# 5 CONJUGATE DIRETION METHODS

## 5.1 INTRODUCTION

- The class of conjugate direction methods can be viewed as being intermediate between the method of steepest descent and Newton's method. The conjugate direction methods have the following properties:

    1. Solve quadratics of $n$ variables in $n$ steps.
    2. The usual implementation, the conjugate gradient algorithm, requires no Hessian matrix evaluations.
    3. No matrix inversion and no storage of an $n \times n$ matrix are required.

- The conjugate direction methods typically perform better than the method of steepest descent, but not as well as Newton's method.

- **Definition 5.1.1** Let $\boldsymbol{Q}$ be a real symmetric $n \times n$ matrix. The directions $\boldsymbol{d}^{(0)}, \boldsymbol{d}^{(1)}, \boldsymbol{d}^{(2)}, \dots, \boldsymbol{d}^{(m)}$ are $\boldsymbol{Q}$-conjugate if for all $i \neq j$, we have $\boldsymbol{d}^{(i)\top} \boldsymbol{Q} \boldsymbol{d}^{(j)} = 0$.

- For a quadratic function of $n$ variables $f(\boldsymbol{x}) = (1/2)\boldsymbol{x}^\top \boldsymbol{Q} \boldsymbol{x} - \boldsymbol{x}^\top \boldsymbol{b}$, $\boldsymbol{x} \in \mathbb{R}^n$, $\boldsymbol{Q} = \boldsymbol{Q}^\top > 0$, the best direction of search is in the $\boldsymbol{Q}$-conjugate direction.

- **Lemma 5.1.1** Let $\boldsymbol{Q}$ be a symmetric positive definite $n \times n$ matrix. If the directions $\boldsymbol{d}^{(0)}, \boldsymbol{d}^{(1)}, \boldsymbol{d}^{(2)}, \dots, \boldsymbol{d}^{(m)} \in \mathbb{R}^n$, $k \leq n - 1$, are nonzero and $\boldsymbol{Q}$-conjugate, then they are linearly independent.

    *Proof.* Let $\alpha_0, \dots, \alpha_k$ be scalars such that

    $$\alpha_0 \boldsymbol{d}^{(0)} + \alpha_1 \boldsymbol{d}^{(1)} + \cdots + \alpha_k \boldsymbol{d}^{(k)} = \boldsymbol{0}.$$

    Premultiplying the equality by $\boldsymbol{d}^{(j)\top} \boldsymbol{Q}$, $0 \leq j \leq k$, yields

    $$\boldsymbol{d}^{(j)\top} \boldsymbol{Q} \alpha_0 \boldsymbol{d}^{(j)\top} \boldsymbol{Q} \boldsymbol{d}^{(0)} + \cdots + \alpha_j \boldsymbol{d}^{(j)\top} \boldsymbol{Q} \boldsymbol{d}^{(j)} + \cdots + \alpha_k \boldsymbol{d}^{(j)\top} \boldsymbol{Q} \boldsymbol{d}^{(k)} = 0$$
    $$\alpha_j \boldsymbol{d}^{(j)\top} \boldsymbol{Q} \boldsymbol{d}^{(j)} = 0,$$

    because all other terms $\boldsymbol{d}^{(j)\top} \boldsymbol{Q} \boldsymbol{d}^{(i)} = 0$, $i \neq j$, by $\boldsymbol{Q}$-conjugacy. But $\boldsymbol{Q} = \boldsymbol{Q}^\top > 0$ and $\boldsymbol{d}^{(j)} \neq \boldsymbol{0}$; hence $\alpha_j = 0$, $j = 0, 1, \dots, k$. Therefore, $\boldsymbol{d}^{(0)}, \boldsymbol{d}^{(1)}, \boldsymbol{d}^{(2)}, \dots, \boldsymbol{d}^{(k)}$, $k \leq n - 1$, are linearly independent.

**Example 5.1.1.** Let

$$\boldsymbol{Q} = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 4 & 2 \\ 1 & 2 & 3 \end{bmatrix}.$$

(a) Show that $\boldsymbol{Q}$ is symmetric positive definite.

(b) Construct a set of $\boldsymbol{Q}$-conjugate vectors $\boldsymbol{d}^{(0)}, \boldsymbol{d}^{(1)}, \boldsymbol{d}^{(2)}$.

