

Chapter 6

Krylov Subspace Methods that are Based on a Projection of the Residual

Remark 6.1 *Idea.* The methods presented in this section determine the iterate $\mathbf{x}^{(k)}$ at the manifold $\mathbf{x}^{(0)} + K_k(\mathbf{r}^{(0)}, A)$ such that the corresponding residual $\mathbf{r}^{(k)}$ is orthogonal to $K_k(\mathbf{r}^{(0)}, A)$. That means, $\mathbf{r}^{(k)}$ is projected into the orthogonal complement $K_k(\mathbf{r}^{(0)}, A)^\perp$ of $K_k(\mathbf{r}^{(0)}, A)$. \square

6.1 General Matrices

Remark 6.2 *Full orthogonalization method.* Let $Q_k = \{\mathbf{q}_1, \dots, \mathbf{q}_k\}$ be an orthonormal basis of $K_k(\mathbf{r}^{(0)}, A)$ computed with Arnoldi's method. The identities (5.1) and (5.4) are valid. It is $\mathbf{q}_1 = \mathbf{r}^{(0)} / \|\mathbf{r}^{(0)}\|_2$. Set $\beta = \|\mathbf{r}^{(0)}\|_2$. By the orthogonality of the columns of Q_k it follows that

$$Q_k^T \mathbf{r}^{(0)} = \beta Q_k^T \mathbf{q}_1 = \beta \mathbf{e}_1. \quad (6.1)$$

Consequently, the iterate $\mathbf{x}^{(k)}$ is given by

$$\mathbf{x}^{(k)} = \mathbf{x}^{(0)} + Q_k \mathbf{y}_k \quad \text{with} \quad \mathbf{y}_k = \tilde{H}_k^{-1}(\beta \mathbf{e}_1), \quad (6.2)$$

since the orthogonal projection of $\mathbf{r}^{(0)}$ into $K_k(\mathbf{r}^{(0)}, A)^\perp$ is unique and the iterate (6.2) fulfills $\mathbf{r}^{(k)} \perp K_k(\mathbf{r}^{(0)}, A)$:

$$\begin{aligned} Q_k^T \mathbf{r}^{(k)} &= Q_k^T \left(\mathbf{b} - A\mathbf{x}^{(0)} - \beta A Q_k \tilde{H}_k^{-1} \mathbf{e}_1 \right) = Q_k^T \mathbf{r}^{(0)} - \beta \underbrace{Q_k^T A Q_k}_{=\tilde{H}_k} \tilde{H}_k^{-1} \mathbf{e}_1 \\ &= Q_k^T \mathbf{r}^{(0)} - \beta \mathbf{e}_1 = \|\mathbf{r}^{(0)}\|_2 \mathbf{e}_1 - \beta \mathbf{e}_1 = \mathbf{0}, \end{aligned}$$

where (5.4) and (6.1) have been used.

The algorithm which is based on this approach is called full orthogonalization method (FOM). Since it is of little relevance in practice, it will not be presented here in detail. Similar to GMRES, an early break down of the Arnoldi process is equivalent of already having computed the solution. FOM possesses the same great problem as GMRES since the whole basis of the $K_k(\mathbf{r}^{(0)}, A)$ has to be stored. In contrast to GMRES, the iterate of FOM is undefined if \tilde{H}_k is singular. This situation can happen, e.g., if A is a symmetric indefinite matrix. \square

6.2 Symmetric Matrices

Remark 6.3 *SYMMLQ for symmetric matrices.* If A is a symmetric matrix, there is a way to perform FOM with a short recurrence, i.e., without having to store the whole basis $\{\mathbf{q}_1, \dots, \mathbf{q}_k\}$ of $K_k(\mathbf{r}^{(0)}, A)$. The resulting method is called SYMMLQ. This method will not be presented here. Instead, the case that A is symmetric and positive definit will be studied in detail. Then, SYMMLQ can be simplified, leading to the famous conjugate gradient (CG) method. \square

