

## 15 Conjugate Gradients

This method for symmetric positive definite matrices is considered to be the “original” Krylov subspace method. It was proposed by Hestenes and Stiefel in 1952, and is motivated by the following theorem.

**Theorem 15.1** *If  $A$  is symmetric positive definite, then solving  $A\mathbf{x} = \mathbf{b}$  is equivalent to minimizing the quadratic form*

$$\varphi(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T A\mathbf{x} - \mathbf{x}^T \mathbf{b}.$$

**Proof** We will consider changes of  $\varphi$  along a certain ray  $\mathbf{x} + \alpha\mathbf{p}$  with  $\alpha \in \mathbb{R}$ , and fixed direction vector  $\mathbf{p} \neq \mathbf{0}$ .

First we show that  $\varphi(\mathbf{x} + \alpha\mathbf{p}) > \varphi(\mathbf{x})$  if  $A$  symmetric positive definite.

$$\begin{aligned} \varphi(\mathbf{x} + \alpha\mathbf{p}) &= \frac{1}{2}(\mathbf{x} + \alpha\mathbf{p})^T A(\mathbf{x} + \alpha\mathbf{p}) - (\mathbf{x} + \alpha\mathbf{p})^T \mathbf{b} \\ &= \frac{1}{2}\mathbf{x}^T A\mathbf{x} + \frac{1}{2}\mathbf{x}^T A(\alpha\mathbf{p}) + \frac{1}{2}(\alpha\mathbf{p})^T A\mathbf{x} + \frac{1}{2}(\alpha\mathbf{p})^T A(\alpha\mathbf{p}) - \mathbf{x}^T \mathbf{b} - \alpha\mathbf{p}^T \mathbf{b} \\ &\stackrel{A^T=A}{=} \underbrace{\frac{1}{2}\mathbf{x}^T A\mathbf{x} - \mathbf{x}^T \mathbf{b}}_{=\varphi(\mathbf{x})} + \alpha\mathbf{p}^T A\mathbf{x} - \alpha\mathbf{p}^T \mathbf{b} + \frac{1}{2}\alpha^2 \mathbf{p}^T A\mathbf{p} \\ &= \varphi(\mathbf{x}) + \alpha\mathbf{p}^T (A\mathbf{x} - \mathbf{b}) + \frac{1}{2}\alpha^2 \underbrace{\mathbf{p}^T A\mathbf{p}}_{>0}. \end{aligned}$$

Thus, we see that  $\varphi$  (as a quadratic function in  $\alpha$  with positive leading coefficient) will have to have a minimum along the ray  $\mathbf{x} + \alpha\mathbf{p}$ .

We now decide what the value of  $\alpha$  at this minimum is. A necessary condition (and also sufficient since the coefficient of  $\alpha^2$  is positive) is

$$\frac{d}{d\alpha}\varphi(\mathbf{x} + \alpha\mathbf{p}) = 0.$$

To this end we compute

$$\frac{d}{d\alpha}\varphi(\mathbf{x} + \alpha\mathbf{p}) = \mathbf{p}^T (A\mathbf{x} - \mathbf{b}) + \alpha\mathbf{p}^T A\mathbf{p},$$

which has its root at

$$\hat{\alpha} = \frac{\mathbf{p}^T (\mathbf{b} - A\mathbf{x})}{\mathbf{p}^T A\mathbf{p}}.$$

The corresponding minimum value is

$$\varphi(\mathbf{x} + \hat{\alpha}\mathbf{p}) = \varphi(\mathbf{x}) - \underbrace{\frac{[\mathbf{p}^T (\mathbf{b} - A\mathbf{x})]^2}{2\mathbf{p}^T A\mathbf{p}}}_{\geq 0}.$$

The last equation shows that

$$\varphi(\mathbf{x} + \hat{\alpha}\mathbf{p}) < \varphi(\mathbf{x}) \quad \text{if and only if} \quad \mathbf{p}^T (\mathbf{b} - A\mathbf{x}) \neq 0,$$

i.e.,  $\mathbf{p}$  is not orthogonal to the residual  $\mathbf{r} = \mathbf{b} - A\mathbf{x}$ .

To see the equivalence with the solution of the linear system  $A\mathbf{x} = \mathbf{b}$  we need to consider two possibilities:

1.  $\mathbf{x}$  is such that  $A\mathbf{x} = \mathbf{b}$ . Then  $\varphi(\mathbf{x} + \hat{\alpha}\mathbf{p}) = \varphi(\mathbf{x})$  and  $\varphi(\mathbf{x})$  is the minimum value.
2.  $\mathbf{x}$  is such that  $A\mathbf{x} \neq \mathbf{b}$ . Then  $\varphi(\mathbf{x} + \hat{\alpha}\mathbf{p}) < \varphi(\mathbf{x})$ , i.e., there exists a direction  $\mathbf{p}$  such that  $\mathbf{p}^T(\mathbf{b} - A\mathbf{v}) \neq 0$  and  $\varphi(\mathbf{x})$  is not the minimum.

