

## Lecture Note 4: The Conjugate Gradient Method

Xianyi Zeng

Department of Mathematical Sciences, UTEP

### 1 A Starting Point: The Steepest Descent Method

If  $A \in \mathbb{R}^{n \times n}$  is symmetric positive-definite, solving:

$$A\mathbf{x} = \mathbf{b}, \quad (1.1)$$

is equivalent to minimizing the function:

$$\phi(\mathbf{x}) = \frac{1}{2} \mathbf{x}^t A \mathbf{x} - \mathbf{x}^t \mathbf{b}. \quad (1.2)$$

Thus we can apply any optimization algorithm to solve this minimization problem and obtain a method for solving (1.1).

At the point, let us consider the *steepest descent method* and select any initial guess  $\mathbf{x}_0$ . With  $\mathbf{x}_k$  available we try to find the direction along which  $\phi(\mathbf{x})$  decreases most rapidly starting from  $\mathbf{x}_k$  and compute the next point  $\mathbf{x}_{k+1}$  by minimizing  $\phi(\mathbf{x})$  in this direction. By Taylor series expansion

$$\phi(\mathbf{x}_k + \alpha \mathbf{d}) = \phi(\mathbf{x}_k) + \alpha \nabla \phi(\mathbf{x}_k)^t \mathbf{d} + O(\|\mathbf{d}\|^2);$$

thus the direction we're looking for is given by  $-\nabla \phi(\mathbf{x}_k) = \mathbf{r}_k$ , the *residual* at the  $k$ -th iteration:

$$\mathbf{r}_k = \mathbf{b} - A\mathbf{x}_k. \quad (1.3)$$

If  $\mathbf{r}_k \neq 0$ , we try to find the next solution point  $\mathbf{x}_{k+1}$  by minimizing  $\phi(\mathbf{x}_k + \alpha \mathbf{r}_k)$  for all  $\alpha \in \mathbb{R}$  (called the *exact line search*). Note that  $\phi(\mathbf{x}_k + \alpha \mathbf{r}_k)$  is a second-degree polynomial in  $\alpha$ :

$$\phi(\mathbf{x}_k + \alpha \mathbf{r}_k) = \frac{1}{2} (\mathbf{r}_k^t A \mathbf{r}_k) \alpha^2 + (\mathbf{r}_k^t A \mathbf{x}_k - \mathbf{r}_k^t \mathbf{b}) \alpha + \phi(\mathbf{x}_k),$$

the solution to the exact line search is:

$$\alpha_k = \frac{\mathbf{r}_k^t \mathbf{r}_k}{\mathbf{r}_k^t A \mathbf{r}_k};$$

the denominator is never zero due to the positive-definiteness of  $A$  and  $\mathbf{r}_k \neq 0$ , thus we can compute  $\mathbf{x}_{k+1}$  as well as  $\mathbf{r}_{k+1}$ . This leads to the following (idealized) algorithm:

$$\mathbf{r}_0 = \mathbf{b} - A\mathbf{x}_0, \quad k = 0; \quad (1.4)$$

while  $\mathbf{r}_k \neq 0$  :

$$\alpha_k = (\mathbf{r}_k^t \mathbf{r}_k) / (\mathbf{r}_k^t A \mathbf{r}_k);$$

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

$$\mathbf{r}_{k+1} = \mathbf{b} - A\mathbf{x}_{k+1};$$

$$k = k + 1;$$

end

$$\mathbf{x} = \mathbf{x}_k.$$

A theoretical bound on the convergence of the steepest descent method is given by:

$$\left(\phi(\mathbf{x}_{k+1}) + \frac{1}{2}\mathbf{b}^t A^{-1}\mathbf{b}\right) \leq \left(1 - \frac{\lambda_{\min}(A)}{\lambda_{\max}(A)}\right) \left(\phi(\mathbf{x}_k) + \frac{1}{2}\mathbf{b}^t A^{-1}\mathbf{b}\right), \quad (1.5)$$

where  $-(\mathbf{b}^t A^{-1}\mathbf{b})/2$  is the theoretical minimum of  $\phi(\mathbf{x})$ , and  $\lambda_{\min}(A)$  and  $\lambda_{\max}(A)$  are the smallest and largest eigenvalues of  $A$ . Thus the global convergence is guaranteed due to the fact that  $0 < \lambda_{\min}(A) < \lambda_{\max}(A)$ .

