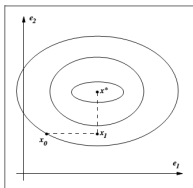


Chapter 10 Conjugate Direction Methods

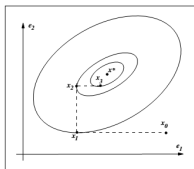
1. Conjugate Direction Method
2. Conjugate Gradient Method
3. Conjugate Gradient Method for Nonquadratic Problems



Conjugate Direction



when successively minimizing along the coordinate directions, we **can** reach the solution of a “**simple**” quadratic function within n iterations.



when successively minimizing along the coordinate directions, we **can not** reach the solution of a “**general**” quadratic function within n iterations.

★ “**general**” quadratic function can be modified by alerting the searching directions.



Conjugate Direction

unconstrained quadratic optimization

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{b}^\top \mathbf{x}, \text{ where } \mathbf{Q} \succ 0, \mathbf{b} \in \mathbb{R}^n.$$

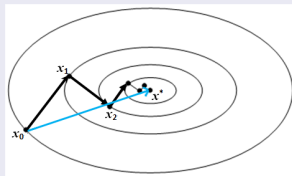
revisit of steepest descent and Newton's methods

- steepest descent method:

$$\begin{cases} \alpha_k = -\frac{\|\mathbf{g}^{(k)}\|^2}{\|\mathbf{g}^{(k)}\|_Q^2} \text{ is exact stepsize.} \\ \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)}, \end{cases}$$

- Newton's method:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}$$

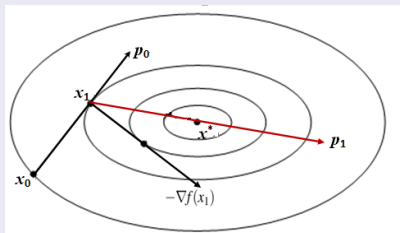


- steepest descent method: less computation and CPU memory at each iteration; but a large number of iterations to reach the minimizer \mathbf{x}^* .
- Newton's method: large computation and CPU memory at each iteration; but only fewer iteration can well approximate the minimizer \mathbf{x}^* .



Conjugate Direction

motivation of conjugate direction



$$\mathbf{x}^* = \mathbf{x}^{(1)} + t\mathbf{p}^{(1)}, \quad t \in \mathbb{R}$$

$$\mathbf{Q}\mathbf{x}^* - \mathbf{b} = \mathbf{Q}\mathbf{x}^{(1)} - \mathbf{b} + t\mathbf{Q}\mathbf{p}^{(1)}$$

$$\nabla f(\mathbf{x}^*) = \nabla f(\mathbf{x}^{(1)}) + t\mathbf{Q}\mathbf{p}_1^{(1)}$$

$$0 = \nabla f(\mathbf{x}_1) + t\mathbf{Q}\mathbf{p}_1^{(1)}$$

$$0 = \mathbf{p}^{(0)\top} \nabla f(\mathbf{x}_1) + t\mathbf{p}^{(0)\top} \mathbf{Q}\mathbf{p}^{(1)}$$

$$0 = \mathbf{p}^{(0)\top} \mathbf{Q}\mathbf{p}^{(1)}$$



Conjugate Direction

Notation:

\mathbb{S}^n : set of $n \times n$ symmetric matrices.

\mathbb{S}_+^n : set of $n \times n$ symmetric positive semidefinite matrices.

\mathbb{S}_{++}^n : set of $n \times n$ symmetric positive definite matrices.

Definition (conjugate direction)

Let $Q \in \mathbb{S}^n$. The vectors $p^{(0)}, p^{(1)}, \dots, p^{(m)}$ in \mathbb{R}^n are Q -conjugate if $p^{(i)\top} Q p^{(j)} = 0$ for all $i \neq j$.

Lemma

Let $Q \in \mathbb{S}_{++}^n$. If the vectors $p^{(0)}, p^{(1)}, \dots, p^{(k)}$ are nonzero and Q -conjugate, then they are linearly independent.

proof: analogously to the proof of orthogonal vectors.

Question: given a $Q \in \mathbb{S}_{++}^n$, how to generate Q -conjugate vectors?