■

The preceding proof actually suggests a rough iterative algorithm:

Take  $\mathbf{x}_0 = \mathbf{0}$ ,  $\mathbf{r}_0 = \mathbf{b}$ ,  $\mathbf{p}_0 = \mathbf{r}_0$

for  $n = 1, 2, 3, \dots$

    Compute a step length

$$\alpha_n = (\mathbf{p}_{n-1}^T \mathbf{r}_{n-1}) / (\mathbf{p}_{n-1}^T A \mathbf{p}_{n-1})$$

    Update the approximate solution

$$\mathbf{x}_n = \mathbf{x}_{n-1} + \alpha_n \mathbf{p}_{n-1}$$

    Update the residual

$$\mathbf{r}_n = \mathbf{r}_{n-1} - \alpha_n A \mathbf{p}_{n-1}$$

    Find a new search direction  $\mathbf{p}_n$

end

Note that at this point we have not specified how to pick the search directions  $\mathbf{p}_n$ . This will be the crucial ingredient in the algorithm.

The formula above for the residual update follows from

$$\begin{aligned} \mathbf{r}_n &= \mathbf{b} - A\mathbf{x}_n = \mathbf{b} - A(\mathbf{x}_{n-1} + \alpha_n \mathbf{p}_{n-1}) \\ &= \mathbf{b} - A\mathbf{x}_{n-1} - \alpha_n A \mathbf{p}_{n-1} \\ &= \mathbf{r}_{n-1} - \alpha_n A \mathbf{p}_{n-1}. \end{aligned}$$

## 15.1 The Steepest Descent Algorithm

An obvious choice for the selection of the search direction is

$$\mathbf{p}_n = -\nabla\varphi(\mathbf{x}_n)$$

since we know from calculus that the direction of largest decrease of  $\varphi$  is in the direction opposite its gradient. Moreover, since  $\varphi(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T A \mathbf{x} - \mathbf{x}^T \mathbf{b}$  we have

$$\nabla\varphi(\mathbf{x}) = A\mathbf{x} - \mathbf{b}.$$

This leads to

**Algorithm** (Steepest Descent)

Take  $\mathbf{x}_0 = \mathbf{0}$ ,  $\mathbf{r}_0 = \mathbf{b}$ ,  $\mathbf{p}_0 = \mathbf{r}_0$

for  $n = 1, 2, 3, \dots$

    Compute a step length

$$\alpha_n = (\mathbf{p}_{n-1}^T \mathbf{r}_{n-1}) / (\mathbf{p}_{n-1}^T A \mathbf{p}_{n-1})$$

    Update the approximate solution

$$\mathbf{x}_n = \mathbf{x}_{n-1} + \alpha_n \mathbf{p}_{n-1}$$

    Update the residual

$$\mathbf{r}_n = \mathbf{r}_{n-1} - \alpha_n A \mathbf{p}_{n-1}$$

    Set the new search direction

$$\mathbf{p}_n = \mathbf{r}_n$$

end

Note that for this choice of search direction the step length  $\alpha$  can also be written as

$$\alpha_n = (\mathbf{r}_{n-1}^T \mathbf{r}_{n-1}) / (\mathbf{p}_{n-1}^T A \mathbf{p}_{n-1}).$$

## 15.2 The Conjugate Gradient Algorithm

It turns out that the “obvious” search directions are not ideal (since they are applied in an iterative fashion). Convergence of the steepest descent algorithm is usually rather slow. It is better to employ so-called *conjugate search directions*. The idea is to somehow remove from the gradient at each step those components parallel to previously used search directions. The resulting algorithm is

**Algorithm** (Conjugate Gradient)

Take  $\mathbf{x}_0 = \mathbf{0}$ ,  $\mathbf{r}_0 = \mathbf{b}$ ,  $\mathbf{p}_0 = \mathbf{r}_0$

for  $n = 1, 2, 3, \dots$

    Compute a step length

$$\alpha_n = (\mathbf{r}_{n-1}^T \mathbf{r}_{n-1}) / (\mathbf{p}_{n-1}^T A \mathbf{p}_{n-1})$$

    Update the approximate solution

$$\mathbf{x}_n = \mathbf{x}_{n-1} + \alpha_n \mathbf{p}_{n-1}$$

    Update the residual

$$\mathbf{r}_n = \mathbf{r}_{n-1} - \alpha_n A \mathbf{p}_{n-1}$$

Compute a gradient correction factor

$$\beta_n = (\mathbf{r}_n^T \mathbf{r}_n) / (\mathbf{r}_{n-1}^T \mathbf{r}_{n-1})$$