Remark 6.4 *Lanczos algorithm for a s.p.d. matrix.* CG will be derived from the Lanczos algorithm 5.12. Starting point is the Cholesky¹ decomposition of \tilde{H}_k

$$\begin{aligned} \tilde{H}_k &= L_k D_k L_k^T & (6.3) \\ &= \begin{pmatrix} 1 & 0 & \cdots & 0 & 0 \\ l_1 & 1 & \cdots & 0 & 0 \\ & \ddots & \ddots & & \\ 0 & 0 & \cdots & l_{k-1} & 1 \end{pmatrix} \begin{pmatrix} d_1 & 0 & \cdots & 0 & 0 \\ 0 & d_2 & \cdots & 0 & 0 \\ & & \ddots & & \\ 0 & 0 & \cdots & & d_k \end{pmatrix} \\ &\quad \times \begin{pmatrix} 1 & l_1 & \cdots & 0 & 0 \\ & \ddots & \ddots & & \\ 0 & 0 & \cdots & 1 & l_{k-1} \\ 0 & 0 & \cdots & 0 & 1 \end{pmatrix}. \end{aligned}$$

Define $\hat{P}_k = Q_k L_k^{-T} = (\hat{\mathbf{p}}_1, \dots, \hat{\mathbf{p}}_k)$. The columns of \hat{P}_k are linear combinations of the columns of Q_k such that $\text{span}\{\hat{\mathbf{p}}_1, \dots, \hat{\mathbf{p}}_k\} \subset K_k(\mathbf{r}^{(0)}, A)$. Since L_k is a non-singular matrix and since the columns of Q_k form a basis of $K_k(\mathbf{r}^{(0)}, A)$, the columns of \hat{P}_k form a basis of $K_k(\mathbf{r}^{(0)}, A)$, too. It is for the iterate (6.2) of FOM

$$\mathbf{x}^{(k)} = \mathbf{x}^{(0)} + \beta Q_k L_k^{-T} D_k^{-1} L_k^{-1} \mathbf{e}_1 = \mathbf{x}^{(0)} + \hat{P}_k \mathbf{y}_k \quad (6.4)$$

with $\mathbf{y}_k = \beta_k D_k^{-1} L_k^{-1} \mathbf{e}_1$. \square

Lemma 6.5 *Columns of \hat{P}_k are A -conjugate.* The columns $\{\hat{\mathbf{p}}_1, \dots, \hat{\mathbf{p}}_k\}$ are mutually A -conjugate, i.e., $\hat{P}_k^T A \hat{P}_k$ is a diagonal matrix.

Proof: Using (5.4) and (6.3) gives

$$\hat{P}_k^T A \hat{P}_k = L_k^{-1} Q_k^T A Q_k L_k^{-T} = L_k^{-1} \tilde{H}_k L_k^{-T} = D_k. \quad \blacksquare$$

Remark 6.6 *First version of a method.* The last column of \hat{P}_k is given by

$$\hat{\mathbf{p}}_k = \mathbf{q}_k - l_{k-1} \hat{\mathbf{p}}_{k-1}, \quad (6.5)$$

which follows immediately from $Q_k = \hat{P}_k L_k^T$. The update \mathbf{y}_k in (6.4) has the form $\mathbf{y}_k = (\mathbf{y}_{k-1}, \eta_k)^T$ with $\eta_k \in \mathbb{R}$, since

$$\begin{aligned} \mathbf{y}_k &= \beta D_k^{-1} L_k^{-1} \mathbf{e}_1 = \beta \underbrace{\begin{pmatrix} D_{k-1} \\ d_k \end{pmatrix}^{-1}}_{\begin{pmatrix} D_{k-1}^{-1} \\ d_k^{-1} \end{pmatrix}} \underbrace{\begin{pmatrix} L_{k-1} & \mathbf{0} \\ l_{k-1} & 1 \end{pmatrix}^{-1}}_{\begin{pmatrix} L_{k-1}^{-1} & \mathbf{0} \\ * & 1 \end{pmatrix}} \mathbf{e}_{1,k} \\ &= \beta \begin{pmatrix} D_{k-1}^{-1} L_{k-1}^{-1} & \mathbf{0} \\ * & d_k^{-1} \end{pmatrix} \mathbf{e}_{1,k} = \begin{pmatrix} \beta D_{k-1}^{-1} L_{k-1}^{-1} \mathbf{e}_{1,k-1} \\ * \end{pmatrix} = \begin{pmatrix} \mathbf{y}_{k-1} \\ \eta_k \end{pmatrix}, \end{aligned}$$

¹André Louis Cholesky (1875 – 1918)

where $\mathbf{e}_{1,k}$ is the first Cartesian unit vector with k components and $\mathbf{e}_{1,k-1}$ the first Cartesian unit vector of length $(k-1)$. This means, the first $(k-1)$ components of \mathbf{y}_k are the components of \mathbf{y}_{k-1} . Now, one needs to find a formula for η_k .