## 2 The Projection Method

An issue with the steepest descent method is that the decay rate of (1.5) can be extremely close to 1 in many practical applications, which means a large number of iterations is needed to achieve a certain accuracy.

The *projection method* can be viewed as an extension of the algorithm (1.4) in the following sense. In every iteration of the steepest descent algorithm, we compute a solution  $\mathbf{x}_{k+1}$  in the one-dimensional affine space  $\mathbf{x}_k + \text{span}(\mathbf{r}_k)$ ; thus it is also contained in the *affine space*:

$$\mathbf{x}_0 + \text{span}(\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_k),$$

which has the dimension no larger than  $\min(n, k+1)$ .

**Remark 1.** An affine space  $\mathcal{A} \subseteq \mathcal{V}$ , where  $\mathcal{V}$  is a vector space, is not necessarily a vector space itself, at least in the sense that it does not necessarily contain a zero vector. A formal definition for the affine space requires that if  $\mathbf{v}_1, \mathbf{v}_2 \in \mathcal{A}$  and  $\alpha \in \mathbb{R}$ , then  $(1-\alpha)\mathbf{v}_1 + \alpha\mathbf{v}_2 \in \mathcal{A}$ . It is not difficult to show that let  $\mathbf{v}_0 \in \mathcal{A}$  be arbitrary, then the set:

$$\{\mathbf{v} : \mathbf{v} + \mathbf{v}_0 \in \mathcal{A}\}$$

is a linear subspace of  $\mathcal{V}$ . Denote this linear subspace by  $\mathcal{S}$ , then we write  $\mathcal{A} = \mathbf{v}_0 + \mathcal{S}$ . One can check that  $\mathcal{S}$  is independent of the particular choice of  $\mathbf{v}_0$ .

**Remark 2.** Let  $\mathbf{v}_1, \dots, \mathbf{v}_k \in \mathcal{V}$ , then  $\text{span}(\mathbf{v}_1, \mathbf{v}_1, \dots, \mathbf{v}_k)$  is the linear vector subspace that is composed of all linear combinations of  $\mathbf{v}_1, \dots, \mathbf{v}_k$ ; its dimension is at most  $\min(k, n)$ .

The idea of the general projection method is to search for a “best” approximation in the affine space  $\mathcal{V}_k = \mathbf{x}_0 + \mathcal{K}_k$ , where  $\mathcal{K}_k$  is a linear subspace spanned by  $k+1$  vectors  $\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_k$ :

$$\mathcal{K}_k = \text{span}(\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_k). \quad (2.1)$$

In order to define the “best” approximation, a general way is to require that the residual  $\mathbf{r}_k$  is orthogonal to another linear space  $\mathcal{L}_k$  spanned by a set of  $k+1$  vectors  $\mathbf{w}_0, \mathbf{w}_1, \dots, \mathbf{w}_k$ . Note that if  $\mathcal{L}_k = A\mathcal{K}_k$  this requirement is equivalent to the minimization problem:

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x} \in \mathbf{x}_0 + \mathcal{K}_k} \|\mathbf{b} - A\mathbf{x}\|_2.$$

Such projection methods thusly involve the following components:

- How to construct the vectors  $\mathbf{v}_k$  and  $\mathbf{w}_k$ .
- How to solve the minimization problem efficiently in the  $k$ -th iteration.

A more extensive discussion will be provided in the next lecture on Krylov space methods. For now let us consider an alternative to define the “best” approximation as finding the minimum of  $\phi(\mathbf{x})$  on the affine space  $\mathbf{x}_0 + \mathcal{K}_k$ .