Set the new search direction

$$\mathbf{p}_n = \mathbf{r}_n + \beta_n \mathbf{p}_{n-1}$$

end

For both the steepest descent and the conjugate gradient algorithm the main computational cost is hidden in the one matrix-vector multiplication that is required per iteration. As with the Arnoldi and GMRES methods, this operation is treated as a “black box” and can be accomplished in  $\mathcal{O}(m)$  to  $\mathcal{O}(m^2)$  operations depending on the structure of  $A$ . In many practical cases the entire (preconditioned) CG algorithm will require only  $\mathcal{O}(m)$  operations. This is very fast.

As mentioned at the beginning of this section, one can also establish a connection to Krylov subspace methods.

**Theorem 15.2** *Let  $A$  be symmetric positive definite. As long as the conjugate gradient method has not yet converged (i.e., as long as  $\mathbf{r}_{n-1} \neq \mathbf{0}$ ) we have*

$$\begin{aligned} \text{span}\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} &= \text{span}\{\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{n-1}\} \\ &= \text{span}\{\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{n-1}\} \\ &= \text{span}\{\mathbf{b}, A\mathbf{b}, \dots, A^{n-1}\mathbf{b}\} = \mathcal{K}_n. \end{aligned}$$

Moreover, the residuals are orthogonal in the usual sense, i.e.,

$$\mathbf{r}_n^T \mathbf{r}_j = 0, \quad j < n,$$

and the search directions are  $A$ -orthogonal (or  $A$ -conjugate), i.e.,

$$\mathbf{p}_n^T A \mathbf{p}_j = 0, \quad j < n.$$

**Proof** An inductive proof of this theorem can be found in the book [Trefethen/Bau].

■

### 15.3 Convergence of the CG Algorithm

Recall that the GMRES algorithm minimizes the 2-norm of the residual,  $\|\mathbf{r}_n\|_2 \rightarrow \min$ . We will now show that the CG algorithm satisfies a different optimality criterion. It minimizes the  $A$ -norm of the error, i.e., if  $\mathbf{e}_n = \mathbf{x}^* - \mathbf{x}_n$ , is the error between the exact solution  $\mathbf{x}^* = A^{-1}\mathbf{b}$  and the  $n$ -th approximation  $\mathbf{x}_n$ , then CG minimizes

$$\|\mathbf{e}_n\|_A = \sqrt{\mathbf{e}_n^T A \mathbf{e}_n}.$$

**Theorem 15.3** *Let  $A$  be symmetric positive definite. If the conjugate gradient algorithm has not yet converged (i.e.,  $\mathbf{r}_{n-1} \neq \mathbf{0}$ ) then  $\mathbf{x}_n$  is the unique vector in  $\mathcal{K}_n$  such that  $\|\mathbf{e}_n\|_A$  is minimized.*

Moreover,  $\|\mathbf{e}_n\|_A \leq \|\mathbf{e}_{n-1}\|_A$  and (if we are using exact arithmetic)  $\mathbf{e}_n = \mathbf{0}$  for some  $n \leq m$ .

**Proof** We will prove the first part only. From the previous theorem we know that the approximate solution  $\mathbf{x}_n$  lies in the Krylov subspace  $\mathcal{K}_n$ . In order to show that  $\mathbf{x}_n$  is the unique minimizer of  $\|\mathbf{e}\|_A$  we consider an arbitrary vector

$$\mathbf{x} = \mathbf{x}_n - \Delta\mathbf{x} \in \mathcal{K}_n$$

and show that in order to minimize  $\|\mathbf{e}\|_A$  we necessarily have  $\Delta\mathbf{x} = \mathbf{0}$ .

If  $\mathbf{x}^*$  is the exact solution of  $A\mathbf{x} = \mathbf{b}$ , then

$$\mathbf{e} = \mathbf{x}^* - \mathbf{x} = \mathbf{x}^* - \mathbf{x}_n + \Delta\mathbf{x} = \mathbf{e}_n + \Delta\mathbf{x}.$$

Therefore,

$$\begin{aligned} \|\mathbf{e}\|_A^2 &= \|\mathbf{e}_n + \Delta\mathbf{x}\|_A^2 \\ &= (\mathbf{e}_n + \Delta\mathbf{x})^T A (\mathbf{e}_n + \Delta\mathbf{x}) \\ &= \mathbf{e}_n^T A \mathbf{e}_n + 2\mathbf{e}_n^T A (\Delta\mathbf{x}) + (\Delta\mathbf{x})^T A (\Delta\mathbf{x}) \end{aligned}$$

since  $A$  is symmetric.