Chong, Edwin Kah Pin., and Stanislaw H. Żak. *An Introduction to Optimization.* Fourth edition., John Wiley & Sons, Inc., 2013.

(a) Note that $\boldsymbol{Q} = \boldsymbol{Q}^\top > 0$. The matrix $\boldsymbol{Q}$ is positive definite because all its leading principal minors are positive:

$$\Delta_1 = 3 > 0,$$

$$\Delta_2 = \det \begin{bmatrix} 3 & 0 \\ 0 & 4 \end{bmatrix} = 12 > 0,$$

$$\Delta_3 = \det \boldsymbol{Q} = 20 > 0.$$

(b) Let $\boldsymbol{d}^{(0)} = [1,0,0]^\top$, $\boldsymbol{d}^{(1)} = \left[d_1^{(1)}, d_2^{(1)}, d_3^{(1)}\right]^\top$, $\boldsymbol{d}^{(2)} = \left[d_1^{(2)}, d_2^{(2)}, d_3^{(2)}\right]^\top$. We require that $\boldsymbol{d}^{(0)\top}\boldsymbol{Q}\boldsymbol{d}^{(1)} = 0$. We have

$$\boldsymbol{d}^{(0)\top}\boldsymbol{Q}\boldsymbol{d}^{(1)} = [1,0,0] \begin{bmatrix} 3 & 0 & 1 \\ 0 & 4 & 2 \\ 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} d_1^{(1)} \\ d_2^{(1)} \\ d_3^{(1)} \end{bmatrix} = 3d_1^{(1)} + d_3^{(1)}.$$

Let $d_1^{(1)} = 1$, $d_2^{(1)} = 0$, $d_3^{(1)} = -3$. Then, $\boldsymbol{d}^{(1)} = [1,0,-3]^\top$, and thus $\boldsymbol{d}^{(0)\top}\boldsymbol{Q}\boldsymbol{d}^{(1)} = 0$.

To find the third vector $\boldsymbol{d}^{(2)}$, which would be $\boldsymbol{Q}$-conjugate with $\boldsymbol{d}^{(0)}$ and $\boldsymbol{d}^{(1)}$, we require that

$$\boldsymbol{d}^{(0)\top}\boldsymbol{Q}\boldsymbol{d}^{(2)} = [1,0,0] \begin{bmatrix} 3 & 0 & 1 \\ 0 & 4 & 2 \\ 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} d_1^{(2)} \\ d_2^{(2)} \\ d_3^{(2)} \end{bmatrix} = 3d_1^{(2)} + d_3^{(2)} = 0,$$

$$\boldsymbol{d}^{(1)\top}\boldsymbol{Q}\boldsymbol{d}^{(2)} = [1,0,-3] \begin{bmatrix} 3 & 0 & 1 \\ 0 & 4 & 2 \\ 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} d_1^{(2)} \\ d_2^{(2)} \\ d_3^{(2)} \end{bmatrix} = -6d_2^{(2)} - 8d_3^{(2)} = 0.$$

If we take $\boldsymbol{d}^{(2)} = [1,4,-3]^\top$, then the resulting set of vectors is mutually conjugate.

## 5.2 THE CONJUGATE DIRECTION ALGORITHM

- We now present the conjugate direction algorithm for minimizing the quadratic function of $n$ variables

$$f(\boldsymbol{x}) = \frac{1}{2}\boldsymbol{x}^\top\boldsymbol{Q}\boldsymbol{x} - \boldsymbol{x}^\top\boldsymbol{b},$$

  where $\boldsymbol{Q} = \boldsymbol{Q}^\top > 0$, $x \in \mathbb{R}^n$. Note that because $\boldsymbol{Q} > 0$, the function $f$ has a global minimizer that can be found by solving $\boldsymbol{Q}\boldsymbol{x} = \boldsymbol{b}$.