Conjugate Direction: Three Methods

Method I

$$\text{Given } Q = \begin{bmatrix} 3 & 0 & 1 \\ 0 & 4 & 2 \\ 1 & 2 & 3 \end{bmatrix},$$

$$\text{let the } Q\text{-conjugate vectors be } \begin{cases} p^{(0)} = [1, 0, 0]^T \\ p^{(1)} = [p_1^{(1)}, p_2^{(1)}, p_3^{(1)}]^T \\ p^{(2)} = [p_1^{(2)}, p_2^{(2)}, p_3^{(2)}]^T \end{cases}$$

$$\xRightarrow{\text{definition}} \begin{cases} p^{(0)T} Q p^{(1)} = 0 \\ p^{(0)T} Q p^{(2)} = 0 \\ p^{(1)T} Q p^{(2)} = 0 \end{cases} \implies \dots \implies \begin{cases} p^{(0)} = [1, 0, 0]^T \\ p^{(1)} = [1, 0, -3]^T \\ p^{(2)} = [1, 4, -3]^T \end{cases}$$



Conjugate Direction: Three Methods

Method II: Gram-Schmidt process

Given $Q \in \mathbb{S}_{++}^n$ and linear independent vectors $\{\mathbf{a}^{(i)}\}_{i=0}^{n-1}$ in \mathbb{R}^n , Gram-Schmidt process can transform $\{\mathbf{a}^{(i)}\}_{i=0}^{n-1}$ into Q -conjugate vectors

Gram-Schmidt process:
$$\begin{cases} \mathbf{p}^{(0)} = \mathbf{a}^{(0)}. \\ \mathbf{p}^{(k)} = \mathbf{a}^{(k)} - \sum_{i=0}^{k-1} \frac{\mathbf{a}^{(k)\top} Q \mathbf{p}^{(i)}}{\mathbf{p}^{(i)\top} Q \mathbf{p}^{(i)}} \mathbf{p}^{(i)}, \quad k = 1, \dots, n-1. \end{cases}$$

The vectors $\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(n-1)}$ are Q -conjugate.

- ★ Particularly, if $Q = I$, Q -conjugate reduces to the orthogonal. Thus, the above process reduces to the Gram-Schmidt orthogonalization.



Conjugate Direction: Three Methods

Method III: eigenvalue decomposition (see exercise 10.4)

Let $Q \in \mathbb{S}^n$. Then, there exists Q -conjugate vectors $\{p^{(1)}, \dots, p^{(n)}\}$ such that each $p^{(i)}$ ($i = 1, \dots, n$) is an eigenvector of Q .

proof. Use the fact that for any $Q \in \mathbb{S}^n$, there exists vectors $\{v^{(1)}, \dots, v^{(n)}\}$ of its eigenvectors such that $v_i^\top v_j = 0$ for all $i, j = 1, \dots, n, i \neq j$.

Theorem

Let $Q \in \mathbb{S}_{++}^n$. If nonzero vectors $\{p^{(1)}, \dots, p^{(n)}\}$ are Q -conjugate and also orthogonal, then each $p^{(i)}, i = 1, \dots, n$, is an eigenvector of Q .

- ★ orthogonal eigenvectors of a positive definite matrix are the Q -conjugate. Conversely, if $\{p^{(1)}, \dots, p^{(n)}\}$ are orthogonal and Q -conjugate, then they are eigenvectors of Q .



Basic Conjugate Direction Algorithm

basic conjugate direction algorithm

To minimizing a quadratic function $\min f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{b}^\top \mathbf{x}$, where $\mathbf{Q} \in \mathbb{S}_{++}^n$ and $\mathbf{b} \in \mathbb{R}^n$, given the initial point $\mathbf{x}^{(0)}$ and \mathbf{Q} -conjugate directions $\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots, \mathbf{p}^{(n-1)}$, the iteration scheme reads

$$\textcircled{1} \quad \mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) = \mathbf{Q}\mathbf{x}^{(k)} - \mathbf{b},$$

$$\textcircled{2} \quad \alpha_k = -\frac{\mathbf{g}^{(k)\top} \mathbf{p}^{(k)}}{\mathbf{p}^{(k)\top} \mathbf{Q} \mathbf{p}^{(k)}},$$

$$\textcircled{3} \quad \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}.$$



Basic Conjugate Direction Algorithm

Example (minimize f by basic conjugate direction algorithms)

$\min f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \mathbf{x} - \mathbf{x}^\top \begin{bmatrix} -1 \\ 1 \end{bmatrix}$. Given the \mathbf{Q} -conjugate directions $\mathbf{p}^{(0)} = [1, 0]^\top$, $\mathbf{p}^{(1)} = [-\frac{3}{8}, \frac{3}{4}]^\top$ and the initial point $\mathbf{x}^{(0)} = [0, 0]^\top$.

Ans: $\because \nabla f(\mathbf{x}) = \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \mathbf{x} - \begin{bmatrix} -1 \\ 1 \end{bmatrix}$, $\nabla^2 f(\mathbf{x}) = \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} =: \mathbf{Q}$

S1. $\mathbf{g}^{(0)} = -\mathbf{b} = [1, -1]^\top$, $\alpha_0 = -\frac{\mathbf{g}^{(0)\top} \mathbf{p}^{(0)}}{\mathbf{p}^{(0)\top} \mathbf{Q} \mathbf{p}^{(0)}} = -\frac{1}{4}$,
 $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{p}^{(0)} = [-\frac{1}{4}, 0]^\top$.

