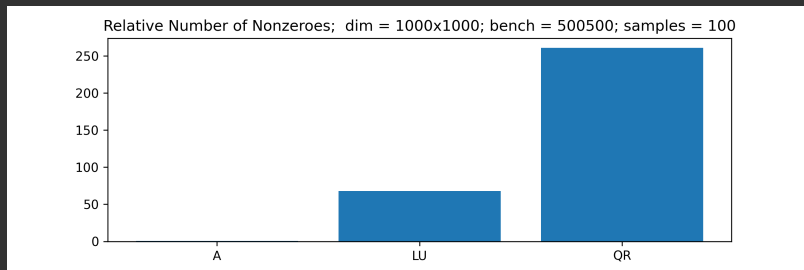
→ We only need to store nonzero entries and their coordinates (see, e.g., CSR from previous sheets).

→ **But:** Direct methods may still produce *dense* ( $\neq$  sparse) byproducts.

```
1  A = genRandomSymSparse(D)
2  # NONZEROES OF A
3  nonzeroesA += [np.count_nonzero(A)]
4  # NONZEROS OF LU
5  lu, piv = linalg.lu_factor(A)
6  nonzeroesLU += [np.count_nonzero(lu)]
7  # NONZEROS OF QR
8  Q, R = linalg.qr(A)
9  nonzeroesQR += [(np.count_nonzero(Q)+np.count_nonzero(R))]
```

10

```
11 # For sym. matrices we only need to store lower or upper triangle
12 bench = int(D*(D+1)*0.5)
13 # Average
14 nonzeroesA = np.round((np.array(nonzeroesA, dtype=float).sum()/N)/bench*100,2)
15 nonzeroesLU = np.round((np.array(nonzeroesLU, dtype=float).sum()/N)/bench*100,2)
16 nonzeroesQR = np.round((np.array(nonzeroesQR, dtype=float).sum()/N)/bench*100,2)
```



**Key idea:** We do not need the full matrix, but only some matrix-vector product  $x \mapsto S_A x$

**Approach:** Using this product we define a sequence  $\{x^0, x^1, x^2, \dots\} \subset \mathbb{R}^n$  such that

- $x^k \rightarrow x^* \in \mathbb{R}^n$  for  $k \rightarrow \infty$
- $x^*$  solves the linear equation  $Ax = b$
- $x^{k+1}$  is a better approximation to  $x^*$  than  $x^k$

Standard classes of such iterations:

- (1) **Linear iterations (splitting methods)**  $\leftarrow$
- (2) **Krylov subspace methods**

**Common advantages:**

- Even a few iteration steps may yield good results in stark contrast to LU (Gaussian Elimination), which has to be performed to the bitter end.
- The matrix  $A$  or its decomposition does not need to be stored! Only a **matrix-vector product**  $S_A x$  has to be provided.
- In general, the overall computational complexity is lower than with *direct* methods (LU, QR,...).

$\rightarrow$  Therefore: Iterative (*=indirect*) methods are to be preferred, when the matrix is large (and sparse).

C.F. Gauß in a letter to Gerling from 1823

(<https://gdz.sub.uni-goettingen.de/id/PPN23601515X?tify>)

Fast jeden Abend mache ich eine neue Auflage des Tableaus, wo immer leicht nachzuhelfen ist. Bei der Einförmigkeit des Messungsgeschäfts gibt dies immer eine angenehme Unterhaltung; man sieht dann auch immer gleich, ob etwas zweifelhaftes eingeschlichen ist, was noch wünschenswerth bleibt, etc. Ich empfehle Ihnen diesen Modus zur Nachahmung. Schwerlich werden Sie je wieder direct eliminiren, wenigstens nicht, wenn Sie mehr als 2 Unbekannte haben. Das indirecte Verfahren lässt sich halb im Schläfe ausführen, oder man kann während desselben an andere Dinge denken.

.....



- He already mentions an iterative method which was later coined *Gauß-Seidel method*.
- In general, one may need a lot of iteration steps but the aim is to keep the iteration instruction simple and fast.

### 5.1.2 A General Framework: Linear Fixed Point Iteration

Linear iterations are of the form

$$x^{k+1} = Mx^k + Nb$$

with  $M, N \in \mathbb{R}^{n \times n}$ . The matrix  $M$  is called the **iteration matrix** and motivates the adjective “linear”.

We first derive a general convergence result and then relate  $M$  and  $N$  to the system  $Ax = b$ .

#### Convergence analysis

**Definition 5.1 (Spectral radius)** Let  $M \in \mathbb{R}^{n \times n}$ . Then the largest eigenvalue of  $M$  in magnitude is called the **spectral radius of  $M$**  and denoted by  $\rho(M)$ , more precisely

$$\rho(M) := \max\{|\lambda_1|, \dots, |\lambda_n|\}.$$

**Theorem 5.2 (Fixed point iteration)** Let  $M, N \in \mathbb{R}^{n \times n}$  and  $b \in \mathbb{R}^n$ . If  $\rho(M) < 1$ , then the sequence

$$x^{k+1} = Mx^k + Nb$$

converges for any starting point  $x^0$  and its limit  $x^* \in \mathbb{R}^n$  is a fixed point of the affine linear function  $x \mapsto Mx + Nb$ , i.e.,

$$x^* = Mx^* + Nb.$$

## Proof for the special case that $M$ is symmetric:

Let us consider the eigenvalues of  $M$

$$|\lambda_n| \leq |\lambda_{n-1}| \leq \dots \leq |\lambda_1| =: \rho(M) < 1.$$

Since  $M$  is assumed to be symmetric, we can apply the theorem on the eigendecomposition of symmetric matrices and find  $Q \in \mathbb{R}^{n \times n}$  orthogonal,  $\Lambda \in \mathbb{R}^{n \times n}$  diagonal, so that

$$M = Q\Lambda Q^T, \text{ where } \Lambda = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}.$$

We will exploit two properties that we can conclude from  $\rho(M) < 1$ . First, we observe that the powers of  $M$  vanish. More precisely,

$$M^2 = M \cdot M = Q\Lambda \underbrace{Q^T Q}_{=I} \Lambda Q^T = Q\Lambda^2 Q^T$$

$\vdots$

$$M^k = Q\Lambda^k Q^T = Q \begin{pmatrix} \lambda_1^k & & 0 \\ & \ddots & \\ 0 & & \lambda_n^k \end{pmatrix} Q^T \xrightarrow[\rho(M) < 1]{k \rightarrow \infty} 0$$

Secondly, we note that the so-called Neumann series converges:  $\sum_{j=0}^{\infty} M^j = (I - M)^{-1}$  (general result, no proof here).