### 3 The Conjugate Gradient Method

Suppose the search spaces  $\mathcal{K}_k$  are spanned by some linearly-independent vectors  $\mathbf{p}_0, \dots, \mathbf{p}_k$ , we denote  $P_k \in \mathbb{R}^{n \times (k+1)}$  as  $P_k = [\mathbf{p}_0 \ \mathbf{p}_1 \ \dots \ \mathbf{p}_k]$ . Hence we may write  $\mathbf{x}_{k+1} = \mathbf{x}_0 + P_k \mathbf{y}_k$  for some coefficient vector  $\mathbf{y}_k \in \mathbb{R}^{k+1}$ . For convenience we also write  $\mathcal{K}_{-1} = \emptyset$  so that  $\mathbf{x}_0$  minimize  $\phi(\mathbf{x})$  in the affine space  $\mathbf{x}_0 + \mathcal{K}_{-1} = \{\mathbf{x}_0\}$ .

Now let  $\mathbf{x}_k$  and  $\mathbf{p}_0, \dots, \mathbf{p}_k$  have already been computed, we consider minimizing  $\phi(\mathbf{x})$  over the space  $\mathbf{x}_0 + \mathcal{K}_k$  and write for now the solution as:

$$\mathbf{x}_{k+1} = \mathbf{x}_0 + P_{k-1} \mathbf{y} + \alpha \mathbf{p}_k, \quad (3.1)$$

for some  $\mathbf{y} \in \mathbb{R}^k$  and  $\alpha \in \mathbb{R}$ . The hope is that  $\mathbf{y} = \mathbf{y}_{k-1}$  so that the minimization problem is equivalent to the line search along the direction  $\mathbf{p}_k$ , an easy sub-problem to solve. To this end we assume (3.1) and obtain

$$\phi(\mathbf{x}_{k+1}) = \phi(\mathbf{x}_0 + P_{k-1} \mathbf{y}) + \alpha \mathbf{y}^t P_{k-1}^t A \mathbf{p}_k + \frac{\alpha^2}{2} \mathbf{p}_k^t A \mathbf{p}_k - \alpha \mathbf{p}_k^t \mathbf{r}_0.$$

If  $\mathbf{p}_k$  is chosen such that  $P_{k-1}^t A \mathbf{p}_k = 0$ , i.e.,  $\mathbf{p}_j^t A \mathbf{p}_k = 0$  for all  $j = 0, \dots, k-1$ , the second term vanishes and the contribution of  $\mathbf{y}$  is contained in the first term. An equivalent statement of this condition is that  $\mathbf{p}_k$  is *A-conjugate* to the subspace  $\mathcal{K}_{k-1}$ .

Assuming this condition is satisfied, we thusly have:

$$\phi(\mathbf{x}_{k+1}) = \phi(\mathbf{x}_0 + P_{k-1} \mathbf{y}) + \frac{\alpha^2}{2} \mathbf{p}_k^t A \mathbf{p}_k - \alpha \mathbf{p}_k^t \mathbf{r}_0,$$

where the minimization over  $(\mathbf{y}, \alpha)$  is decoupled into minimizing over  $\mathbf{y}$  and over  $\alpha$  separately. The first part leads to  $\mathbf{y} = \mathbf{y}_{k-1}$  as desired; and the second part leads to:

$$\alpha_k = \frac{\mathbf{p}_k^t \mathbf{r}_0}{\mathbf{p}_k^t A \mathbf{p}_k} = \frac{\mathbf{p}_k^t \mathbf{r}_k}{\mathbf{p}_k^t A \mathbf{p}_k}. \quad (3.2)$$

Here the second equality comes from:

$$\mathbf{p}_k^t \mathbf{r}_k = \mathbf{p}_k^t (\mathbf{b} - A(\mathbf{x}_0 + P_{k-1} \mathbf{y}_{k-1})) = \mathbf{p}_k^t \mathbf{r}_0 - (\mathbf{p}_k^t A P_{k-1}) \mathbf{y}_{k-1} = \mathbf{p}_k^t \mathbf{r}_0; \quad (3.3)$$

and one can compute the next iteration as  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$ .