S2. $\mathbf{g}^{(1)} = \mathbf{Q} \mathbf{x}^{(1)} - \mathbf{b} = [0, -\frac{3}{2}]^\top$, $\alpha_1 = -\frac{\mathbf{g}^{(1)\top} \mathbf{p}^{(1)}}{\mathbf{p}^{(1)\top} \mathbf{Q} \mathbf{p}^{(1)}} = 2$,
 $\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{p}^{(1)} = [-1, \frac{3}{2}]^\top$.

$\because f$ is a 2D quadratic function,

$\therefore \mathbf{x}^* = \mathbf{x}^{(2)}$ is a minimizer.



Basic Conjugate Direction Algorithm

Theorem

For any initial point $\mathbf{x}^{(0)}$, the basic conjugate direction algorithm converges to the unique \mathbf{x}^* (indeed, $\mathbf{x}^* = \mathbf{Q}^{-1}\mathbf{b}$) in n steps, i.e., $\mathbf{x}^{(n)} = \mathbf{x}^*$.

proof. \because The vectors $\{\mathbf{p}^{(i)}\}_{i=0}^{n-1}$ are \mathbf{Q} -conjugate,
 $\therefore \{\mathbf{p}^{(i)}\}_{i=0}^{n-1}$ are linear independent, $\therefore \{\mathbf{p}^{(i)}\}_{i=0}^{n-1}$ is a basis of \mathbb{R}^n .
 \therefore The vector $\mathbf{x}^* - \mathbf{x}^{(0)} \in \mathbb{R}^n$ can be linearly expressed as

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

Premultiply $\mathbf{p}^{(k)\top} \mathbf{Q}$ to both sides of the above equation

$$\mathbf{p}^{(k)\top} \mathbf{Q}(\mathbf{x}^* - \mathbf{x}^{(0)}) = \beta_k \mathbf{p}^{(k)\top} \mathbf{Q} \mathbf{p}^{(k)} \implies \beta_k = \frac{\mathbf{p}^{(k)\top} \mathbf{Q}(\mathbf{x}^* - \mathbf{x}^{(0)})}{\mathbf{p}^{(k)\top} \mathbf{Q} \mathbf{p}^{(k)}}.$$

On the other hand, the iterate $\mathbf{x}^{(k)}$ can be written as

$$\mathbf{x}^{(k)} - \mathbf{x}^{(0)} = \alpha_0 \mathbf{p}^{(0)} + \cdots + \alpha_{k-1} \mathbf{p}^{(k-1)}. \quad (2)$$

$$\therefore \mathbf{x}^{(n)} - \mathbf{x}^{(0)} = \alpha_0 \mathbf{p}^{(0)} + \cdots + \alpha_{n-1} \mathbf{p}^{(n-1)} \text{ with } \alpha_k = -\frac{\mathbf{g}^{(k)\top} \mathbf{p}^{(k)}}{\mathbf{p}^{(k)\top} \mathbf{Q} \mathbf{p}^{(k)}}. \quad (3)$$

By comparing (1) and (3), $\mathbf{x}^{(n)} = \mathbf{x}^{(0)} \iff \boxed{\beta_k = \alpha_k}$ for $k = 0, \dots, n-1$.



Basic Conjugate Direction Algorithm

Indeed, the numerators of β_k and α_k are equal.

$$\begin{aligned} & \therefore \mathbf{p}^{(k)\top} \mathbf{Q}(\mathbf{x}^* - \mathbf{x}^{(0)}) \\ &= \mathbf{p}^{(k)\top} \mathbf{Q}(\mathbf{x}^* - \mathbf{x}^{(k)} + \mathbf{x}^{(k)} - \mathbf{x}^{(0)}) \\ &= \mathbf{p}^{(k)\top} \mathbf{Q}(\mathbf{x}^* - \mathbf{x}^{(k)}) + \underbrace{\mathbf{p}^{(k)\top} \mathbf{Q}(\mathbf{x}^{(k)} - \mathbf{x}^{(0)})}_{=0, \text{ (why?)}} \\ &= \mathbf{p}^{(k)\top} (\mathbf{b} - \mathbf{Q}\mathbf{x}^{(k)}) = -\mathbf{p}^{(k)\top} \mathbf{g}^{(k)}. \\ & \therefore \beta_k = -\frac{\mathbf{p}^{(k)\top} \mathbf{g}^{(k)}}{\mathbf{p}^{(k)\top} \mathbf{Q}\mathbf{p}^{(k)}} = \alpha_k, \text{ which implies that } \mathbf{x}^* = \mathbf{x}^{(n)}. \end{aligned}$$

Lemma (accurate stepsize: $\alpha_k = \arg \min f(\mathbf{x}^{(k)} + \alpha \mathbf{p}^{(k)})$)