Now we are in the position to proof the convergence result. Therefore let  $x^0 \in \mathbb{R}^n$ , then

$$\begin{aligned} x^{k+1} &= M \cdot \underbrace{x^k}_{= Mx^{k-1} + Nb} + Nb = \underbrace{M^k x^0}_{\xrightarrow{k \rightarrow \infty} 0} + \underbrace{\left( \sum_{j=0}^k M^j \right)}_{\xrightarrow{k \rightarrow \infty} (I-M)^{-1}} \cdot Nb \xrightarrow{k \rightarrow \infty} (I - M)^{-1} Nb =: x^* \Leftrightarrow Nb = (I - M)x^* \end{aligned}$$

Which is equivalent to  $x^* = Mx^* + Nb$ . Thus, the limit  $x^*$  is a fixed point.

### 5.1.3 Splitting Methods

We now apply Theorem 5.2 to the linear system  $Ax = b$  by reformulating it as a fixed point problem. We also explain the idea of *preconditioning*.

**General scheme:**

$$x^{k+1} = Mx^k + Nb$$

#### First Approach

We rewrite the original system into a fixed point problem:

$$Ax = b \Leftrightarrow (x - x) + Ax = b \Leftrightarrow x = (I - A)x + Ib.$$

Then, if  $\rho(I - A) < 1$ , we can conclude by Theorem 5.2 that the fixed point iteration converges

$$x^{k+1} = (I - A)x^k + b \xrightarrow{k \rightarrow \infty} x^*, \text{ with } Ax^* = b.$$

#### Preconditioning

We see that the success of this procedure relies on the condition  $\rho(I - A) < 1$ . Unfortunately this is not generally the case, even not for invertible matrices (e.g.,  $A = 2I \in \text{GL}_n(\mathbb{R})$ , then  $\rho(I - A) = 1$ ). An idea to increase our chances is to first *precondition* the original system: Let  $N \in \mathbb{R}^{n \times n}$  be invertible, then

$$Ax = b \Leftrightarrow \underbrace{NA}_{=: \tilde{A}} x = \underbrace{Nb}_{=: \tilde{b}} \Leftrightarrow \tilde{A}x = \tilde{b} \Leftrightarrow x = (I - \tilde{A})x^k + \tilde{b} = (I - NA)x + Nb.$$

Now for this preconditioned fixed point problem: If  $\rho(I - NA) < 1$ , we can conclude by Theorem 5.2 that the fixed point iteration converges

$$x^{k+1} = (I - NA)x^k + Nb \xrightarrow{k \rightarrow \infty} x^*, \text{ with } Ax^* = b.$$



## What is a good preconditioner?

The optimal choice would be  $N = A^{-1}$ , because in this case we obtain the solution after one step:

$$x^1 = (I - NA)x^0 + Nb = A^{-1}b.$$

Access to the inverse would make our whole endeavor in this section dispensable. However, it serves as a compass to construct good preconditioners, which are mappings  $x \mapsto Nx$  for which  $N \approx A^{-1}$ .

The advantage and aim of preconditioners are not only to higher our chances for convergence, but also to accelerate the convergence. In fact, by using a good preconditioner  $N \approx A^{-1}$ , we cannot only make sure that the convergence criteria

$$\rho(I - \tilde{A}) = \rho(I - NA) \ll 1$$

is satisfied, but we can also show that the smaller  $\rho(I - NA)$ , the faster is the convergence.

## How to construct a preconditioner?

- Very simple preconditioning: **“Relaxation”**

Often it already suffices to make the steps small enough. This idea corresponds to the preconditioner

$$N := \theta I, \text{ with sufficiently small } 0 < \theta \begin{cases} = 1 : & \text{no relaxation} \\ < 1 : & \text{relaxation} \\ > 1 : & \text{over-relaxation} \end{cases}$$

The iteration matrix is then given by

$$M = I - \theta A.$$

- Next idea: **“Splitting based”**

Recall that the best preconditioner is  $N = A^{-1}$ , which is not accessible. However, what about inverting just a (significant) part of  $A$ ? In order to find such parts which are easy to invert, we consider the splitting of the matrix in the form:

The diagram shows the matrix splitting  $A = L + D + U$  using a grid-based matrix representation. The matrix  $A$  is split into three parts:  $L$  (lower triangular part, highlighted in green),  $D$  (diagonal part, highlighted in blue), and  $U$  (upper triangular part, highlighted in red). The equation is written as  $A = L + D + U$  with the corresponding matrix diagrams.

For example,  $D$  and  $(L + D)$  are diagonal and lower triangular, respectively. Thus they are easy to invert (i.e., to solve for) and we may obtain  $A^{-1} \approx D^{-1}$  or  $A^{-1} \approx (L + D)^{-1}$ , respectively.