- **Basic Conjugate Direction Algorithm** Given a starting point $\boldsymbol{x}^{(0)}$ and $\boldsymbol{Q}$-conjugate directions $\boldsymbol{d}^{(0)}, \boldsymbol{d}^{(1)}, \boldsymbol{d}^{(2)}, \dots, \boldsymbol{d}^{(n-1)}$; for $k \geq 0$,

$$g^{(k)} = \nabla f\big(x^{(k)}\big) = Qx^{(k)} - b,$$
$$\alpha_k = -\frac{g^{(k)\top} d^{(k)}}{d^{(k)\top} Q d^{(k)}},$$
$$x^{(k+1)} = x^{(k)} + \alpha_k d^{(k)}.$$

- **Theorem 5.2.1** For any starting point $x^{(0)}$ the basic conjugate direction algorithm converges to the unique $x^*$ (that solves $Qx = b$) in $n$ steps; that is, $x^{(n)} = x^*$.

  *Proof.* Consider $x^* - x^{(0)} \in \mathbb{R}^n$. Because $d^{(i)}$ are linearly independent, there exist constants $\beta_i$, $i = 0, \dots, n-1$, such that

  $$x^* - x^{(0)} = \beta_0 d^{(0)} + \cdots + \beta_{n-1} d^{(n-1)}.$$

  Now premultiply both sides of this equation by $d^{(k)\top} Q$, $0 \le k < n$, to obtain

  $$d^{(k)\top} Q\big(x^* - x^{(0)}\big) = \beta_k d^{(k)\top} Q d^{(k)},$$

  where the terms $d^{(k)\top} Q d^{(i)} = 0$, $k \ne i$, by the $Q$-conjugate property. Hence,

  $$\beta_k = \frac{d^{(k)\top} Q\big(x^* - x^{(0)}\big)}{d^{(k)\top} Q d^{(k)}}.$$

  Note that

  $$x^{(k)} = x^{(k-1)} + \alpha_{k-1} d^{(k-1)}$$
  $$x^{(k)} = x^{(k-2)} + \alpha_{k-2} d^{(k-2)} + \alpha_{k-1} d^{(k-1)}$$
  $$x^{(k)} = x^{(0)} + \alpha_0 d^{(0)} + \cdots + \alpha_{k-2} d^{(k-2)} + \alpha_{k-1} d^{(k-1)}$$
  $$x^{(k)} - x^{(0)} = \alpha_0 d^{(0)} + \cdots + \alpha_{k-2} d^{(k-2)} + \alpha_{k-1} d^{(k-1)}$$

  Now, writing
  $$x^* - x^{(0)} = \big(x^* - x^{(k)}\big) + \big(x^{(k)} - x^{(0)}\big)$$

  and premultiplying the above by $d^{(k)\top} Q$, we obtain

  $$\begin{aligned}
  d^{(k)\top} Q\big(x^* - x^{(0)}\big) &= d^{(k)\top} Q\big(x^* - x^{(k)}\big) + d^{(k)\top} Q\big(x^{(k)} - x^{(0)}\big) \\
  &= d^{(k)\top} Q\big(x^* - x^{(k)}\big) + d^{(k)\top} Q\big(\alpha_0 d^{(0)} + \cdots + \alpha_{k-1} d^{(k-1)}\big) \\
  &= d^{(k)\top} Q\big(x^* - x^{(k)}\big) \\
  &= d^{(k)\top}\big(Qx^* - Qx^{(k)}\big) \\
  &= d^{(k)\top}\big(b - Qx^{(k)}\big) \\
  &= -d^{(k)\top} g^{(k)}
  \end{aligned}$$

  because $g^{(k)} = Qx^{(k)} - b$ and $Qx^* = b$. Thus,

  $$\beta_k = \frac{d^{(k)\top} Q\big(x^* - x^{(0)}\big)}{d^{(k)\top} Q d^{(k)}} = -\frac{d^{(k)\top} g^{(k)}}{d^{(k)\top} Q d^{(k)}} = \alpha_k$$

  and $x^* = x^{(n)}$, which completes the proof.

**Example 5.2.1** Find the minimizer of

$$f(x_1, x_2) = \frac{1}{2}\boldsymbol{x}^\top \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \boldsymbol{x} - \boldsymbol{x}^\top \begin{bmatrix} -1 \\ 1 \end{bmatrix}, \boldsymbol{x} \in \mathbb{R}^2,$$

using the conjugate direction method with the initial point $\boldsymbol{x}^{(0)} = [0,0]^\top$, and $\boldsymbol{Q}$-conjugate directions $\boldsymbol{d}^{(0)} = [1,0]^\top$ and $\boldsymbol{d}^{(1)} = [-3/8, 3/4]^\top$.