In exact line search, the searching direction orthogonal to gradient, i.e., $\mathbf{g}^{(k+1)\top} \mathbf{p}^{(k)} = 0$. (see Chapter 8 for proof)



Basic Conjugate Direction Algorithm

Lemma

In conjugate direction method, given a k satisfying $0 \leq k \leq n-1$, then $\mathbf{g}^{(k+1)\top} \mathbf{p}^{(i)} = 0$ for all $0 \leq i \leq k$.

proof. $\because \mathbf{g}^{(k)} = \mathbf{Q}\mathbf{x}^{(k)} - \mathbf{b}$.

$$\therefore \mathbf{g}^{(k+1)} - \mathbf{g}^{(k)} = \mathbf{Q}(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) \xrightarrow{\text{recursion}} \alpha_k \mathbf{Q}\mathbf{p}^{(k)}.$$

(By induction) It follows from the above lemma that $\mathbf{g}^{(1)\top} \mathbf{p}^{(0)} = 0$.

(Induction hypothesis) If the assertion holds for $k-1$, i.e.,

$$\mathbf{g}^{(k)\top} \mathbf{p}^{(i)} = 0 \text{ for all } 0 \leq i \leq k-1.$$

We now prove that assertion holds for k , i.e., $\mathbf{g}^{(k+1)\top} \mathbf{p}^{(i)} = 0$ for all $0 \leq i \leq k$.

- It follows from the above lemma on accurate step, $\mathbf{g}^{(k+1)\top} \mathbf{p}^{(k)} = 0$.
- Furthermore, it suffice to prove $\mathbf{g}^{(k+1)\top} \mathbf{p}^{(i)} = 0$ for all $0 \leq i \leq k-1$.

$$\because \mathbf{g}^{(k+1)} = \mathbf{g}^{(k)} + \alpha_k \mathbf{Q}\mathbf{p}^{(k)}.$$

$$\therefore \mathbf{g}^{(k+1)\top} \mathbf{p}^{(i)} = \mathbf{g}^{(k)\top} \mathbf{p}^{(i)} + \alpha_k \mathbf{p}^{(k)\top} \mathbf{Q}\mathbf{p}^{(i)} = 0 \text{ for all } 0 \leq i \leq k-1$$

(why?)

Thus, given a $0 \leq k \leq n-1$ $\mathbf{g}^{(k+1)\top} \mathbf{p}^{(i)} = 0$ for all $0 \leq i \leq k$.

$$\star \mathbf{g}^{(k+1)} \perp \text{span}[\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots, \mathbf{p}^{(k)}].$$



Expanding Subspace Theorem

Theorem

The conjugate direction method satisfies $f(\mathbf{x}^{(k+1)}) = \min_{\alpha} f(\mathbf{x}^{(k)} + \alpha \mathbf{p}^{(k)})$, and it can also satisfy $f(\mathbf{x}^{(k+1)}) = \min_{a_0, \dots, a_k} f(\mathbf{x}^{(0)} + \sum_{i=0}^k a_i \mathbf{p}^{(i)})$. In another word, let $\mathcal{V}_k = \mathbf{x}^{(0)} + \text{span}[\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots, \mathbf{p}^{(k)}]$. Then $f(\mathbf{x}^{(k+1)}) = \min_{\mathbf{x} \in \mathcal{V}_k} f(\mathbf{x})$.

proof. let matrix $\mathbf{D}^{(k)} = [\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots, \mathbf{p}^{(k)}]$.

$$\therefore \mathcal{V}_k = \mathbf{x}^{(0)} + \text{span}[\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots, \mathbf{p}^{(k)}] = \mathbf{x}^{(0)} + \mathcal{R}(\mathbf{D}^{(k)}).$$

$$\therefore \mathbf{x}^{(k+1)} = \mathbf{x}^{(0)} + \sum_{i=0}^k \alpha_i \mathbf{p}^{(i)} = \mathbf{x}^{(0)} + \mathbf{D}^{(k)} \boldsymbol{\alpha}, \text{ where } \boldsymbol{\alpha} = [\alpha_0, \dots, \alpha_k]^T.$$

$$\therefore \mathbf{x}^{(k+1)} \in \mathbf{x}^{(0)} + \mathcal{R}(\mathbf{D}^{(k)}) = \mathcal{V}_k.$$

For any $\mathbf{x} \in \mathcal{V}_k$, there exists an $\boldsymbol{\alpha} \in \mathbb{R}^{k+1}$ such that $\mathbf{x} = \mathbf{x}^{(0)} + \mathbf{D}^{(k)} \boldsymbol{\alpha}$.



Expanding Subspace Theorem

Let $\phi_k(\alpha) := f(\mathbf{x}^{(0)} + \mathbf{D}^{(k)}\alpha)$.

Then, ϕ_k is a quadratic function with unique minimum.