## General scheme:

$$x^{k+1} = (I - NA)x^k + Nb = x^k - N(Ax^k - b)$$

| Preconditioner              | Iteration Matrix  | Iteration Instruction  | Method Name                             |
|-----------------------------|---|--|---|
| $N$                         | $M = I - NA$  | $x^{k+1} = x^k - N(Ax^k - b)$  |   |
| $\theta I$                  | $M_{Rich} = I - \theta A$   | $x^{k+1} = x^k - \theta(Ax^k - b)$   | (relax.) <b>Richardson</b>              |
| $\theta D^{-1}$             | $M_{Jac} = I - \theta D^{-1}A$<br>$= (1 - \theta)I - \theta D^{-1}(L + U)$                      | $x^{k+1} = x^k - \theta D^{-1}(Ax^k - b)$<br><i>Element-based:</i><br>$x_i^{k+1} = (1 - \theta)x_i^k + \frac{\theta}{a_{ii}} \left( b_i - \sum_{j \neq i} a_{ij} x_j^{k+1} \right)$  | (weighted) <b>Jacobi</b>                |
| $(L + D)^{-1}$              | $M_{GS} = I - (L + D)^{-1}A$<br>$= (L + D)^{-1}U$   | $x^{k+1} = x^k - (L + D)^{-1}(Ax^k - b)$<br><i>Element-based:</i><br>$x_i^{k+1} = \frac{1}{a_{ii}} \left( b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{k+1} - \sum_{j=i+1}^n a_{ij} x_j^k \right)$                                       | <b>Gauß-Seidel</b>                      |
| $\theta(\theta L + D)^{-1}$ | $M_{SOR} = I - \theta(\theta L + D)^{-1}A$<br>$= (\theta L + D)^{-1}((1 - \theta)D - \theta U)$ | $x^{k+1} = x^k - \theta(\theta L + D)^{-1}(Ax^k - b)$<br><i>Element-based:</i><br>$x_i^{k+1} = (1 - \theta)x_i^k + \frac{\theta}{a_{ii}} \left( b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{k+1} - \sum_{j=i+1}^n a_{ij} x_j^k \right)$ | <b>Successive Over-Relaxation (SOR)</b> |

## What remains: When does $\rho(M) < 1$ hold?

Idea: Derive some (possibly easy-to-compute) conditions which are sufficient for  $\rho(M) < 1$   
→ As usual, we often need to assume some properties for  $A$  (e.g., symmetry).

| Method  | Condition  |
|---|--|
| <b>(relax.) Richardson</b><br>$M_{Rich} = I - \theta A$                               | <u>With relaxation (<math>\theta \neq 1</math>):</u><br>If $A$ is symmetric and positive definite (spd), then:<br>$\rho(M_{Rich}) < 1 \Leftrightarrow 0 < \theta < \frac{2}{\lambda_{\max}(A)}$  |
| <b>(weighted) Jacobi</b><br>$M_{Jac} = I - \theta D^{-1}A$                            | <u>No relaxation (<math>\theta = 1</math>):</u><br>If $A$ is strictly diagonally dominant (i.e., $ a_{ii}  > \sum_{i \neq j}  a_{ij} $ ), then $\rho(M_{Jac}) < 1$<br><br><u>With relaxation (<math>\theta \neq 1</math>):</u><br>If $A$ is spd, then:<br>$\rho(M_{Jac}) < 1 \Leftrightarrow 0 < \theta < \frac{2}{\lambda_{\max}(D^{-1}A)}$ |
| <b>Gauß-Seidel</b><br>$M_{GS} = I - (L + D)^{-1}A$                                    | <ul style="list-style-type: none"><li>• If <math>A</math> is strictly diagonally dominant, then <math>\rho(M_{GS}) &lt; 1</math></li><li>• If <math>A</math> is spd, then <math>\rho(M_{GS}) &lt; 1</math></li></ul>   |
| <b>Successive Over-Relaxation (SOR)</b><br>$M_{SOR} = I - \theta(\theta L + D)^{-1}A$ | If $A$ is spd, then $\rho(M_{SOR}) < 1$ for $0 < \theta < 2$   |

**Remark:** For many matrix classes (e.g., symmetric), we find

$$\rho(M_{GS}) \leq \rho(M_{Jac}) \leq \rho(M_{Rich}) \quad (\theta = 1).$$

## Example

Let us consider the symmetric tridiagonal matrix

$$A = \begin{pmatrix} 3 & -1 & & 0 \\ -1 & \ddots & \ddots & \\ & \ddots & \ddots & -1 \\ 0 & & -1 & 3 \end{pmatrix} = \underbrace{\begin{pmatrix} 0 & & & 0 \\ -1 & \ddots & & \\ & \ddots & \ddots & \\ 0 & & -1 & 0 \end{pmatrix}}_L + \underbrace{\begin{pmatrix} 3 & & & 0 \\ & \ddots & & \\ 0 & & 3 & \\ & & & \end{pmatrix}}_D + \underbrace{\begin{pmatrix} 0 & -1 & & 0 \\ & \ddots & \ddots & \\ & & \ddots & -1 \\ 0 & & & 0 \end{pmatrix}}_U$$

For large  $n$  this matrix is extremely sparse:

$$\frac{\# \text{ nonzero entries}}{n^2} = \frac{n + 2(n-1)}{n^2} \xrightarrow{n \rightarrow \infty} 0.$$

Based on the convergence criteria from the table above, we now check which methods would converge here. Since  $|3| > |-1| + |-1|$  it is easily verified that  $A$  is strictly diagonally dominant. Therefore Jacobi and Gauß-Seidel converge.

What about the Richardson method? For simplicity we consider  $n = 3$  and compute the eigenvalues of  $A$  and find

$$\begin{aligned} \det(A - \lambda I) &= \det \begin{pmatrix} 3-\lambda & -1 & 0 \\ -1 & 3-\lambda & -1 \\ 0 & -1 & 3-\lambda \end{pmatrix} \\ &\stackrel{\text{Sarrus}}{=} -(3-\lambda)^3 + (3-\lambda) + (3-\lambda) \\ &= (3-\lambda) \cdot (2 - (3-\lambda)^2) \stackrel{!}{=} 0 \Leftrightarrow \lambda = 3, (3-\lambda) = \pm\sqrt{2} \Leftrightarrow \lambda = 3, \lambda = 3 \pm \sqrt{2} \end{aligned}$$

Since all eigenvalues of this symmetric matrix are positive, we have that  $A$  is positive definite. Noting that  $\lambda_{\max}(A) = 3 + \sqrt{2}$ , we apply the above convergence criteria and find that the relaxed Richardson method converges, if the the relaxation parameter is chosen from the interval

$$0 < \theta < \frac{2}{\lambda_{\max}(A)} = \frac{2}{3 + \sqrt{2}} \approx 0.45 < 1.$$