Now the problem reduces to finding  $\mathbf{p}_k$  such that it is *A-conjugate* to  $\mathcal{K}_{k-1}$ . Let us denote this requirement as  $\mathbf{p}_k \in (A\mathcal{K}_{k-1})^\perp$ , where  $A\mathcal{K}_{k-1}$  is the linear space obtained by pre-multiplying any member of  $\mathcal{K}_{k-1}$  by the matrix  $A$  and the superscript  $\perp$  denotes the orthogonal (w.r.t. the usual Euclidean norm) complement. First, we want to know whenever the current residual  $\mathbf{r}_k$  is non-zero, the next non-zero search direction  $\mathbf{p}_k \in (A\mathcal{K}_{k-1})^\perp$  so that  $\mathbf{p}_k^t \mathbf{r}_k \neq 0$  (hence  $\alpha_k \neq 0$ ) can always be found.

To this end, let us use the method of induction and first assume  $\mathbf{r}_0 \neq 0$ . Note that  $\mathcal{K}_{-1} = \emptyset$  thus  $(A\mathcal{K}_{-1})^\perp = \mathbb{R}^n$ , and we can simply set  $\mathbf{p}_0 = \mathbf{r}_0 \neq 0$ . Now suppose we work up to the point that  $\mathbf{r}_k \neq 0$  and want to find a suitable  $\mathbf{p}_k$ . Note that  $\mathbf{x} = A^{-1}\mathbf{b}$ , the exact solution, does not belong to  $\mathbf{x}_0 + \mathcal{K}_{k-1}$  since  $\mathbf{r}_k \neq 0$ , we have:

$$A^{-1}\mathbf{b} \notin \mathbf{x}_0 + \mathcal{K}_{k-1} \Rightarrow \mathbf{r}_0 = \mathbf{b} - A\mathbf{x}_0 \notin A\mathcal{K}_{k-1}. \quad (3.4)$$

This means that we can find  $\mathbf{p}_k \in (A\mathcal{K}_{k-1})^\perp$  such that  $\mathbf{p}_k^t \mathbf{r}_0 \neq 0$ ; and  $\mathbf{p}_k^t \mathbf{r}_k \neq 0$  follows from the fact that  $\mathbf{p}_k \in (A\mathcal{K}_{k-1})^\perp$  and a similar argument using (3.3).

From the previous proof, we know that if  $\mathbf{r}_k \neq 0$  then  $\mathbf{r}_k \notin A\mathcal{K}_{k-1}^\perp$ . Following this observation, the *conjugate gradient method* finds  $\mathbf{p}_k$  by removing from  $\mathbf{r}_k$  the latter's projection onto  $(A\mathcal{K}_{k-1})^\perp$ . That is, we write:

$$\mathbf{r}_k = \gamma_k \mathbf{p}_k + \mathbf{z}_k, \quad \mathbf{z}_k \in A\mathcal{K}_{k-1}, \quad \mathbf{p}_k \in (A\mathcal{K}_{k-1})^\perp \text{ and } \gamma_k \neq 0. \quad (3.5)$$

More details about this decomposition as well as an explicit formula to compute  $\gamma_k \mathbf{p}_k$  and  $\mathbf{z}_k$  from a basis of  $A\mathcal{K}_{k-1}$  is described in Exercise 1.

Because  $\mathbf{r}_k - \mathbf{r}_0 = A(\mathbf{x}_0 - \mathbf{x}_k) \in A\mathcal{K}_{k-1}$ , it is not difficult to see that the orthogonal decomposition of  $\mathbf{r}_0$  w.r.t.  $A\mathcal{K}_{k-1}$  is given by:

$$\mathbf{r}_0 = \gamma_k \mathbf{p}_k + (\mathbf{z}_k + \mathbf{r}_0 - \mathbf{r}_k).$$

Now let  $\mathbf{p}_0 = \mathbf{r}_0$ , we can then use induction to show that:

$$\mathcal{K}_k = \text{span}(\mathbf{r}_0, A\mathbf{r}_0, \dots, A^k \mathbf{r}_0). \quad (3.6)$$