From the definition of \mathbf{y}_k it follows that $L_k D_k \mathbf{y}_k = \beta \mathbf{e}_1$, i.e.,

$$L_k \begin{pmatrix} y_{k,1} d_1 \\ y_{k,2} d_2 \\ \vdots \\ \eta_k d_k \end{pmatrix} = \begin{pmatrix} \beta \\ 0 \\ \vdots \\ 0 \end{pmatrix}.$$

Hence, $l_{k-1} y_{k,k-1} d_{k-1} + \eta_k d_k = 0$ and

$$\eta_k = -\frac{l_{k-1} y_{k,k-1} d_{k-1}}{d_k}, \quad \text{if } k \geq 2. \quad (6.6)$$

The first component, η_1 is given by

$$\eta_1 = \beta D_1^{-1} L_1^{-1} \mathbf{e}_1 = \frac{\beta}{d_1}.$$

Inserting all terms into (6.4) gives

$$\begin{aligned} \mathbf{x}^{(k)} &= \mathbf{x}^{(0)} + \hat{P}_k \mathbf{y}_k = \mathbf{x}^{(0)} + \left(\hat{P}_{k-1} \hat{\mathbf{p}}_k \right) \begin{pmatrix} \mathbf{y}_{k-1} \\ \eta_k \end{pmatrix} = \mathbf{x}^{(0)} + \hat{P}_{k-1} \mathbf{y}_{k-1} + \eta_k \hat{\mathbf{p}}_k \\ &= \mathbf{x}^{(k-1)} + \eta_k \hat{\mathbf{p}}_k. \end{aligned} \quad (6.7)$$

Thus, the new iterate can be computed with (6.5) and (6.6). This approach shows that a short recurrence is possible. However, it is not yet optimal and it can be simplified. \square

Remark 6.7 *Optimal version of the method.* It holds for the residual, using (6.4), that

$$\mathbf{r}^{(k)} = \mathbf{b} - A\mathbf{x}^{(k)} = \mathbf{b} - A\mathbf{x}^{(0)} - A\hat{P}_k \mathbf{y}_k = \mathbf{r}^{(0)} - A\mathbf{z}$$

with some vector $\mathbf{z} \in K_k(\mathbf{r}^{(0)}, A) = \text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_k\}$ since the columns of \hat{P}_k form a basis of $K_k(\mathbf{r}^{(0)}, A)$, see Remark 6.4. This representation shows first that

$$\mathbf{r}^{(k)} \in K_{k+1}(\mathbf{r}^{(0)}, A) = \text{span}\{\mathbf{q}_1, \dots, \mathbf{q}_k, \mathbf{q}_{k+1}\}.$$

By construction, see Remark 6.2, it is also $\mathbf{r}^{(k)} \perp K_k(\mathbf{r}^{(0)}, A)$. These two properties imply that $\mathbf{r}^{(k)} \in \text{span}\{\mathbf{q}_{k+1}\}$ such that

$$\mathbf{r}^{(k)} = \pm \left\| \mathbf{r}^{(k)} \right\|_2 \mathbf{q}_{k+1}.$$

Using (6.7), the residual vector $\mathbf{r}^{(k)}$ can be computed recursively by

$$\mathbf{r}^{(k)} = \mathbf{b} - A\mathbf{x}^{(k)} = \mathbf{b} - A \left(\mathbf{x}^{(k-1)} + \eta_k \hat{\mathbf{p}}_k \right) = \mathbf{r}^{(k-1)} - \eta_k A \hat{\mathbf{p}}_k. \quad (6.8)$$