## Final Remarks

- Using the Richardson iteration we only need the evaluation of the matrix vector product  $x \mapsto Ax$  to solve the system  $Ax = b$ .
- All iterations of the form  $x^{k+1} = x^k - N(Ax^k - b)$  can be seen as preconditioned Richardson iterations. Next semester you will learn a strong correspondence between the Richardson method and the gradient method, as well as preconditioned Richardson method and Newton-type methods to minimize functions of the form  $f(x) := \frac{1}{2}x^\top Ax - b^\top x$ .
- By looking at the element-wise formulas for the Jacobi (with  $\theta = 1$ ) and Gauß-Seidel method (see the orange formulas in the table above) we can understand these methods as alternating methods. Consider for example the simple case  $n = 2$  and write out these formulas. Then you will see that we alternately compute the components  $x_1$  and  $x_2$ .
- Similar to a block LU factorization, we can consider blocks  $A_{ij}$  (= matrices) instead of numbers  $a_{ij}$  in the element-wise formulas of Jacobi and Gauß-Seidel. Thereby, one obtains the block Jacobi and block Gauß-Seidel method. These are highly related to so called additive and multiplicative Schwarz methods, respectively.
- In practice, splitting methods (or more precisely the preconditioner  $N$ ) are mainly used as preconditioners for Krylov subspace methods, which we will address in the next section.

## 5.2 Krylov Subspace Methods

### 5.2.1 Krylov Subspaces

Situation:

- $A \in GL_n(\mathbb{R})$  (invertible,  $n$  typically large, but sparse)
- $b \in \mathbb{R}^n$

We want to solve  $Ax = b$ , but we cannot work with the matrix as a (dense) “array”. Instead we only have access to the mapping

$$x \mapsto Ax \quad (\text{matrix-vector product}).$$

Then given  $b \in \mathbb{R}^n$ , we can produce the vectors

$$b, Ab, A^2b, \dots, A^k b.$$

There is not much more to consider. Let us collect all linear combinations of these vectors and give it a name:

**Definition 5.3 (Krylov\* subspaces)** Let  $A \in \mathbb{R}^{n \times n}$  and  $b \in \mathbb{R}^n$ , then the set

$$K_r(A, b) := \text{span}\{b, Ab, A^2b, \dots, A^{r-1}b\}$$

is called **Krylov Subspace of order  $r \geq 1$  generated by  $A$  and  $b$ .**

In order to develop an iterative scheme based on this definition, we look deeper into the Krylov subspaces and first collect some insightful observations.

---

\*named after the Russian engineer Alexei Krylov who developed the idea in a paper published around 1931.

## Remarks

i) If  $b = 0$ , then  $A^k b = 0$  for all  $A \in \mathbb{R}^{n \times n}$  and all  $k \in \mathbb{N}$ , so that

$$K_r(A, 0) = \{0\} \quad \forall A \in \mathbb{R}^{n \times n}, r \geq 1.$$

ii) If  $b \neq 0$  and  $A = I$ , then  $A^k b = b$  for all  $k \in \mathbb{N}$ , so that

$$K_r(I, b) = \text{span}\{b\} \quad \forall b \neq 0, r \geq 1.$$

iii) If  $b$  is an eigenvector of  $A$  to the eigenvalue  $\lambda \in \sigma(A)$ , then  $A^k b = \lambda^k b$ , so that  $A^k b$  and  $b$  are linearly dependent, implying

$$K_r(A, b) = \text{span}\{b\} \quad \forall \text{ eigenvectors } b \text{ of } A.$$

iv) Insight from the power method:

Recall: Let  $|\lambda_1| > |\lambda_2| \geq \dots |\lambda_n|$  and  $b^T v_1 \neq 0$ , then  $\frac{A^k b}{\|A^k b\|} \rightarrow v_1$ .

Thus, for large  $k$  we have that the  $A^k b$  point into a “similar direction” – more precisely, into the direction of  $v_1$ . With other words the  $A^k b$  become more and more linearly dependent.

v) Dimension of Krylov subspaces:

- Since  $(K_r(A, b)) \subset \mathbb{R}^n$  is spanned by  $r$  vectors, we clearly have  $\dim(K_r(A, b)) \leq \min(r, n)$ .
- For  $A \in GL_n(\mathbb{R})$  one can show that

$$b, Ab, \dots, A^{r-1}b \text{ are independent} \quad \forall r \leq r_{\max},$$

where  $r_{\max}$  is the maximal dimension a Krylov subspace generated by  $A$  and  $b$  can have, i.e.,  
 $r_{\max} := \max_{s \leq n} (\dim K_s(A, b)).$



vi) Next we state the crucial result, which forms the basis for the development of Krylov subspace methods:

**Lemma 5.4** *Let  $A \in GL_n(\mathbb{R})$  and  $b \in \mathbb{R}^n$ , then there exists an order  $r \leq n$  so that for the solution  $x^*$  of  $Ax = b$  we have*

$$x^* = A^{-1}b \in K_r(A, b).$$

*In particular, we find coefficients  $\beta_0, \dots, \beta_{r-1}$ , such that*

$$x^* = \sum_{j=0}^{r-1} \beta_j A^j b.$$

*Proof.* Consider the  $(n+1)$  vectors  $b, Ab, \dots, A^{n-1}b, A^n b$ , which are necessarily dependent in  $\mathbb{R}^n$ . Thus due to the dependence we find  $\alpha_0, \alpha_1, \dots, \alpha_n \in \mathbb{R}$ , which are not all zero, so that

$$\begin{aligned} 0 &= \sum_{j=0}^n \alpha_j A^j b = \alpha_0 A^0 b + \alpha_1 A^1 b + \dots + \alpha_n A^n b \\ &\stackrel{(k:=\text{smallest index with } \alpha_k \neq 0)}{=} \underbrace{\alpha_k}_{\neq 0} A^k b + \alpha_{k+1} A^{k+1} b + \dots + \alpha_n A^n b \\ &\stackrel{(A^{-(k+1)} \cdot |)}{\Leftrightarrow} 0 = \alpha_k \underbrace{A^{-1} b}_{=x^*} + \alpha_{k+1} A^0 b + \dots + \alpha_n A^{n-k-1} b \\ &\stackrel{(\alpha_k \neq 0)}{\Leftrightarrow} x^* = \frac{\alpha_{k+1}}{-\alpha_k} A^0 b + \frac{\alpha_{k+2}}{-\alpha_k} A^1 b + \dots + \frac{\alpha_n}{-\alpha_k} A^{n-k-1} b = \sum_{j=0}^{r-1} \beta_j A^j b, \end{aligned}$$

with  $\beta_j := \frac{\alpha_{k+1+j}}{-\alpha_k}$  and  $r := n - k$ . □

## Idea of Krylov Subspace Methods

- In words, the latter lemma states that the sought-after solution can be found in a Krylov subspace of some order  $r \leq n$ . However we neither generally know  $r$  nor the linear coefficients  $\beta_j$ .
- This fact leads to the following general iterative scheme of so-called Krylov subspace methods:  
In each iteration step  $1 \leq r \leq n$ , minimize the residual  $\|Ax - b\|_2$  over all  $x \in K_r(A, b)$ , i.e.,

$$x_r := \operatorname{argmin}_{x \in K_r(A, b) \subset \mathbb{R}^n} \|Ax - b\|_2^2.$$