Indeed, the choice of  $\mathbf{p}_0$  indicates  $\mathcal{K}_0 = \text{span}(\mathbf{r}_0)$  and now we assume (3.6) is true for  $k-1$  ( $k \geq 1$ ), then we have:

$$\mathbf{p}_k = \frac{1}{\gamma_k} [\mathbf{r}_0 - (\mathbf{z}_k + \mathbf{r}_0 - \mathbf{r}_k)] \in \text{span}(\mathbf{r}_0, A\mathcal{K}_{k-1}) = \text{span}(\mathbf{r}_0, A\mathbf{r}_0, \dots, A^k \mathbf{r}_0),$$

and (3.6) for the case  $k$  follows immediately.

Next let's take a closer look at  $\mathbf{r}_k = \mathbf{b} - A\mathbf{x}_k = \mathbf{b} - A(\mathbf{x}_0 + P_{k-1}\mathbf{y}_{k-1})$ , where  $\mathbf{y}_{k-1}$  minimize  $\phi(\mathbf{x}_0 + P_{k-1}\mathbf{y})$  over all  $\mathbf{y} \in \mathbb{R}^k$ . The solution to this minimization problem is not difficult to compute:

$$\mathbf{y}_{k-1} = (P_{k-1}^t A P_{k-1})^{-1} P_{k-1}^t A \mathbf{r}_0.$$

It is then computed that:

$$P_{k-1}^t \mathbf{r}_k = P_{k-1}^t (\mathbf{r}_0 - A P_{k-1} \mathbf{y}_{k-1}) = 0,$$

or equivalently  $\mathbf{r}_k \in \mathcal{K}_{k-1}^\perp$ . Combining this result and (3.5), we obtain again by induction that  $\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_k$  are orthogonal to each other (hence are also linearly independent) and:

$$\mathcal{K}_k = \text{span}(\mathbf{r}_0, \mathbf{r}_1, \dots, \mathbf{r}_k). \quad (3.7)$$

There remains an important step of deriving a formula to compute  $\mathbf{z}_k$  efficiently. Let us assume for now  $k \geq 2$ . Note that  $\mathbf{z}_k \in A\mathcal{K}_{k-1}$ , one of whose basis is given by  $A P_{k-1}$ , hence:

$$\mathbf{z}_k = A P_{k-1} \mathbf{y}_{k-1}, \quad \mathbf{y}_{k-1} = \arg \min_{\mathbf{y} \in \mathbb{R}^k} \|\mathbf{r}_k - A P_{k-1} \mathbf{y}\|.$$

Let us write  $\mathbf{y} \in \mathbb{R}^k$  and  $P_{k-1}$  as

$$\mathbf{y} = \begin{bmatrix} \mathbf{w} \\ \beta \end{bmatrix}, \quad \mathbf{w} \in \mathbb{R}^{k-1}, \beta \in \mathbb{R}; \quad P_{k-1} = \begin{bmatrix} P_{k-2} & \mathbf{p}_{k-1} \end{bmatrix},$$

---

<sup>1</sup>Otherwise if  $\mathbf{r}_k \in A\mathcal{K}_{k-1}$ , then for all  $\mathbf{p} \in (A\mathcal{K}_{k-1})^\perp$  we must have  $\mathbf{p}^t \mathbf{r}_k = 0$ , contradiction.

then the vector whose 2-norm is to be minimized is the same as:

$$\mathbf{r}_k - AP_{k-1}\mathbf{y} = \mathbf{r}_k - AP_{k-2}\mathbf{w} - \beta A\mathbf{p}_{k-1}.$$

Noticing that:

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

we further compute:

$$\mathbf{r}_k - AP_{k-1}\mathbf{y} = \mathbf{r}_k - AP_{k-2}\mathbf{w} - \frac{\beta}{\alpha_{k-1}}(\mathbf{r}_{k-1} - \mathbf{r}_k) = \left(1 + \frac{\beta}{\alpha_{k-1}}\right)\mathbf{r}_k - AP_{k-2}\mathbf{w} - \frac{\beta}{\alpha_{k-1}}\mathbf{r}_{k-1}.$$