Note that the optimal is $\boldsymbol{x}^* = \boldsymbol{Q}^{-1}\boldsymbol{b} = \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix}^{-1} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1/2 & -1/2 \\ -1/2 & 1 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ 3/2 \end{bmatrix}.$

To begin, we have

$$\boldsymbol{g}^{(0)} = -\boldsymbol{b} = [1, -1]^\top,$$

and hence

$$\alpha_0 = -\frac{\boldsymbol{g}^{(0)\top}\boldsymbol{d}^{(0)}}{\boldsymbol{d}^{(0)\top}\boldsymbol{Q}\boldsymbol{d}^{(0)}} = -\frac{[1,-1]\begin{bmatrix} 1 \\ 0 \end{bmatrix}}{[1,0]\begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix}\begin{bmatrix} 1 \\ 0 \end{bmatrix}} = -\frac{1}{4}.$$

Thus,

$$\boldsymbol{x}^{(1)} = \boldsymbol{x}^{(0)} + \alpha_0 \boldsymbol{d}^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \frac{1}{4}\begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -1/4 \\ 0 \end{bmatrix}.$$

To find $\boldsymbol{x}^{(2)}$, we compute

$$\boldsymbol{g}^{(1)} = \boldsymbol{Q}\boldsymbol{x}^{(1)} - \boldsymbol{b} = \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix}\begin{bmatrix} -1/4 \\ 0 \end{bmatrix} - \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ -3/2 \end{bmatrix}$$

and

$$\alpha_1 = -\frac{\boldsymbol{g}^{(1)\top}\boldsymbol{d}^{(1)}}{\boldsymbol{d}^{(1)\top}\boldsymbol{Q}\boldsymbol{d}^{(1)}} = -\frac{\left[0, -\frac{3}{2}\right]\begin{bmatrix} -3/8 \\ 3/4 \end{bmatrix}}{\left[-\frac{3}{8}, \frac{3}{4}\right]\begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix}\begin{bmatrix} -3/8 \\ 3/4 \end{bmatrix}} = 2.$$

Therefore,

$$\boldsymbol{x}^{(2)} = \boldsymbol{x}^{(1)} + \alpha_0 \boldsymbol{d}^{(1)} = \begin{bmatrix} -1/4 \\ 0 \end{bmatrix} + 2\begin{bmatrix} -3/8 \\ 3/4 \end{bmatrix} = \begin{bmatrix} -1 \\ 3/2 \end{bmatrix}.$$

Because $f$ is a quadratic function in two variables, $\boldsymbol{x}^{(2)} = \boldsymbol{x}^*$.

**Example 5.2.2.** Show that $\alpha_k = \arg\min f\big(x^{(k)} + \alpha d^{(k)}\big)$.

Note that $\phi_k(\alpha) = f\big(x^{(k)} + \alpha d^{(k)}\big)$ is a quadratic function in $\alpha$:

$$f\big(x^{(k)} + \alpha d^{(k)}\big) = \frac{1}{2}\big(x^{(k)} + \alpha d^{(k)}\big)^{\top} Q\big(x^{(k)} + \alpha d^{(k)}\big) - \big(x^{(k)} + \alpha d^{(k)}\big)^{\top} b$$

$$= \frac{1}{2}\alpha^2 d^{(k)\top} Q d^{(k)} + \alpha c_1 + c_2.$$

In addition, the coefficient of the $\alpha^2$ term in $\phi_k$ is $d^{(k)\top} Q d^{(k)} > 0$. Hence, $\phi_k'(\alpha_k) = 0$ implies $\alpha_k$ is a minimizer of $\phi_k'(\cdot)$.