Thereby we generate a sequence  $x_1, x_2, \dots, x_n$ .

- Since  $r \leq n$ , we know that in exact arithmetic (neglecting presence of rounding errors) at latest the step  $r = n$  ( $x_n$ ) gives us the exact solution, so that one may call these methods “direct”. However  $n$  is typically very large ( $\geq 10^5$ ) and in many examples it turns out that a few steps may already give very good approximations.
- Such methods to solve  $Ax = b$  have been developed in the early 1950s; the most famous are:
  - **Conjugate Gradient (CG)** (Hestenes, Stiefel; 1952)  
for  $A$  symmetric and positive definite
  - **MINimal RESidual (MINRES)** (Paige, Saunders; 1975)  
for  $A$  symmetric and invertible
  - **General Minimal RESidual (GMRES)** (Schultz, Saad; 1986)  
for  $A \in GL_n(\mathbb{R})$

→ In this chapter we will derive the GMRES method.

## Comparison to other methods

| Least squares   | Krylov subspace   | Splitting  |
|---|---|--|
| $\hat{x} := \operatorname{argmin}_{x \in \mathbb{R}^n} \ Ax - b\ _2^2$<br>$A\hat{x} = \operatorname{argmin}_{w \in \operatorname{Im}(A)} \ w - b\ _2^2$ | $x_r := \operatorname{argmin}_{x \in K_r(A, b)} \ Ax - b\ _2^2$<br>$Ax_r = \operatorname{argmin}_{w \in AK_r(A, b)} \ w - b\ _2^2$    | $x_r := Mx_{r-1} + Nb$   |
| $\rightsquigarrow$ Projection of $b$ onto<br>$\operatorname{Im}(A) = A\mathbb{R}^n = \mathbb{R}^n$  | $\rightsquigarrow$ Projection of $b$ onto<br>$AK_r(A, b) = \operatorname{span}\{Ab, A^2b, \dots, A^rb\}$<br>$\subseteq K_{r+1}(A, b)$ |  |
| Yields exact solution, i.e.,<br>$\hat{x} = A^{-1}b$   | In theory: Yields a finite sequence<br>with $x_n = A^{-1}b$   | Typically yields an infinite sequence<br>with $x^r \xrightarrow{(k \rightarrow \infty)} A^{-1}b$ |

As a preparation for what follows, we will first derive an iterative method to find orthonormal bases for Krylov subspaces:

### 5.2.2 The Arnoldi Iteration

- **Situation:** Let us consider the  $r$ -th Krylov subspace

$$K_r(A, b) = \text{span}\{A^0b, A^1b, \dots, A^{r-1}b\}$$

with  $r \leq r_{\max} = \max_{s \leq n} (\dim(K_s(A, b)))$ , so that all  $A^{j-1}b$  are independent and  $\dim(K_r(A, b)) = r$ .

- **Aim:** Find an orthonormal basis  $\{q_1, \dots, q_r\}$  of  $K_r(A, b)$ .
- **Idea:** Apply the Gram–Schmidt orthogonalization process to the linearly independent vectors  $c_j := A^{j-1}b$ ,  $1 \leq j \leq r$ .
- **Recall Gram–Schmidt:** Let  $c_1, \dots, c_r \in \mathbb{R}^n$ , where  $r \leq n$ , be linearly independent vectors. Then an orthonormal basis  $\{q_1, \dots, q_r\}$  of  $\text{span}(c_1, \dots, c_r) = \text{Im}(C)$  can be found by the following iterative scheme:

$$q_1 := \frac{c_1}{\|c_1\|}$$

$$\text{“subtracting projections:” } \hat{q}_j := c_j - \sum_{\ell=1}^{j-1} q_\ell^\top c_j \cdot q_\ell, \quad r_{\ell j} := q_\ell^\top c_j$$

$$\text{“normalization:” } q_j := \frac{\hat{q}_j}{\|\hat{q}_j\|_2}, \quad r_{jj} := \|\hat{q}_j\|_2$$

The matrix perspective: Putting the vectors  $q_j$  and  $r_j$  (which are computed step by step) into matrices, say  $Q$  and  $R$ , then we obtain the (reduced) QR-decomposition of  $C$ , i.e., a matrix  $Q = [q_1, \dots, q_r] \in \mathbb{R}^{n \times r}$  with orthonormal columns and an upper triangular matrix  $R \in \mathbb{R}^{r \times r}$  with  $r_{ii} \neq 0$ , so that  $C = QR$ , which implies (also see Lemma 1.47)  $\text{Im}(C) = \text{Im}(QR) = \text{Im}(Q) = \text{span}(q_1, \dots, q_r)$ .