Because both  $AP_{k-2}\mathbf{w}$  and  $\mathbf{r}_{k-1}$  belongs to  $\mathcal{K}_{k-1}$ , using the previous result that  $\mathbf{r}_k \in \mathcal{K}_{k-1}^\perp$  we conclude that:

$$\|\mathbf{r}_k - AP_{k-1}\mathbf{y}\|^2 = \left(1 + \frac{\beta}{\alpha_{k-1}}\right)^2 \|\mathbf{r}_k\|^2 + \left\|AP_{k-2}\mathbf{w} + \frac{\beta}{\alpha_{k-1}}\mathbf{r}_{k-1}\right\|^2.$$

Let  $\mathbf{y}_{k-1}$  have the components  $\mathbf{w}_{k-1}$  and  $\beta_{k-1}$ , then  $\mathbf{w}_{k-1}$  solves the minimization problem:

$$-\frac{\alpha_{k-1}}{\beta_{k-1}}\mathbf{w}_{k-1} = \arg \min_{\mathbf{w} \in \mathbb{R}^{k-1}} \|\mathbf{r}_{k-1} - AP_{k-2}\mathbf{w}\|,$$

or equivalently:

$$-\frac{\alpha_{k-1}}{\beta_{k-1}}AP_{k-2}\mathbf{w}_{k-1} = \arg \min_{\mathbf{z} \in A\mathcal{K}_{k-2}} \|\mathbf{r}_{k-1} - \mathbf{z}\| = \mathbf{z}_{k-1} = \mathbf{r}_{k-1} - \gamma_{k-1}\mathbf{p}_{k-1}.$$

This leads to a formula of  $\mathbf{z}_k$ :

$$\mathbf{z}_k = AP_{k-2}\mathbf{w}_{k-1} + \beta_{k-1}A\mathbf{p}_{k-1} = -\frac{\beta_{k-1}}{\alpha_{k-1}}(\mathbf{r}_{k-1} - \gamma_{k-1}\mathbf{p}_{k-1}) + \beta_{k-1}A\mathbf{p}_{k-1} = -\frac{\beta_{k-1}}{\alpha_{k-1}}(\mathbf{r}_k - \gamma_{k-1}\mathbf{p}_{k-1}).$$

In the view of (3.5),  $\mathbf{p}_k$  is a linear combination of  $\mathbf{r}_k$  and  $\mathbf{p}_{k-1}$ . This statement is clearly also true for  $k < 2$ .

**Remark 3.** *Strictly speaking, we need to show first that  $\beta_{k-1} \neq 0$  in the preceding argument. But this is actually fairly straightforward by an argument of contradiction.*

At last, without loss of generality, we can always scale the search direction appropriately so that:

$$\mathbf{p}_k = \mathbf{r}_k + s_k\mathbf{p}_{k-1}, \quad (3.8)$$

where the scalar  $s_k$  can be evaluated by pre-multiplying both sides of (3.8) by  $\mathbf{p}_{k-1}^t A$ :

$$0 = \mathbf{p}_{k-1}^t A\mathbf{p}_k = \mathbf{p}_{k-1}^t A\mathbf{r}_k + s_k\mathbf{p}_{k-1}^t A\mathbf{p}_{k-1},$$

here the first identity is due to the fact that  $\mathbf{p}_k$  is  $A$ -conjugate to  $\mathcal{K}_{k-1}$ . Thus:

$$s_k = -\frac{\mathbf{p}_{k-1}^t A\mathbf{r}_k}{\mathbf{p}_{k-1}^t A\mathbf{p}_{k-1}}. \quad (3.9)$$