Now, apply the chain rule to get

$$\frac{d\phi_k}{d\alpha}(\alpha_k) = \nabla f\big(x^{(k)} + \alpha_k d^{(k)}\big)^{\top} d^{(k)}$$

$$= \nabla f\big(x^{(k+1)}\big)^{\top} d^{(k)}$$

$$= \big(Q x^{(k+1)} - b\big)^{\top} d^{(k)}$$

$$= g^{(k+1)\top} d^{(k)}.$$

Note that $g^{(k+1)\top} d^{(k)} = 0$. To see this,

$$g^{(k+1)\top} d^{(k)} = \big(Q x^{(k+1)} - b\big)^{\top} d^{(k)}$$

$$= x^{(k+1)\top} Q^{\top} d^{(k)} - b^{\top} d^{(k)}$$

$$= \left(x^{(k)} - \left(\frac{g^{(k)\top} d^{(k)}}{d^{(k)\top} Q d^{(k)}}\right) d^{(k)}\right)^{\top} Q^{\top} d^{(k)} - b^{\top} d^{(k)}$$

$$= x^{(k)\top} Q d^{(k)} - \left(\frac{g^{(k)\top} d^{(k)}}{d^{(k)\top} Q d^{(k)}}\right) d^{(k)\top} Q d^{(k)} - b^{\top} d^{(k)}$$

$$= \big(Q x^{(k)} - b\big)^{\top} d^{(k)} - \left(\frac{g^{(k)\top} d^{(k)}}{d^{(k)\top} Q d^{(k)}}\right) d^{(k)\top} Q d^{(k)}$$

$$= g^{(k)\top} d^{(k)} - g^{(k)\top} d^{(k)}$$

$$= 0$$

Hence, $\alpha_k = \arg\min f\big(x^{(k)} + \alpha d^{(k)}\big)$.

## 5.3 THE CONJUGATE GRADIENT ALGORITHM

- The conjugate gradient algorithm does not use prespecified conjugate directions, but instead computes the directions as a linear combination, in such a way that all the directions are mutually $Q$-conjugate.

- Consider the quadratic function

$$f(x) = \frac{1}{2}x^{\top} Q x - x^{\top} b, \qquad x \in \mathbb{R}^n,$$

where $\boldsymbol{Q} = \boldsymbol{Q}^\top > 0$. Our first search direction from an initial point $\boldsymbol{x}^{(0)}$ is in the direction of steepest descent; that is,

$$\boldsymbol{d}^{(0)} = -\boldsymbol{g}^{(0)} = -\nabla f\big(\boldsymbol{x}^{(0)}\big) = -\boldsymbol{Q}\boldsymbol{x}^{(0)} + \boldsymbol{b}.$$

Thus,

$$\boldsymbol{x}^{(1)} = \boldsymbol{x}^{(0)} + \alpha_0 \boldsymbol{d}^{(0)},$$

where

$$\alpha_0 = \arg\min_{\alpha \geq 0} f\big(\boldsymbol{x}^{(0)} + \alpha_0 \boldsymbol{d}^{(0)}\big) = -\frac{\boldsymbol{g}^{(0)\top}\boldsymbol{d}^{(0)}}{\boldsymbol{d}^{(0)\top}\boldsymbol{Q}\boldsymbol{d}^{(0)}}.$$

In the next stage, we search in a direction $\boldsymbol{d}^{(1)}$ that is $\boldsymbol{Q}$-conjugate to $\boldsymbol{d}^{(0)}$. We choose $\boldsymbol{d}^{(1)}$ as a linear combination of $\boldsymbol{g}^{(1)}$ and $\boldsymbol{d}^{(0)}$. In general, at the $(k+1)$th step, we choose $\boldsymbol{d}^{(k+1)}$ to be a linear combination of $\boldsymbol{g}^{(k+1)}$ and $\boldsymbol{d}^{(k)}$. Specifically, we choose

$$\boldsymbol{d}^{(k+1)} = -\boldsymbol{g}^{(k+1)} + \beta_k \boldsymbol{d}^{(k)}, \quad k = 0, 1, 2, \dots$$

The coefficients $\beta_k$, $k = 1, 2, \dots$, are chosen in such a way that $\boldsymbol{d}^{(k+1)}$ is $\boldsymbol{Q}$-conjugate to $\boldsymbol{d}^{(0)}, \boldsymbol{d}^{(1)}, \dots, \boldsymbol{d}^{(k)}$. This is accomplished by choosing $\beta_k$ to be

$$\beta_k = \frac{\boldsymbol{g}^{(k+1)\top}\boldsymbol{Q}\boldsymbol{d}^{(k)}}{\boldsymbol{d}^{(k)\top}\boldsymbol{Q}\boldsymbol{d}^{(k)}}.$$

- The conjugate gradient algorithm is summarized below.