By the chain rules,

$$D\phi_k(\alpha) = \nabla f(\mathbf{x}^{(0)} + \mathbf{D}^{(k)}\alpha)^\top \mathbf{D}^{(k)} = \nabla f(\mathbf{x}^{(k+1)})^\top \mathbf{D}^{(k)} = \mathbf{g}^{(k+1)\top} \mathbf{D}^{(k)}. \quad (4)$$

It follows from the above Lemma 10.2 that $\mathbf{g}^{(k+1)\top} \mathbf{D}^{(k)} = \mathbf{0}^\top$.

$\therefore (4) \implies D\phi_k(\alpha) = 0 \implies \alpha$ is the minimum of ϕ_k , i.e.,

$$f(\mathbf{x}^{(k+1)}) = \min_{\alpha} f(\mathbf{x}^{(0)} + \mathbf{D}^{(k)}\alpha) = \min_{\mathbf{x} \in \mathcal{V}_k} f(\mathbf{x}).$$

★ As k increases, the subspace $\text{span}[\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots, \mathbf{p}^{(k)}]$ continues to expand until it fills the entire \mathbb{R}^n (provided that the vectors $\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots$ are linear independent).

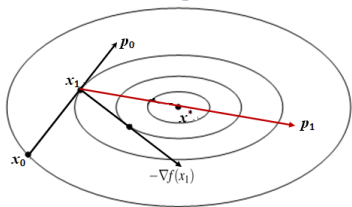


Conjugate Gradient Method

motivation of conjugate gradient (CG) method

- CG method does not use prespecified conjugate directions, but computes the directions as method progresses. At k th iteration, $\mathbf{p}^{(k)}$ is calculated as a linear combination of $\mathbf{p}^{(k-1)}$ and the current gradient $\mathbf{g}^{(k)}$, in such a way that all the directions are \mathbf{Q} -conjugate.
- CG method exploits the fact that for a quadratic function of n variables, we can locate the minimizer by performing n searches along conjugate directions.

Analysis: given initial point $\mathbf{x}^{(0)}$ and exact stepsize criterion.



- **step 0:** use gradient $\mathbf{p}^{(0)} = -\mathbf{g}^{(0)}$ as the search direction to produce the new iterate

$$\alpha_0 = -\frac{\mathbf{g}^{(0)\top} \mathbf{p}^{(0)}}{\mathbf{p}^{(0)\top} \mathbf{Q} \mathbf{p}^{(0)}}$$

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{p}^{(0)}.$$



Conjugate Gradient Method

Analysis: given initial point $\mathbf{x}^{(0)}$ and exact stepsize criterion.

- **step 1:** let $\mathbf{p}^{(1)}$ be the linear combination of $\mathbf{g}^{(1)}$ and $\mathbf{p}^{(0)}$, i.e.,

$$\mathbf{p}^{(1)} = -\mathbf{g}^{(1)} + \beta_0 \mathbf{p}^{(0)}$$

$$\underline{\underline{Q\text{-conjugacy of } \mathbf{p}^{(1)}, \mathbf{p}^{(0)}}} \Rightarrow \beta_0 = \frac{\mathbf{g}^{(1)\top} \mathbf{Q} \mathbf{d}^{(0)}}{\mathbf{p}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}}.$$

use this $\mathbf{p}^{(1)}$ as the search direction to generate the new iterate

$$\alpha_1 = -\frac{\mathbf{g}^{(1)\top} \mathbf{p}^{(1)}}{\mathbf{p}^{(1)\top} \mathbf{Q} \mathbf{p}^{(1)}}$$

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{p}^{(1)}.$$

- **step k :** compute the search direction

$$\mathbf{p}^{(k)} = -\mathbf{g}^{(k)} + \beta_{k-1} \mathbf{p}^{(k-1)}, \text{ where } \beta_{k-1} = \frac{\mathbf{g}^{(k)\top} \mathbf{Q} \mathbf{d}^{(k-1)}}{\mathbf{p}^{(k-1)\top} \mathbf{Q} \mathbf{d}^{(k-1)}}. \quad (5)$$



Conjugate Gradient Method

Proposition

The vectors $\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots, \mathbf{p}^{(n-1)}$ produced by (5) are Q -conjugate.

proof. (induction)

(i) First, prove $\mathbf{p}^{(0)\top} Q \mathbf{p}^{(1)} = 0$.

$$\therefore \mathbf{p}^{(0)\top} Q \mathbf{p}^{(1)} = \mathbf{p}^{(0)\top} Q (-\mathbf{g}^{(1)} + \beta_0 \mathbf{p}^{(0)}) \stackrel{\beta_0 = \frac{\mathbf{g}^{(1)\top} Q \mathbf{p}^{(0)}}{\mathbf{p}^{(0)\top} Q \mathbf{p}^{(0)}}}{=} 0.$$

(ii) Inductive hypothesis: $\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots, \mathbf{p}^{(k)}$ are Q -conjugate.

$$\stackrel{\text{By Lemma 10.2}}{\implies} \mathbf{g}^{(k+1)} \perp \text{span}[\mathbf{p}^{(0)}, \mathbf{p}^{(1)}, \dots, \mathbf{p}^{(k)}].$$

We first prove $\mathbf{g}^{(k+1)\top} \mathbf{g}^{(j)} = 0$ for all $0 \leq j \leq k$.