Now we consider the specific choice  $c_j := A^{j-1}b$ . By inserting these  $c_j$  into the above scheme, we obtain the so called **Arnoldi iteration**:

$$q_1 := \frac{b}{\|b\|_2} \quad (\rightarrow \text{ orthonormal basis for } K_1(A, b) = \text{span}(b), b \neq 0)$$

**for**  $j = 2, \dots, r$  :

$$\text{"subtracting projections:" } \hat{q}_j := Aq_{j-1} - \sum_{\ell=1}^{j-1} q_\ell^\top (Aq_{j-1}) \cdot q_\ell, \quad h_{\ell,j-1} := q_\ell^\top (Aq_{j-1}) \quad (*)$$

$$\text{"normalization:" } q_j := \frac{\hat{q}_j}{\|\hat{q}_j\|_2}, \quad h_{j,j-1} := \|\hat{q}_j\| \quad (**)$$

All in all, the Arnoldi process yields an orthonormal basis for  $K_r(A, b)$ , so that

$$K_r(A, b) = \text{span}(b, Ab, \dots, A^{r-1}b) = \text{span}(q_1, \dots, q_r)$$

or with  $Q_r := [q_1, \dots, q_r] \in \mathbb{R}^{n \times r}$  in matrix form

$$K_r(A, b) = \text{Im}(Q_r).$$

Remarks:

- Rearranging (\*) and (\*\*) easily gives

$$Aq_{j-1} = \|\hat{q}_j\|_2 \cdot q_j + \sum_{\ell=1}^{j-1} q_\ell^\top (Aq_{j-1}) \cdot q_\ell = \sum_{\ell=1}^j h_{\ell,j-1} q_\ell.$$

With other words,  $Aq_{j-1}$  is a linear combination of  $q_1, \dots, q_j$ .

- Also observe for  $j \leq r$ : The first  $q_1, \dots, q_{j-1}$  are a basis for  $K_{j-1}(A, b)$ . Thus, assumed these are given, then in order to find an orthonormal basis for  $K_j(A, b)$  we just need to compute one more vector, namely  $q_j$ .

### 5.2.3 GMRES

(with Arnoldi and  $x_0 = 0$ )

Let us recall the general idea of Krylov subspace methods: In each iteration step we compute

$$x_r := \operatorname{argmin}_{x \in K_r(A, b)} \|Ax - b\|_2^2.$$

Now we approach this minimization problem in a specific procedure resulting in the GMRES method:

We first find an orthonormal basis  $\{q_1, \dots, q_r\}$  of  $K_r(A, b)$  so that  $K_r(A, b) = \operatorname{span}(q_1, \dots, q_r) = \operatorname{Im}(Q_r)$ , where  $Q_r := [q_1, \dots, q_r] \in \mathbb{R}^{n \times r}$ . Then, since by definition  $x_r \in K_r(A, b)$ , the minimization problem can be rephrased as

$$x_r = \operatorname{argmin}_{x \in \operatorname{Im}(Q_r)} \|Ax - b\|_2^2 \stackrel{(x_r = Q_r c_r)}{=} Q_r \cdot \underbrace{(\operatorname{argmin}_{c \in \mathbb{R}^r} \|AQ_r c - b\|_2^2)}_{=: c_r, \text{ standard least squares problem}}.$$

The least squares problem can then be solved with the help of the QR-decomposition of the design matrix  $AQ_r$ , say  $\tilde{Q}_r \tilde{R}_r := AQ_r$ , so that the corresponding normal equation reads as

$$\tilde{R}_r c_r = \tilde{Q}_r^T b \quad (\text{when is } \tilde{R}_r \text{ invertible?})$$

Thus, GMRES boils down to the three steps. For  $1 \leq r \leq n$ , do

- Step 1: Find an orthonormal basis for  $K_r(A, b)$ .
- Step 2: Find  $\tilde{Q}_r \tilde{R}_r := AQ_r$ .
- Step 3: Solve  $\tilde{R}_r c_r = \tilde{Q}_r^T b$  and set  $x_r := Q_r c_r$ .

Crucial: We can iteratively compute step 1 and step 2, i.e., use  $Q_{r-1}$ ,  $\tilde{R}_{r-1}$ ,  $\tilde{Q}_{r-1}$  to obtain  $Q_r$ ,  $\tilde{R}_r$ ,  $\tilde{Q}_r$ !

## GMRES – Step 1 “Find an orthonormal basis for $K_r(A, b)$ ”

Given  $Q_{r-1} = (q_1, \dots, q_{r-1}) \in \mathbb{R}^{n \times (r-1)}$  with  $\text{Im}(Q_{r-1}) = K_{r-1}(A, b)$  from the previous iteration step, let us compute

$$q_r := \text{Arnoldi\_step}(q_1, \dots, q_{r-1}; Aq_{r-1}).$$

(“orthogonalizing  $Aq_{r-1}$  against all  $q_1, \dots, q_{r-1}$  by subtracting projections and normalization”)

We recall the  $r$ -th Arnoldi iteration step in an algorithmic fashion:

Given  $q_1(= \frac{b}{\|b\|}), \dots, q_{r-1}$ , then  $q_r$  is computed by:

```
 $q_r := \text{Arnoldi\_step}(q_1, \dots, q_{r-1}; v = Aq_{r-1}):$   
  for  $\ell = 1, \dots, r-1$  do  
     $h_{\ell, r-1} = q_\ell^T v$   
     $v = v - h_{\ell, r-1} q_\ell$     (subtracting projections)  
  end  
   $h_{r, r-1} = \|v\|$   
  if  $h_{r, r-1} \neq 0$  then  
     $q_r = \frac{v}{h_{r, r-1}}$     (normalization)  
    return  $q_r$   
  end  
  return  $v$ 
```

In matrix notation this can be summed up as follows:

$$A \cdot \underbrace{\begin{pmatrix} | & & | \\ q_1 & \cdots & q_{r-1} \\ | & & | \end{pmatrix}}_{=: Q_{r-1} \in \mathbb{R}^{n \times (r-1)}} = \underbrace{\begin{pmatrix} | & & | \\ Aq_1 & \cdots & Aq_{r-1} \\ | & & | \end{pmatrix}}_{\text{Arnoldi: } Aq_j = \sum_{\ell=1}^{j+1} h_{\ell,j} q_\ell} = \underbrace{\begin{pmatrix} | & & | & | \\ q_1 & \cdots & q_{r-1} & q_r \\ | & & | & | \end{pmatrix}}_{=: Q_r \in \mathbb{R}^{n \times r}} \underbrace{\begin{pmatrix} h_{1,1} & h_{1,2} & \cdots & & \\ h_{2,1} & h_{2,2} & & & \vdots \\ 0 & h_{3,2} & & & \\ \vdots & 0 & \ddots & h_{r-1,r-2} & h_{r-1,r-1} \\ 0 & 0 & 0 & h_{r,r-2} & h_{r,r-1} \end{pmatrix}}_{=: H_{r,r-1} \in \mathbb{R}^{r \times (r-1)}}$$

$$\stackrel{Q_{r-1}^T}{\Leftrightarrow} Q_{r-1}^T A Q_{r-1} = Q_{r-1}^T Q_r H_{r,r-1} = \begin{pmatrix} & | & 0 \\ I_{(r-1) \times (r-1)} & | & \vdots \\ & | & 0 \end{pmatrix} \begin{pmatrix} H_{r-1} \\ - & - & - \\ \text{row } r \text{ of } H_{r,r-1} \end{pmatrix} = H_{r-1}$$