In the end, the detailed algorithm for the conjugate gradient method is given by the following:

$$\begin{aligned}
& \mathbf{r}_0 = \mathbf{b} - A\mathbf{x}_0, \quad k = 0; \\
& \text{while } \mathbf{r}_k \neq 0 : \\
& \quad \text{if } k = 0 \\
& \quad \quad \mathbf{p}_0 = \mathbf{r}_0; \\
& \quad \text{else} \\
& \quad \quad s_k = -\mathbf{p}_{k-1}^t A \mathbf{r}_k / \mathbf{p}_{k-1}^t A \mathbf{p}_{k-1}; \\
& \quad \quad \mathbf{p}_k = \mathbf{r}_k + s_k \mathbf{p}_{k-1}; \\
& \quad \text{end} \\
& \quad \alpha_k = \mathbf{p}_k^t \mathbf{r}_k / \mathbf{p}_k^t A \mathbf{p}_k; \\
& \quad \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k; \\
& \quad \mathbf{r}_{k+1} = \mathbf{b} - A\mathbf{x}_{k+1}; \\
& \quad k = k + 1; \\
& \text{end} \\
& \mathbf{x} = \mathbf{x}_k.
\end{aligned} \tag{3.10}$$

## 4 Further Analysis

There are several steps of the algorithm (3.10) that can be improved for computational efficiency. In particular, there are four matrix-vector multiplications in each iteration, one of which can be re-used between adjacent loops; hence algorithm (3.10) has in average three matrix-vector multiplications per iteration. This is the major computational cost with the method.

Now let us derive an equivalent method that only requires one matrix-vector multiplication per iteration. Particularly, pre-multiplying  $\mathbf{r}_k = \mathbf{r}_{k-1} - \alpha_{k-1} A \mathbf{p}_{k-1}$  by  $\mathbf{r}_k^t$  we obtain:

$$\mathbf{r}_k^t \mathbf{r}_k = -\alpha_{k-1} \mathbf{r}_k^t A \mathbf{p}_{k-1};$$

and pre-multiplying by  $\mathbf{p}_{k-1}^t = \mathbf{r}_{k-1}^t - s_{k-1} \mathbf{p}_{k-2}^t$  we obtain:

$$0 = \mathbf{p}_{k-1}^t \mathbf{r}_k = (\mathbf{r}_{k-1}^t - s_{k-1} \mathbf{p}_{k-2}^t) \mathbf{r}_{k-1} - \alpha_{k-1} \mathbf{p}_{k-1}^t A \mathbf{p}_{k-1} = \mathbf{r}_{k-1}^t \mathbf{r}_{k-1} - \alpha_{k-1} \mathbf{p}_{k-1}^t A \mathbf{p}_{k-1}.$$

Thus we can compute  $s_k$  and  $\alpha_k$  instead as:

$$s_k = -\frac{\mathbf{p}_{k-1}^t A \mathbf{r}_k}{\mathbf{p}_{k-1}^t A \mathbf{p}_{k-1}} = \frac{\mathbf{r}_k^t \mathbf{r}_k}{\mathbf{r}_{k-1}^t \mathbf{r}_{k-1}}, \quad \alpha_k = \frac{\mathbf{r}_k^t \mathbf{r}_k}{\mathbf{p}_k^t A \mathbf{p}_k}. \tag{4.1}$$

Finally, the next residual is updated from the previous one by:

$$\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k A \mathbf{p}_k, \tag{4.2}$$

where the product of  $A \mathbf{p}_k$  is already available in computing  $\alpha_k$ . To this end, the following version

of the method only requires one matrix-vector multiplication per iteration:

$$\begin{aligned}
& \mathbf{r}_0 = \mathbf{b} - A\mathbf{x}_0, \quad k = 0; \\
& \text{while } \mathbf{r}_k \neq 0 : \\
& \quad \text{if } k = 0 \\
& \quad \quad \mathbf{p}_0 = \mathbf{r}_0; \\
& \quad \text{else} \\
& \quad \quad s_k = \mathbf{r}_k^t \mathbf{r}_k / \mathbf{r}_{k-1}^t \mathbf{r}_{k-1}; \\
& \quad \quad \mathbf{p}_k = \mathbf{r}_k + s_k \mathbf{p}_{k-1}; \\
& \quad \text{end} \\
& \quad \alpha_k = \mathbf{r}_k^t \mathbf{r}_k / \mathbf{p}_k^t A \mathbf{p}_k; \\
& \quad \mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k; \\
& \quad \mathbf{r}_{k+1} = \mathbf{r}_k - \alpha_k A \mathbf{p}_k; \\
& \quad k = k + 1; \\
& \text{end} \\
& \mathbf{x} = \mathbf{x}_k.
\end{aligned} \tag{4.3}$$

Using the structure of  $\mathcal{K}_k$ , we see that with exact arithmetics the algorithm (4.3) converges in at most  $n$  steps; and if  $A = I + B$  such that  $\text{rank} B = r$ , it converges in at most  $r + 1$  steps.