Next we realize that

$$\mathbf{e}_n^T A = (\mathbf{x}^* - \mathbf{x}_n)^T A = (A^{-1}\mathbf{b} - \mathbf{x}_n)^T A = \mathbf{b}^T - \mathbf{x}_n^T A = \mathbf{r}_n^T,$$

and observe that

$$\mathbf{r}_n^T \Delta\mathbf{x} = 0$$

since  $\Delta\mathbf{x} \in \mathcal{K}_n = \text{span}\{\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_{n-1}\}$  and  $\mathbf{r}_n^T \mathbf{r}_j = 0$  for  $j < n$  by the previous theorem.

This leaves us with

$$\begin{aligned} \|\mathbf{e}\|_A^2 &= \mathbf{e}_n^T A \mathbf{e}_n + 2 \underbrace{\mathbf{e}_n^T A (\Delta\mathbf{x})}_{=\mathbf{r}_n^T \Delta\mathbf{x}} + (\Delta\mathbf{x})^T A (\Delta\mathbf{x}) \\ &= \|\mathbf{e}_n\|_A^2 + (\Delta\mathbf{x})^T A (\Delta\mathbf{x}). \end{aligned}$$

Note that the quadratic form  $(\Delta\mathbf{x})^T A (\Delta\mathbf{x})$  is certainly non-negative since  $A$  is positive definite. Moreover, it is zero only if  $\Delta\mathbf{x} = \mathbf{0}$ .

Thus, the  $A$ -norm of the error is minimized if  $\Delta\mathbf{x} = \mathbf{0}$ , i.e., for the CG approximate solution  $\mathbf{x}_n$ . ■

We can come to the same conclusion with the following argument:

$$\begin{aligned} \|\mathbf{e}_n\|_A^2 &= \mathbf{e}_n^T A \mathbf{e}_n = (\mathbf{x}^* - \mathbf{x}_n)^T A (\mathbf{x}^* - \mathbf{x}_n) \\ &= (\mathbf{x}^*)^T \underbrace{A \mathbf{x}^*}_{=\mathbf{b}} - 2\mathbf{x}_n^T \underbrace{A \mathbf{x}^*}_{=\mathbf{b}} + \mathbf{x}_n^T A \mathbf{x}_n \\ &= (\mathbf{x}^*)^T \mathbf{b} + \mathbf{x}_n^T A \mathbf{x}_n - 2\mathbf{x}_n^T \mathbf{b} \\ &= (\mathbf{x}^*)^T \mathbf{b} + 2\varphi(\mathbf{x}_n). \end{aligned}$$

Here  $\varphi(\mathbf{x}_n)$  is the same quadratic form used earlier. Since  $(\mathbf{x}^*)^T \mathbf{b}$  is a constant we see that minimizing the  $A$ -norm of the error is equivalent to minimizing the quadratic form  $\varphi(\mathbf{x}_n)$ .

For the *convergence rate* of the CG algorithm one can show that

$$\|\mathbf{e}_n\|_A \leq \|\mathbf{e}_0\|_A \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^n,$$

where  $\kappa = \kappa_2(A)$  the 2-norm condition number of  $A$ . Since

$$\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} = 1 - \frac{2}{\sqrt{\kappa} + 1}$$

we see that convergence will be very slow if  $\kappa$  is large. This shows that preconditioning efforts for the CG algorithm are aimed at reducing the condition number of  $A$ .

For a moderate size  $\kappa$  it turns out that one can expect convergence of the CG algorithm in  $\mathcal{O}(\sqrt{\kappa})$  iterations. In fact, in practice the CG algorithm often converges faster than predicted by this upper bound.

**Remark** It is possible to interpret the conjugate gradient method as an analogue of Lanczos iteration for linear systems. Since we claimed earlier that Lanczos iteration is a special case of Arnoldi iteration for symmetric matrices, it turns out that the  $(n+1)$ -term recursion we derived earlier for Arnoldi iteration turn into a 3-term recursion. One can indeed show that this 3-term recursion is hidden inside the CG algorithm.

Convergence of the CG algorithm is illustrated in the MATLAB code `CGDemo.m`. The symmetric test matrix is constructed as follows. Initially it contains ones on the main diagonal and random numbers uniformly distributed in  $[-1, 1]$  in the off-diagonal positions. Then any off-diagonal entry with  $|a_{ij}| > \tau$  is set to zero, where  $\tau$  is a parameter. For small values of  $\tau$  the matrix is positive definite and very sparse, and the CG algorithm converges rapidly. For larger values, such as  $\tau = 0.2$  the matrix is no longer positive definite, and the CG algorithm does not converge. We also note that for these test matrices, preconditioning does not improve convergence of the CG algorithm.