For a given $0 \leq j \leq k$, we have $\mathbf{p}^{(j)} = -\mathbf{g}^{(j)} + \beta_{j-1} \mathbf{p}^{(j-1)}$.

$$\therefore 0 = \mathbf{g}^{(k+1)\top} \mathbf{p}^{(j)} = -\mathbf{g}^{(k+1)\top} \mathbf{g}^{(j)} + \beta_{j-1} \underbrace{\mathbf{g}^{(k+1)\top} \mathbf{p}^{(j-1)}}_{0 \text{ (why?)}}$$

$$\therefore \mathbf{g}^{(k+1)\top} \mathbf{g}^{(j)} = 0 \text{ for all } 0 \leq j \leq k.$$



Conjugate Gradient Method

We now prove $\mathbf{p}^{(k+1)\top} \mathbf{Qp}^{(j)} = 0, j = 0, \dots, k$

$$\begin{aligned}\therefore \mathbf{p}^{(k+1)\top} \mathbf{Qp}^{(j)} &= (-\mathbf{g}^{(k+1)} + \beta_k \mathbf{p}^{(k)})^\top \mathbf{Qp}^{(j)} \\ &= -\mathbf{g}^{(k+1)\top} \mathbf{Qp}^{(j)} + \beta_k \mathbf{p}^{(k)\top} \mathbf{Qp}^{(j)}\end{aligned}$$

Case 1: when $j < k$, inductive hypothesis $\implies \mathbf{p}^{(k)\top} \mathbf{Qp}^{(j)} = 0$.

$$\therefore \mathbf{p}^{(k+1)\top} \mathbf{Qp}^{(j)} = -\mathbf{g}^{(k+1)\top} \mathbf{Qp}^{(j)}.$$

Since we know that

$$\mathbf{g}^{(j+1)} = \mathbf{g}^{(j)} + \alpha_j \mathbf{Qp}^{(j)} \text{ and } \mathbf{g}^{(k+1)\top} \mathbf{g}^{(i)} = 0, i = 0, \dots, k$$

$$\therefore \mathbf{p}^{(k+1)\top} \mathbf{Qp}^{(j)} = -\mathbf{g}^{(k+1)\top} \frac{(\mathbf{g}^{(j+1)} - \mathbf{g}^{(j)})}{\alpha_j} = 0.$$

$$\therefore \mathbf{p}^{(k+1)\top} \mathbf{Qp}^{(j)} = 0, j = 0, \dots, k-1.$$

Case 2: when $j = k$, $\mathbf{p}^{(k+1)\top} \mathbf{Qp}^{(j)} \stackrel{\because j=k}{=} \mathbf{p}^{(k+1)\top} \mathbf{Qp}^{(k)} =$
 $(-\mathbf{g}^{(k+1)} + \beta_k \mathbf{p}^{(k)})^\top \mathbf{Qp}^{(k)} \stackrel{\because \beta_k}{=} 0.$

Overall, it follows from cases 1 and 2 that $\mathbf{p}^{(k+1)\top} \mathbf{Qp}^{(j)} = 0, j = 0, \dots, k$



Pseudo-Code of CG Method

Pseudo-code of CG method

Require: initial point $\mathbf{x}^{(0)}$, tolerance ϵ ,
set $k = 0$ and $\mathbf{p}^{(0)} = -\nabla f(\mathbf{x}^{(0)})$

1: **repeat**

$$2: \alpha_k = -\frac{\mathbf{g}^{(k)\top} \mathbf{p}^{(k)}}{\mathbf{p}^{(k)\top} \mathbf{Q} \mathbf{p}^{(k)}};$$

$$3: \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)};$$

$$4: \mathbf{g}^{(k+1)} = \nabla f(\mathbf{x}^{(k+1)});$$

$$5: \beta_k = \frac{\mathbf{g}^{(k+1)\top} \mathbf{Q} \mathbf{p}^{(k)}}{\mathbf{p}^{(k)\top} \mathbf{Q} \mathbf{p}^{(k)}};$$

$$6: \mathbf{p}^{(k+1)} = -\mathbf{g}^{(k+1)} + \beta_k \mathbf{p}^{(k)};$$

7: **until** $\|\mathbf{g}^{(k)}\| < \epsilon$.