  1. Set $k := 0$; select the initial point $\boldsymbol{x}^{(0)}$.

  2. $\boldsymbol{g}^{(0)} = \nabla f\big(\boldsymbol{x}^{(0)}\big)$. If $\boldsymbol{g}^{(0)} = \boldsymbol{0}$, stop; else, set $\boldsymbol{d}^{(0)} = -\boldsymbol{g}^{(0)}$.

  3. $\alpha_k = -\dfrac{\boldsymbol{g}^{(k)\top}\boldsymbol{d}^{(k)}}{\boldsymbol{d}^{(k)\top}\boldsymbol{Q}\boldsymbol{d}^{(k)}}$.

  4. $\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} + \alpha_k \boldsymbol{d}^{(k)}$.

  5. $\boldsymbol{g}^{(k+1)} = \nabla f\big(\boldsymbol{x}^{(k+1)}\big)$. If $\boldsymbol{g}^{(k+1)} = \boldsymbol{0}$, stop.

  6. $\beta_k = \dfrac{\boldsymbol{g}^{(k+1)\top}\boldsymbol{Q}\boldsymbol{d}^{(k)}}{\boldsymbol{d}^{(k)\top}\boldsymbol{Q}\boldsymbol{d}^{(k)}}$.

  7. $\boldsymbol{d}^{(k+1)} = -\boldsymbol{g}^{(k+1)} + \beta_k \boldsymbol{d}^{(k)}$.

  8. Set $k := k + 1$; go to step 3.

- **Proposition 5.3.1** In the conjugate gradient algorithm, the directions $\boldsymbol{d}^{(0)}, \boldsymbol{d}^{(1)}, \dots, \boldsymbol{d}^{(n-1)}$ are $\boldsymbol{Q}$-conjugate.

**Example 5.3.1** Find the minimizer of

$$f(x_1, x_2, x_3) = \frac{3}{2}x_1^2 + 2x_2^2 + \frac{3}{2}x_3^2 + x_1 x_3 + 2x_2 x_3 - 3x_1 - x_3.$$

using the conjugate gradient method with the initial point $\boldsymbol{x}^{(0)} = [0, 0, 0]^\mathsf{T}$.

Note that the optimal is $\boldsymbol{x}^* = \boldsymbol{Q}^{-1}\boldsymbol{b} = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 4 & 2 \\ 1 & 2 & 3 \end{bmatrix}^{-1} \begin{bmatrix} 3 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.4 & -0.05 & -0.2 \\ 0.1 & 0.4 & -0.3 \\ -0.2 & -0.3 & 0.6 \end{bmatrix} \begin{bmatrix} 3 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}.$

We can represent $f$ as

$$f(\boldsymbol{x}) = \frac{1}{2}\boldsymbol{x}^\mathsf{T}\boldsymbol{Q}\boldsymbol{x} - \boldsymbol{x}^\mathsf{T}\boldsymbol{b},$$

where

$$\boldsymbol{Q} = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 4 & 2 \\ 1 & 2 & 3 \end{bmatrix}, \quad \boldsymbol{b} = \begin{bmatrix} 3 \\ 0 \\ 1 \end{bmatrix}.$$

We have

$$\boldsymbol{g}(\boldsymbol{x}) = \nabla f(\boldsymbol{x}) = \boldsymbol{Q}\boldsymbol{x} - \boldsymbol{b} = [3x_1 + x_3 - 3, 4x_2 + 2x_3, x_1 + 2x_2 + 3x_3 - 1]^\mathsf{T}.$$

Hence,

$$\boldsymbol{g}^{(0)} = [-3, 0, -1]^\mathsf{T},$$

$$\boldsymbol{d}^{(0)} = -\boldsymbol{g}^{(0)},$$

$$\alpha_0 = -\frac{\boldsymbol{g}^{(0)\mathsf{T}}\boldsymbol{d}^{(0)}}{\boldsymbol{d}^{(0)\mathsf{T}}\boldsymbol{Q}\boldsymbol{d}^{(0)}} = \frac{10}{36} = 0.2778$$

and

$$\boldsymbol{x}^{(1)} = \boldsymbol{x}^{(0)} + \alpha_0 \boldsymbol{d}^{(0)} = [0.8333, 0, 0.2778]^\mathsf{T}.$$