Setting

$$\mathbf{q}_{k+1} = \frac{\mathbf{r}^{(k)}}{\left\| \mathbf{r}^{(k)} \right\|_2} \quad (6.9)$$

and denoting $\mathbf{p}_k = \left\| \mathbf{r}^{(k-1)} \right\|_2 \hat{\mathbf{p}}_k$, one obtains with (6.7) and (6.8)

$$\mathbf{r}^{(k)} = \mathbf{r}^{(k-1)} - \frac{\eta_k}{\left\| \mathbf{r}^{(k-1)} \right\|_2} A \mathbf{p}_k = \mathbf{r}^{(k-1)} - \nu_k A \mathbf{p}_k, \quad (6.10)$$

$$\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} + \frac{\eta_k}{\left\| \mathbf{r}^{(k-1)} \right\|_2} \mathbf{p}_k = \mathbf{x}^{(k-1)} + \nu_k \mathbf{p}_k. \quad (6.11)$$

With (6.5) and (6.9), one gets

$$\begin{aligned}\mathbf{p}_k &= \left\| \mathbf{r}^{(k-1)} \right\|_2 \hat{\mathbf{p}}_k = \left\| \mathbf{r}^{(k-1)} \right\|_2 \left(\mathbf{q}_k - l_{k-1} \frac{\mathbf{p}_{k-1}}{\left\| \mathbf{r}^{(k-2)} \right\|_2} \right) \\ &= \mathbf{r}^{(k-1)} - \frac{\left\| \mathbf{r}^{(k-1)} \right\|_2 l_{k-1}}{\left\| \mathbf{r}^{(k-2)} \right\|_2} \mathbf{p}_{k-1} = \mathbf{r}^{(k-1)} + \mu_k \mathbf{p}_{k-1}.\end{aligned}\quad (6.12)$$

Now, formulas for ν_k and μ_k are needed. Multiplying (6.12) from the left hand side with $\mathbf{p}_k^T A$ gives

$$\mathbf{p}_k^T A \mathbf{p}_k = \mathbf{p}_k^T A \mathbf{r}^{(k-1)} + \mu_k \left\| \mathbf{r}^{(k-1)} \right\|_2^2 \underbrace{\hat{\mathbf{p}}_k^T A \hat{\mathbf{p}}_{k-1}}_{=0, \text{Lemma 6.5}} = \mathbf{p}_k^T A \mathbf{r}^{(k-1)}.\quad (6.13)$$

Multiplying (6.10) from left with $(\mathbf{r}^{(k-1)})^T$ and using $\mathbf{r}^{(j)} = c\mathbf{q}_j, j = 1, \dots, k$, and the orthonormality of the vectors \mathbf{q}_j leads to

$$\underbrace{(\mathbf{r}^{(k-1)})^T \mathbf{r}^{(k)}}_{=0} = (\mathbf{r}^{(k-1)})^T \mathbf{r}^{(k-1)} - \nu_k (\mathbf{r}^{(k-1)})^T A \mathbf{p}_k,$$

which gives together with (6.13)

$$\nu_k = \frac{(\mathbf{r}^{(k-1)})^T \mathbf{r}^{(k-1)}}{(\mathbf{r}^{(k-1)})^T A \mathbf{p}_k} = \frac{(\mathbf{r}^{(k-1)})^T \mathbf{r}^{(k-1)}}{\mathbf{p}_k^T A \mathbf{p}_k}.\quad (6.14)$$

Now, multiplying (6.12) from left with $\mathbf{p}_{k-1}^T A$ leads to, using Lemma 6.5,

$$\underbrace{\mathbf{p}_{k-1}^T A \mathbf{p}_k}_{=0} = \mathbf{p}_{k-1}^T A \mathbf{r}^{(k-1)} + \mu_k \mathbf{p}_{k-1}^T A \mathbf{p}_{k-1},$$

which gives

$$\mu_k = -\frac{\mathbf{p}_{k-1}^T A \mathbf{r}^{(k-1)}}{\mathbf{p}_{k-1}^T A \mathbf{p}_{k-1}}.$$