other choices of β_k

$$\begin{aligned} \beta_k &= \frac{\mathbf{g}^{(k+1)\top} \mathbf{Q} \mathbf{p}^{(k)}}{\mathbf{p}^{(k)\top} \mathbf{Q} \mathbf{p}^{(k)}} \\ &= \frac{\mathbf{g}^{(k+1)\top} [\mathbf{g}^{(k+1)} - \mathbf{g}^{(k)}]}{\mathbf{p}^{(k)\top} [\mathbf{g}^{(k+1)} - \mathbf{g}^{(k)}]} \quad (\text{Hestenes-Stiefel}) \\ &= \frac{\mathbf{g}^{(k+1)\top} [\mathbf{g}^{(k+1)} - \mathbf{g}^{(k)}]}{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}} \quad (\text{Polak-Ribiere}) \\ &= \frac{\mathbf{g}^{(k+1)\top} \mathbf{g}^{(k+1)}}{\mathbf{g}^{(k)\top} \mathbf{g}^{(k)}} \quad (\text{Fletcher-Reeves}) \\ &= \frac{\mathbf{g}^{(k+1)\top} \mathbf{g}^{(k+1)}}{\mathbf{p}^{(k)\top} [\mathbf{g}^{(k+1)} - \mathbf{g}^{(k+1)}]} \quad (\text{Dai-Yuan}). \end{aligned}$$

★ For quadratic problems, all the above formulae of β_k are equivalent. But, for nonquadratic problems, they yield different solution.

★ Powell proves that the **Fletcher-Peeves** formula is more stable.



Conjugate Gradient Method

Example (minimize f by CG method)

$f(x_1, x_2) = x_1^2 + 2x_2^2$ with initial point $\mathbf{x}^{(0)} = [5, 5]^\top$.

- **Step 1:** $\mathbf{g}^{(0)} = \nabla f(\mathbf{x}^{(0)}) = [10, 20]^\top$, $\mathbf{p}^{(0)} = -\mathbf{g}^{(0)} = [-10, -20]^\top$.

$$\alpha_0 = -\frac{\mathbf{g}^{(0)\top} \mathbf{p}^{(0)}}{\mathbf{p}^{(0)\top} \mathbf{Q} \mathbf{p}^{(0)}} = \frac{5}{18},$$

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{p}^{(0)} = \left[\frac{20}{9}, \frac{-5}{9}\right]^\top.$$

- **Step 2:** $\mathbf{g}^{(1)} = \nabla f(\mathbf{x}^{(1)}) = \left[\frac{-40}{9}, \frac{20}{9}\right]^\top$.

$$\beta_0 = \frac{\mathbf{g}^{(1)\top} \mathbf{Q} \mathbf{d}^{(0)}}{\mathbf{p}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} = \frac{4}{81}.$$

$$\mathbf{p}^{(1)} = -\mathbf{g}^{(1)} + \beta_0 \mathbf{p}^{(0)} = \left[\frac{400}{81}, \frac{100}{81}\right]^\top.$$

$$\alpha_1 = -\frac{\mathbf{g}^{(1)\top} \mathbf{p}^{(1)}}{\mathbf{p}^{(1)\top} \mathbf{Q} \mathbf{p}^{(1)}} = \frac{9}{20},$$

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{p}^{(1)} = [0, 0]^\top.$$

$\therefore \mathbf{g}^{(3)} = \nabla f(\mathbf{x}^{(3)}) = \mathbf{0}$, $\therefore \mathbf{x}^* = \mathbf{x}^{(2)}$ is the minimizer.



Conjugate Gradient Method

Example (minimize f by CG method)

$f(x_1, x_2, x_3) = \frac{3}{2}x_1^2 + 2x_2^2 + \frac{3}{2}x_3^2 + x_1x_3 + 2x_2x_3 - 3x_1 - x_3$ with initial point $\mathbf{x}^{(0)} = [0, 0, 0]^\top$.

- **Step 1:** $\mathbf{g}^{(0)} = \nabla f(\mathbf{x}^{(0)}) = [-3, 0, -1]^\top$, $\mathbf{p}^{(0)} = -\mathbf{g}^{(0)}$.

$$\alpha_0 = -\frac{\mathbf{g}^{(0)\top} \mathbf{p}^{(0)}}{\mathbf{p}^{(0)\top} \mathbf{Q} \mathbf{p}^{(0)}} = \frac{10}{36}, \quad \mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{p}^{(0)} = [0.8333, 0, 0.2778]^\top.$$

- **Step 2:** $\mathbf{g}^{(1)} = \nabla f(\mathbf{x}^{(1)}) = [-0.2222, 0.5556, 0.6667]^\top$.