The next stage yields

$$\boldsymbol{g}^{(1)} = \nabla f(\boldsymbol{x}^{(1)}) = [-0.2222, 0.5556, 0.6667]^\mathsf{T},$$

$$\beta_0 = \frac{\boldsymbol{g}^{(1)\mathsf{T}}\boldsymbol{Q}\boldsymbol{d}^{(0)}}{\boldsymbol{d}^{(0)\mathsf{T}}\boldsymbol{Q}\boldsymbol{d}^{(0)}} = 0.08025.$$

We can now compute

$$\boldsymbol{d}^{(1)} = -\boldsymbol{g}^{(1)} + \beta_0 \boldsymbol{d}^{(0)} = [0.4630, -0.5556, -0.5864]^\mathsf{T}.$$

Hence,

$$\alpha_1 = -\frac{g^{(1)\top}d^{(1)}}{d^{(1)\top}Qd^{(1)}} = 0.2187$$

and

$$x^{(2)} = x^{(1)} + \alpha_1 d^{(1)} = [0.9364, -0.1215, 0.1495]^\top.$$

To perform the third iteration, we compute

$$g^{(2)} = \nabla f(x^{(2)}) = [-0.04673, -0.1869, 0.1402]^\top,$$

$$\beta_1 = \frac{g^{(2)\top}Qd^{(1)}}{d^{(1)\top}Qd^{(1)}} = 0.07075,$$

$$d^{(2)} = -g^{(2)} + \beta_1 d^{(1)} = [0.07948, 0.1476, -0.1817]^\top.$$

Hence,

$$\alpha_2 = -\frac{g^{(2)\top}d^{(2)}}{d^{(2)\top}Qd^{(2)}} = 0.8231,$$

and

$$x^{(3)} = x^{(2)} + \alpha_2 d^{(2)} = [1.000, 0.000, 0.000]^\top.$$

Note that

$$g^{(3)} = \nabla f(x^{(3)}) = 0,$$

as expected, because $f$ is a quadratic function of three variables. Hence, $x^* = x^{(3)}$.

## 5.4 THE CONJUGATE GRADIENT ALGORITHM FOR NONQUADRATIC PROBLEMS

- The algorithm can be extended to general nonlinear functions by interpreting $f(x) = (1/2)x^\top Qx - x^\top b$ as a second-order Taylor series approximation of the objective function.

- However, for a general nonlinear function the Hessian is a matrix that has to be reevaluated at each iteration of the algorithm. Thus, an efficient implementation of the conjugate gradient algorithm that eliminates the Hessian evaluation at each step is desirable.

- Observe that $Q$ appears only in the computation of the scalars $\alpha_k$ and $\beta_k$. Because
$$\alpha_k = \arg\min_{\alpha \geq 0} f(x^{(k)} + \alpha d^{(k)}),$$

the closed-form formula for $\alpha_k$ in the algorithm can be replaced by a numerical line search procedure.

- The modifications are all based on algebraically manipulating the formula $\beta_k$ in such a way that $\boldsymbol{Q}$ is eliminated.

- Note that $\boldsymbol{Q}\boldsymbol{d}^{(k)}$ and $(\boldsymbol{g}^{(k+1)} - \boldsymbol{g}^{(k)})/\alpha_k$ are equal in the quadratic case:

$$\boldsymbol{x}^{(k+1)} = \boldsymbol{x}^{(k)} + \alpha_k \boldsymbol{d}^{(k)}$$
$$\boldsymbol{Q}\boldsymbol{x}^{(k+1)} = \boldsymbol{Q}\boldsymbol{x}^{(k)} + \alpha_k \boldsymbol{Q}\boldsymbol{d}^{(k)}$$
$$\boldsymbol{Q}\boldsymbol{x}^{(k+1)} - \boldsymbol{b} = \boldsymbol{Q}\boldsymbol{x}^{(k)} - \boldsymbol{b} + \alpha_k \boldsymbol{Q}\boldsymbol{d}^{(k)}$$
$$\boldsymbol{g}^{(k+1)} = \boldsymbol{g}^{(k)} + \alpha_k \boldsymbol{Q}\boldsymbol{d}^{(k)}$$
$$\boldsymbol{Q}\boldsymbol{d}^{(k)} = (\boldsymbol{g}^{(k+1)} - \boldsymbol{g}^{(k)})/\alpha_k$$