Finally, we state without proof an estimate on the error of the conjugate gradient method:

$$\|\mathbf{x} - \mathbf{x}_k\|_A \leq 2 \|\mathbf{x} - \mathbf{x}_0\|_A \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k, \tag{4.4}$$

where  $\kappa = \kappa_2(A)$  is the condition number of  $A$  in the induced  $L^2$ -norm. More details can be found in David G. Luenberger's work [1].

## Exercises

**Exercise 1.** Let  $\mathcal{K}$  be a  $k$ -dimensional subspace of  $\mathbb{R}^n$  with a basis  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$ . Then the orthogonal complement of  $\mathcal{K}$  is defined as:

$$\mathcal{K}^\perp = \{\mathbf{w} \in \mathbb{R}^n : \mathbf{w}^t \mathbf{v} = 0 \quad \forall \mathbf{v} \in \mathcal{K}\}. \tag{4.5}$$

Let  $\mathbf{u} \in \mathbb{R}^n$  be arbitrary, its **orthogonal decomposition** with respect to  $\mathcal{K}$  is:

$$\mathbf{u} = \mathbf{w} + \mathbf{v}, \quad \text{such that } \mathbf{w} \in \mathcal{K}^\perp \text{ and } \mathbf{v} \in \mathcal{K}. \tag{4.6}$$

(i) Show that this decomposition is unique.

Usually we denote this unique  $\mathbf{v}$  by  $\mathcal{P}\mathbf{u}$  and  $\mathbf{w}$  by  $(\mathcal{I} - \mathcal{P})\mathbf{u}$  where  $\mathcal{P} : \mathbb{R}^n \rightarrow \mathcal{K}$  is known as the projection of  $\mathbb{R}^n$  onto  $\mathcal{K}$  and  $\mathcal{I} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is the identity map.

(ii) Prove that  $\mathbf{v} = \mathcal{P}\mathbf{u}$  minimizes the distance between any vector of  $\mathcal{K}$  and  $\mathbf{u}$ , i.e.,

$$\mathbf{v} = \mathcal{P}\mathbf{u} = \arg \min_{\mathbf{v}' \in \mathcal{K}} \|\mathbf{v}' - \mathbf{u}\|. \tag{4.7}$$

**Hint:** We just need to show that  $\|\mathbf{v}' - \mathbf{u}\| \geq \|\mathbf{v} - \mathbf{u}\|$  for all  $\mathbf{v}' \in \mathcal{K}$ . This is rewritten as  $\|(\mathbf{v}' - \mathbf{v}) - \mathbf{w}\| \geq \|\mathbf{w}\|$ . To proceed, use  $\mathbf{w} \in \mathcal{K}^\perp$  to show that  $\|(\mathbf{v}' - \mathbf{v}) - \mathbf{w}\|^2 = \|\mathbf{v}' - \mathbf{v}\|^2 + \|\mathbf{w}\|^2$ .

Now let  $V = [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_k] \in \mathbb{R}^{n \times k}$ , then any vector of  $\mathcal{K}$  can be written as  $V\mathbf{y}$  with some  $\mathbf{y} \in \mathbb{R}^k$ .  
 (iii) Use this representation and (4.7) to show:

$$\mathcal{P}\mathbf{u} = V(V^t V)^{-1} V^t \mathbf{u}. \quad (4.8)$$

## References

- [1] David G. Luenberger. Introduction to Linear and Nonlinear Programming. Addison Wesley Publishing Company, 1973.