- We observe that  $H_{r-1} \in \mathbb{R}^{(r-1) \times (r-1)}$  has only one subdiagonal. Such matrices are called (upper) **Hessenberg matrices**.
- In particular, for  $r = n + 1$ , we find

$$Q_n^T A Q_n = H_n \in \mathbb{R}^{n \times n}.$$

With other words, with the help of the Arnoldi iteration we can find in finitely many steps (at most  $n$  steps) an upper Hessenberg  $H_n$ , which is orthogonally similar to  $A$  and thus has the same eigenvalues. That is why one can tailor the QR-algorithm (which was an algorithm to compute eigenvalues) to matrices of Hessenberg structure. The Arnoldi iteration is therefore also considered an eigenvalue algorithm. Furthermore, if  $A$  is symmetric, so is  $H_n$ , which then becomes a tridiagonal matrix.



## GMRES – Step 2 “Compute $QR$ -decomposition $\tilde{Q}_r \tilde{R}_r := A Q_r$ ”

Again, we want to rely on the computations from the previous steps, i.e.,  $\tilde{Q}_{r-1} = [\tilde{q}_1, \dots, \tilde{q}_{r-1}] \in \mathbb{R}^{n \times (r-1)}$  and  $\tilde{R}_{r-1} = [\tilde{r}_1, \dots, \tilde{r}_{r-1}] \in \mathbb{R}^{(r-1) \times (r-1)}$  with  $\tilde{Q}_{r-1} \tilde{R}_{r-1} = A Q_{r-1}$ . Therefore, let us first observe that

$$\tilde{Q}_r \tilde{R}_r \stackrel{!}{=} A Q_r = A \cdot [Q_{r-1} | q_r] = [A Q_{r-1} | A q_r] = [\tilde{Q}_{r-1} \tilde{R}_{r-1} | A q_r] \in \mathbb{R}^{n \times r}.$$

Let us use one Arnoldi step to orthogonalize  $A q_r$  against  $\tilde{q}_1, \dots, \tilde{q}_{r-1}$  to obtain

$$\tilde{q}_r, \tilde{r}_r := \text{Arnoldi\_step}(\tilde{q}_1, \dots, \tilde{q}_{r-1}; A q_r),$$

where the vector  $\tilde{q}_r$  and the coefficients of  $\tilde{r}_r = (\tilde{r}_{1,r}, \dots, \tilde{r}_{r,r})^\top \in \mathbb{R}^r$  are computed via

$$\hat{q}_r := A q_r - \sum_{j=1}^{r-1} \tilde{r}_{j,r} \cdot \tilde{q}_j, \quad \tilde{r}_{j,r} := \tilde{q}_j^\top (A q_r), \quad \tilde{q}_r := \frac{\hat{q}_r}{\|\hat{q}_r\|_2}, \quad \tilde{r}_{rr} := \|\hat{q}_r\|_2,$$

from which we can conclude

$$A q_r = \tilde{r}_{rr} \tilde{q}_r + \sum_{j=1}^{r-1} \tilde{r}_{j,r} \tilde{q}_j = \sum_{j=1}^r \tilde{r}_{j,r} \tilde{q}_j. \quad (*)$$

Then let us define potential candidates for the sought-after  $QR$ -decomposition as follows:

$$\tilde{Q}_r := [\tilde{Q}_{r-1} | \tilde{q}_r], \quad \tilde{R}_r := \begin{pmatrix} \tilde{R}_{r-1} & | \\ - & - & \tilde{r}_r \\ 0 \cdots 0 & | \end{pmatrix} \in \mathbb{R}^{r \times r}.$$

By  $(*)$  we can write  $A q_r = \tilde{Q}_r \tilde{r}_r$  and indeed find, as desired,

$$A Q_r = (A Q_{r-1} | A q_r) = (\tilde{Q}_{r-1} \tilde{R}_{r-1} | \tilde{Q}_r \tilde{r}_r) = \tilde{Q}_r \tilde{R}_r.$$

### GMRES – Step 3 “Solve a triangular system”

The last step is easily performed by invoking a routine to solve the upper triangular system

$$\tilde{R}_r c_r = \tilde{Q}_r^T b$$

via backward substitution. Then we find our  $r$ -th iterate by setting

$$x_r := Q_r c_r.$$

## 5.2.4 Summary and final Remarks: GMRES with Arnoldi

i) We only need one matrix-vector product  $v \mapsto Av$  in each step! This product can be delivered to the solver GMRES as some sort of black box function.

ii) **An initial guess**  $x_0 \neq 0$ :

Consider

$$x_r := x_0 + p_r, \quad \text{for some } x_0 \neq 0,$$

then

$$Ax_r = b \Leftrightarrow A(x_0 + p_r) = b \Leftrightarrow Ap_r = b - Ax_0 = \hat{b}.$$

Then we could solve the auxiliary system to obtain  $p := \text{GMRES}(A, \hat{b})$  and set  $x := x_0 + p$ . The above algorithm can therefore easily be extended to allow for an initial guess  $x_0$  other than the zero vector. Also note that the Krylov subspaces are then generated based on the initial residual  $\hat{b} = b - Ax_0$ .

iii) **Preconditioning:**

Let us consider  $N \in \text{GL}_n(\mathbb{R})$ ,  $N \approx A^{-1}$ , where  $v \mapsto Nv$  is also easy to compute. Then

$$Ax = b \Leftrightarrow \underbrace{NA}_{=: \tilde{A}} x = \underbrace{Nb}_{=: \tilde{b}} \Leftrightarrow \tilde{A}x = \tilde{b}.$$

Then we could solve the equivalent system to obtain  $x := \text{GMRES}(\tilde{A}, \tilde{b})$ .