- The **Hestenes-Stiefel formula** for $\beta_k$ is based on replacing the term $\boldsymbol{Q}\boldsymbol{d}^{(k)}$ by the term $(\boldsymbol{g}^{(k+1)} - \boldsymbol{g}^{(k)})/\alpha_k$:

$$\beta_k = \frac{\boldsymbol{g}^{(k+1)\top}\boldsymbol{Q}\boldsymbol{d}^{(k)}}{\boldsymbol{d}^{(k)\top}\boldsymbol{Q}\boldsymbol{d}^{(k)}} = \frac{\boldsymbol{g}^{(k+1)\top}(\boldsymbol{g}^{(k+1)} - \boldsymbol{g}^{(k)})/\alpha_k}{\boldsymbol{d}^{(k)\top}(\boldsymbol{g}^{(k+1)} - \boldsymbol{g}^{(k)})/\alpha_k} = \frac{\boldsymbol{g}^{(k+1)\top}(\boldsymbol{g}^{(k+1)} - \boldsymbol{g}^{(k)})}{\boldsymbol{d}^{(k)\top}(\boldsymbol{g}^{(k+1)} - \boldsymbol{g}^{(k)})}.$$

**Example 5.4.1** Perform one iteration of conjugate gradient descent with Hestenes-Stiefel formula modification to find the minimizer of $f(x_1, x_2, x_3) = (x_1 - 4)^4 + (x_2 - 3)^2 + 4(x_3 + 5)^4$. The initial point is $x^{(0)} = [4, 2, -1]^\top$.

$$g(x) = \nabla f(x) = [4(x_1 - 4)^3, 2(x_2 - 3), 16(x_3 + 5)^3]^\top.$$

**Iteration 1**

$$\boldsymbol{g}^{(0)} = \nabla f(\boldsymbol{x}^{(0)}) = [0, -2, 1024]^\top, \qquad \boldsymbol{d}^{(0)} = -\boldsymbol{g}^{(0)}, \qquad \alpha_0 = \operatorname*{argmin}_\alpha f(\boldsymbol{x}^{(0)} + \alpha_0 \boldsymbol{d}^{(0)})$$

Using the secant method, we obtain $\alpha_0 = 3.967 \times 10^{-3}$ and

$$\boldsymbol{x}^{(1)} = \boldsymbol{x}^{(0)} - \alpha_0 \boldsymbol{g}^{(0)} = [4, 2, -1]^\top - 3.967 \times 10^{-3}[0, -2, 1024]^\top$$
$$= [4.000, 2.007934, -5.061568]^\top \approx [4.000, 2.008, -5.062]^\top.$$

$$\boldsymbol{g}^{(1)} = \nabla f(\boldsymbol{x}^{(1)}) = [4(4 - 4)^3, 2(2.008 - 3), 16(-5.062 + 5)^3]^\top$$
$$\approx [0, -1.984, -0.003813]^\top$$

$$\beta_0 = \frac{\boldsymbol{g}^{(1)\top}(\boldsymbol{g}^{(1)} - \boldsymbol{g}^{(0)})}{\boldsymbol{d}^{(0)\top}(\boldsymbol{g}^{(1)} - \boldsymbol{g}^{(0)})}$$
$$= \frac{[0, -1.984, -0.003813]([0, -1.984, -0.003813]^\top - [0, -2, 1024]^\top)}{-[0, -2, 1024]([0, -1.984, -0.003813]^\top - [0, -2, 1024]^\top)}$$
$$\approx 0.000003693$$

$$\boldsymbol{d}^{(1)} = -\boldsymbol{g}^{(1)} + \beta_0 \boldsymbol{d}^{(0)} = -[0, -1.984, -0.003813]^\top + 0.000003693(-[0, -2, 1024]^\top)$$
$$\approx [0, 1.984, 0.00003137]^\top$$