$$\beta_0 = \frac{\mathbf{g}^{(1)\top} \mathbf{Q} \mathbf{d}^{(0)}}{\mathbf{p}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} = 0.08025.$$

$$\mathbf{p}^{(1)} = -\mathbf{g}^{(1)} + \beta_0 \mathbf{p}^{(0)} = [0.4630, -0.5556, -0.5864]^\top.$$

$$\alpha_1 = -\frac{\mathbf{g}^{(1)\top} \mathbf{p}^{(1)}}{\mathbf{p}^{(1)\top} \mathbf{Q} \mathbf{p}^{(1)}} = 0.2187,$$

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{p}^{(1)} = [0.9346, -0.1215, 0.1495]^\top.$$

- **Step 3:** $\mathbf{g}^{(2)} = \nabla f(\mathbf{x}^{(2)}) = [-0.04673, -0.1869, 0.1402]^\top$.

$$\beta_1 = \frac{\mathbf{g}^{(2)\top} \mathbf{Q} \mathbf{d}^{(1)}}{\mathbf{p}^{(1)\top} \mathbf{Q} \mathbf{d}^{(1)}} = 0.07075.$$

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

$$\alpha_2 = -\frac{\mathbf{g}^{(2)\top} \mathbf{d}^{(2)}}{\mathbf{d}^{(2)\top} \mathbf{Q} \mathbf{d}^{(2)}} = 0.8231, \quad \mathbf{x}^{(3)} = \mathbf{x}^{(2)} + \alpha_2 \mathbf{p}^{(2)} = [1.0, 0.0, 0.0]^\top.$$

$$\therefore \mathbf{g}^{(3)} = \nabla f(\mathbf{x}^{(3)}) = \mathbf{0}, \quad \therefore \mathbf{x}^* = \mathbf{x}^{(3)} \text{ is the minimizer.}$$



Extension of CG to Nonquadratic Problems

Extension: Theoretically, Q in CG method is replaced by **Hessian** of objective function. However, for general nonlinear functions, each iteration must refresh the Hessian, which requires **a large amount of computation**.

Pseudo-code of CG method for nonquadratic problems

1. Set $k = 0$, select the initial value $\mathbf{x}^{(0)}$.
2. $\mathbf{g}^{(0)} = \nabla f(\mathbf{x}^{(0)})$, if $\mathbf{g}^{(0)} = \mathbf{0}$, stop; else, set $\mathbf{p}^{(0)} = -\mathbf{g}^{(0)}$.
3. $\alpha_k = -\frac{\mathbf{g}^{(k)\top} \mathbf{p}^{(k)}}{\mathbf{p}^{(k)\top} \mathbf{Q} \mathbf{p}^{(k)}} \longrightarrow \boxed{\alpha_k = \arg \min_{a \geq 0} f(\mathbf{x}^{(k)} + a \mathbf{p}^{(k)})}$.
4. $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}$.
5. $\mathbf{g}^{(k+1)} = \nabla f(\mathbf{x}^{(k+1)})$. If $\mathbf{g}^{(k+1)} = \mathbf{0}$, stop.
6. $\beta_k = \frac{\mathbf{g}^{(k+1)\top} \mathbf{Q} \mathbf{p}^{(k)}}{\mathbf{p}^{(k)\top} \mathbf{Q} \mathbf{p}^{(k)}} \longrightarrow \boxed{\text{HS/PR/FR formula}}$.
7. $\mathbf{p}^{(k+1)} = -\mathbf{g}^{(k+1)} + \beta_k \mathbf{p}^{(k)}$.
8. Set $k = k + 1$, go to step 3.

- ★ The search directional conjugacy and orthogonality disappear in the nonquadratic conjugate gradient method. **[reason?]**
- ★ No longer has the n-step convergence. Generally need a **restrat** operation.



Extension of CG to Nonquadratic Problems

- ① CG method is related to the **Krylov subspace method**. Krylov subspace method (proposed by Hestenes, Stiefel and Lanczos), which was rated as one of the ten **most influential methods** in science and engineering of the 20th century.
- ② Characteristics of CG method.
 - ① For the n D quadratic problem, the results can **be obtained in n steps**.
 - ② CG method **not need to calculate the Hessian matrix**.
 - ③ There is **no need to store $n \times n$ matrices** and **not perform inverse operations**.
- ③ In terms of calculation efficiency, **CG and gradient descent method** is between the fastest descent method and the Newton method. CG method is suitable to solve with large-scale problems.



Exercise

Exercise in the textbook: 10.4, 10.10.

Note: 10.4(b) should modify:

Q is a positive and definite matrix, $\{\mathbf{p}^{(1)}, \dots, \mathbf{p}^{(n)}\}$ is an orthogonal vectors, if $\{\mathbf{p}^{(1)}, \dots, \mathbf{p}^{(n)}\}$ is Q -conjugated, so $\{\mathbf{p}^{(1)}, \dots, \mathbf{p}^{(n)}\}$ must be the eigenvector of matrix Q .