To simplify this expression, multiply (6.10) from left with $(\mathbf{r}^{(k)})^T$ such that one obtains, using $\mathbf{r}^{(k)} \perp \mathbf{r}^{(k-1)}$,

$$(\mathbf{r}^{(k)})^T \mathbf{r}^{(k)} = 0 - \nu_k (\mathbf{r}^{(k)})^T A \mathbf{p}_k \implies \nu_k = -\frac{(\mathbf{r}^{(k)})^T \mathbf{r}^{(k)}}{(\mathbf{r}^{(k)})^T A \mathbf{p}_k}.$$

This expression gives together with (6.14)

$$-\frac{(\mathbf{r}^{(k)})^T A \mathbf{p}_k}{\mathbf{p}_k^T A \mathbf{p}_k} = \frac{(\mathbf{r}^{(k)})^T \mathbf{r}^{(k)}}{(\mathbf{r}^{(k-1)})^T \mathbf{r}^{(k-1)}}$$

such that

$$\mu_k = \frac{(\mathbf{r}^{(k-1)})^T \mathbf{r}^{(k-1)}}{(\mathbf{r}^{(k-2)})^T \mathbf{r}^{(k-2)}}.\quad (6.15)$$

The evaluation of this expression requires only two inner products but not any matrix-vector product. These considerations lead to Algorithm 6.8. \square

Algorithm 6.8 Conjugate Gradient (CG). Given a symmetric positive definite matrix $A \in \mathbb{R}^{n \times n}$, a right hand side $\mathbf{b} \in \mathbb{R}^n$, an initial iterate $\mathbf{x}^{(0)} \in \mathbb{R}^n$ and a tolerance $\varepsilon > 0$.

1. $\mathbf{r}^{(0)} = \mathbf{b} - A\mathbf{x}^{(0)}$
2. $\mathbf{p}_1 = \mathbf{r}^{(0)}$
3. $k = 0$
4. **while** $\|\mathbf{r}^{(k)}\|_2 > \varepsilon$
5. $k = k + 1$
6. $\mathbf{s} = A\mathbf{p}_k$
7. $\nu_k = \frac{(\mathbf{r}^{(k-1)})^T \mathbf{r}^{(k-1)}}{\mathbf{p}_k^T \mathbf{s}}$ % (6.14)
8. $\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} + \nu_k \mathbf{p}_k$ % (6.11)
9. $\mathbf{r}^{(k)} = \mathbf{r}^{(k-1)} - \nu_k \mathbf{s}$ % (6.10)
10. $\mu_{k+1} = \frac{(\mathbf{r}^{(k)})^T \mathbf{r}^{(k)}}{(\mathbf{r}^{(k-1)})^T \mathbf{r}^{(k-1)}}$ % (6.15)
11. $\mathbf{p}_{k+1} = \mathbf{r}^{(k)} + \mu_{k+1} \mathbf{p}_k$ % (6.12)
12. **endwhile**

□

Remark 6.9 *First publication of CG.* The CG method has been published the first time by Hestenes² and Stiefel³ in Hestenes and Stiefel (1952). □

Remark 6.10 *Costs of CG.* The costs of one CG iteration are:

- one matrix-vector multiplication, line 6,
- three additions of vectors in \mathbb{R}^n , lines 8, 9, 11,
- three multiplications of vectors with a scalar, lines 8, 9, 11,
- two inner product of vectors, lines 7, 10. The inner product $(\mathbf{r}^{(k-1)})^T \mathbf{r}^{(k-1)}$ is already known from the previous iteration.

One has to store four vectors: $\mathbf{x}^{(k)}, \mathbf{r}^{(k)}, \mathbf{p}_k, \mathbf{s}$. In comparison with the conjugate residual method, CG needs one vector update ($2n$ flops) less and one has to store one vector less. Since both schemes exhibit in general a similar convergence, CG is generally preferred.