What is actually a good precondition? Quite a bit of research has been put into this question over the last decades. What one can say: Krylov subspace methods converge quickly if the eigenvalues appear in clusters away from zero.

#### iv) Restarted GMRES:

- We need to store  $Q_r := (q_1, \dots, q_r)$ ,  $\tilde{Q}_r := (\tilde{q}_1, \dots, \tilde{q}_r)$  and  $\tilde{R}_r := (\tilde{r}_1, \dots, \tilde{r}_r)$ . Thus, if the above algorithm runs many iterations, then  $Q_r, \tilde{Q}_r$  may become large (usually dense) matrices ( $n \gg 1$ ), then the advantage of the sparsity of  $A$  is lost.
- Idea: Perform only a fixed number of iterations, say  $m = 20 - 40$  and then restart GMRES with initial guess  $x_m$ . The resulting algorithm is sometimes coined GMRES(m). You will later learn about a similar idea in the context of the L-BFGS method.
- Remark: If  $A$  is symmetric and positive definite, then one can show that we do not need the whole history  $q_1, \dots, q_r$ . This is the power of the CG-method which you will investigate next semester.

#### v) Another perspective and variant: Let us consider

$$\min_{x_r \in K_r(A, b)} \|Ax_r - b\|_2^2 = \min_{c_r \in \mathbb{R}^r} \|AQ_r c_r - b\|_2^2 = \min_{c_r \in \mathbb{R}^r} \|Q_{r+1} H_{r+1, r} c_r - b\|_2^2$$

$$\stackrel{(\#)}{=} \min_{c_r \in \mathbb{R}^r} \underbrace{\left\| H_{r+1, r} c_r - \begin{pmatrix} \|b\| \\ 0 \\ \vdots \\ 0 \end{pmatrix} \right\|_2^2}_{\text{we can evaluate this residual without computing } x_r!} \quad (*)$$

$$(\#): q_1 := \frac{b}{\|b\|} \Rightarrow b = q_1 \|b\| \Rightarrow Q_{r+1}^T b = \begin{pmatrix} \|b\| \\ 0 \\ \vdots \\ 0 \end{pmatrix} \Rightarrow b = Q_{r+1} \begin{pmatrix} \|b\| \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

Let  $\hat{Q}_{r+1}\hat{R}_r := H_{r+1,r}$  be a  $QR$ -decomposition, then the normal equation associated to (\*) could be solved by:

$$\hat{R}_r c_r = \hat{Q}_{r+1}^T \begin{pmatrix} \|b\| \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

This results in a different variant of the GMRES and the following adaptations would apply:

- In step 2: We could derive such a  $QR$ -decomposition  $\hat{Q}_{r+1}\hat{R}_r := H_{r+1,r}$  from our arrays by setting

$$\hat{Q}_{r+1} := Q_{r+1}^T \tilde{Q}_r, \quad \hat{R}_r := \tilde{R}_r$$

because then

$$\tilde{Q}_r \tilde{R}_r \stackrel{!}{=} A Q_r = Q_{r+1} H_{r+1,r} \Leftrightarrow Q_{r+1}^T \tilde{Q}_r \tilde{R}_r = H_{r+1,r} \Leftrightarrow \hat{Q}_{r+1} \hat{R}_r = H_{r+1,r} \quad (*).$$

- In step 3: Our normal equation could then be written as

$$\tilde{R}_r c_r = \tilde{Q}_r^T b \Leftrightarrow \tilde{R}_r c_r = (Q_{r+1}^T \tilde{Q}_r)^T \begin{pmatrix} \|b\| \\ 0 \\ \vdots \\ 0 \end{pmatrix} \Leftrightarrow \hat{R}_r c_r = \hat{Q}_{r+1} \begin{pmatrix} \|b\| \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

All in all we obtain the following algorithm:

**INPUT:**  $A \in GL_n(\mathbb{R})$ ,  $b \in \mathbb{R}^n$

**OUTPUT:** approximation  $x_r \in K_r(A, b)$  to the exact solution  $A^{-1}b$

GMRES( $A, b, x_0 = 0$ , tol = 1e-6, maxiter=None,  $N = I$ ):

$b = b - Ax_0$  //account for initial guess

$A = NA$ ,  $b = Nb$  //account for preconditioner

    //Initialization:

$q_1 := \frac{b}{\|b\|_2}$ ,  $Q_1 := [q_1]$

$v := Aq_1$ ,  $\tilde{q}_1 := \frac{v}{\|v\|_2}$ ,  $\tilde{Q}_1 := [\tilde{q}_1]$ ,  $\tilde{R}_1 = [\|v\|_2]$

**for**  $r = 2, \dots, \min(n, \text{maxiter})$  **do**

        //STEP 1: use Arnoldi to find column  $q_r$  by orthogonalizing  $v$  against  $q_1, \dots, q_{r-1}$

$q_r, h_{r-1} := \text{Arnoldi\_step}(Q_{r-1}; v)$  //we don't need  $h_{r-1}$

$Q_r := [Q_{r-1}, q_r]$

$v := Aq_r$

        //STEP 2: use Arnoldi to find columns  $\tilde{q}_r, \tilde{r}_r$  by orthogonalizing  $v$  against  $\tilde{q}_1, \dots, \tilde{q}_{r-1}$

$\tilde{q}_r, \tilde{r}_r := \text{Arnoldi\_step}(\tilde{Q}_{r-1}; v)$

$\tilde{Q}_r := [\tilde{Q}_{r-1}, \tilde{q}_r]$ ,  $\tilde{R}_r := [\tilde{R}_{r-1}, \tilde{r}_r]$

        //STEP 3: solve auxiliary least squares problems to obtain coordinates

$c_r := \text{solve\_triangular}(\tilde{R}_r, \tilde{Q}_r^T b)$

$x_r := Q_r c_r$

        //Attention: Evaluate the original residual here:

**if**  $\|N^{-1}(Ax_r - b)\|_2 < \text{tol}$  **then**

            | break

**end**

**end**