Altogether, CG is in general the best performing iterative scheme without a multigrid component for solving linear systems of equations with a symmetric positive definite matrix. □

Definition 6.11 Energy norm. Let $A \in \mathbb{R}^{n \times n}$ be s.p.d., then A induces a vector norm by

$$\|\mathbf{x}\|_A = (\mathbf{x}, A\mathbf{x})^{1/2} \quad \forall \mathbf{x} \in \mathbb{R}^n,$$

the so-called energy norm. □

Theorem 6.12 Minimization of the error in the energy norm. Let $A \in \mathbb{R}^{n \times n}$ be symmetric and positive definite. The iterate

$$\mathbf{x}^{(k)} = \mathbf{x}^{(0)} + \left\| \mathbf{r}^{(0)} \right\|_2 Q_k \tilde{H}_k^{-1} \mathbf{e}_1$$

is well defined and it is the solution of

$$\min_{\mathbf{y} \in \{\mathbf{x}^{(0)} + K_k(\mathbf{r}^{(0)}, A)\}} \|\mathbf{x} - \mathbf{y}\|_A,$$

²Magnus Rudolph Hestenes (1906 – 1991)

³Eduard L. Stiefel (1909 – 1978)

where \mathbf{x} is the solution of (1.1). The corresponding residual $\mathbf{r}^{(k)}$ is orthogonal to $K_k(\mathbf{r}^{(0)}, A)$, i.e., $Q_k^T \mathbf{r}^{(k)} = \mathbf{0}$.

Proof: The non-singularity of the matrix \tilde{H}_k , Lemma 5.14, induces that the iterate $\mathbf{x}^{(k)}$ is well defined. The orthogonality of $\mathbf{r}^{(k)}$ and $K_k(\mathbf{r}^{(0)}, A)$ follows by the construction of the method, see Remark 6.2.

Let $\mathbf{y} \in \{\mathbf{x}^{(0)} + K_k(\mathbf{r}^{(0)}, A)\}$, $\mathbf{y} \neq \mathbf{x}^{(k)}$, and denote $\mathbf{z} = \mathbf{y} - \mathbf{x}^{(k)} \in K_k(\mathbf{r}^{(0)}, A)$. Using the symmetry of A , the orthogonality of $\mathbf{r}^{(k)}$ to $\mathbf{z} \in K_k(\mathbf{r}^{(0)}, A)$ and the positive definiteness of A gives

$$\begin{aligned}
& \|\mathbf{x} - \mathbf{y}\|_A^2 \\
&= (\mathbf{x} - \mathbf{y})^T A (\mathbf{x} - \mathbf{y}) = (\mathbf{x} - \mathbf{x}^{(k)} - \mathbf{z})^T A (\mathbf{x} - \mathbf{x}^{(k)} - \mathbf{z}) \\
&= (\mathbf{x} - \mathbf{x}^{(k)})^T A (\mathbf{x} - \mathbf{x}^{(k)}) - \mathbf{z}^T A (\mathbf{x} - \mathbf{x}^{(k)}) - \underbrace{(\mathbf{x} - \mathbf{x}^{(k)})^T A \mathbf{z}}_{= ((\mathbf{x} - \mathbf{x}^{(k)})^T A \mathbf{z})^T \in \mathbb{R}} + \mathbf{z}^T A \mathbf{z} \\
&= (\mathbf{x} - \mathbf{x}^{(k)})^T A (\mathbf{x} - \mathbf{x}^{(k)}) - 2\mathbf{z}^T A (\mathbf{x} - \mathbf{x}^{(k)}) + \mathbf{z}^T A \mathbf{z} \\
&= (\mathbf{x} - \mathbf{x}^{(k)})^T A (\mathbf{x} - \mathbf{x}^{(k)}) - \underbrace{2\mathbf{z}^T \mathbf{r}^{(k)}}_{=0} + \underbrace{\mathbf{z}^T A \mathbf{z}}_{>0} \\
&> (\mathbf{x} - \mathbf{x}^{(k)})^T A (\mathbf{x} - \mathbf{x}^{(k)}) = \|\mathbf{x} - \mathbf{x}^{(k)}\|_A^2.
\end{aligned}$$

■

Remark 6.13 *On the energy norm.* To minimize the error in the energy norm is more natural than to minimize the Euclidean norm of the residual since

$$\|\mathbf{r}^{(k)}\|_2 = \|A(\mathbf{x} - \mathbf{x}^{(k)})\|_2 = \|\mathbf{x} - \mathbf{x}^{(k)}\|_{A^2}.$$

The energy norm is the natural measure for the error.

In the literature, one can find the derivation of the CG method also with the starting point of trying to minimize the error in the energy norm. One finds that the most simple approach, the steepest descent method, converges very slowly. Considerations on improving the iterative scheme lead finally to the CG method. □