

Lemma 1.7.3 For any $k, \ell \geq 0$, $k \neq \ell$, we have $\langle f'(\mathbf{x}_k), f'(\mathbf{x}_\ell) \rangle = 0$.

Proof: Let $k \geq i$, and consider

$$\phi(\boldsymbol{\lambda}) = f\left(\mathbf{x}_0 + \sum_{j=1}^k \lambda_j f'(\mathbf{x}_{j-1})\right).$$

From the previous lemma, there is a $\boldsymbol{\lambda}^*$ such that $\mathbf{x}_k = \mathbf{x}_0 + \sum_{j=1}^k \lambda_j^* f'(\mathbf{x}_{j-1})$. Moreover, $\boldsymbol{\lambda}^*$ is the minimum of the function $\phi(\boldsymbol{\lambda})$. Therefore,

$$\frac{\partial \phi}{\partial \lambda_i}(\boldsymbol{\lambda}^*) = \langle f'(\mathbf{x}_k), f'(\mathbf{x}_{i-1}) \rangle = 0.$$

■

Corollary 1.7.4 The sequence generated by the conjugate gradient method for the convex quadratic function is finite.

Proof: Since the number of orthogonal directions in \mathbb{R}^n cannot exceed n .

■

Let us define $\boldsymbol{\delta}_i = \mathbf{x}_{i+1} - \mathbf{x}_i$. It is clear that $\mathcal{L}_k = \text{Lin}\{\boldsymbol{\delta}_0, \boldsymbol{\delta}_1, \dots, \boldsymbol{\delta}_{k-1}\}$.

Lemma 1.7.5 For any $k, \ell \geq 0$, $k \neq \ell$, $\langle \mathbf{A}\boldsymbol{\delta}_k, \boldsymbol{\delta}_\ell \rangle = 0$.

Proof: Let $k > \ell$. Then

$$\langle \mathbf{A}\boldsymbol{\delta}_k, \boldsymbol{\delta}_\ell \rangle = \langle \mathbf{A}(\mathbf{x}_{k+1} - \mathbf{x}_k), \boldsymbol{\delta}_\ell \rangle = \langle f'(\mathbf{x}_{k+1}) - f'(\mathbf{x}_k), \mathbf{x}_{\ell+1} - \mathbf{x}_\ell \rangle = 0,$$

due to Lemma 1.7.3.

■

The vectors $\{\boldsymbol{\delta}_i\}$ are called *conjugate* with respect to matrix \mathbf{A} .

Now, let us be more precise with the conjugate gradient method. We will define the next iterations as follows:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - h_k f'(\mathbf{x}_k) + \sum_{j=0}^{k-1} \lambda_j \boldsymbol{\delta}_j$$

Using the previous properties, we arrive that

$$\lambda_j = 0, \quad (j = 0, 1, \dots, k-2), \quad \lambda_{k-1} = \frac{h_k \|f'(\mathbf{x}_k)\|^2}{\langle f'(\mathbf{x}_k) - f'(\mathbf{x}_{k-1}), \boldsymbol{\delta}_{k-1} \rangle}. \quad (1.6)$$

Thus

$$\mathbf{x}_{k+1} = \mathbf{x}_k - h_k \mathbf{p}_k$$

where

$$\mathbf{p}_k = f'(\mathbf{x}_k) - \frac{\|f'(\mathbf{x}_k)\|^2 \mathbf{p}_{k-1}}{\langle f'(\mathbf{x}_k) - f'(\mathbf{x}_{k-1}), \mathbf{p}_{k-1} \rangle}.$$

Finally, we can present the Conjugate Gradient Method

Conjugate Gradient Method	
Step 0:	Let $\mathbf{x}_0 \in \mathbb{R}^n$, compute $f(\mathbf{x}_0)$, $f'(\mathbf{x}_0)$ and set $\mathbf{p}_0 := f'(\mathbf{x}_0)$, $k := 0$
Step 1:	Find $\mathbf{x}_{k+1} := \mathbf{x}_k - h_k \mathbf{p}_k$ by “approximate line search” on the scalar h_k
Step 2:	Compute $f(\mathbf{x}_{k+1})$ and $f'(\mathbf{x}_{k+1})$
Step 3:	Compute the coefficient β_{k+1}
Step 4:	Set $\mathbf{p}_{k+1} := f'(\mathbf{x}_{k+1}) - \beta_{k+1} \mathbf{p}_k$, $k := k + 1$ and go to Step 1

The most popular choices for the coefficient β_k are:

1. *Hestenes-Stiefel (1952)*: $\beta_{k+1} = \frac{\langle f'(\mathbf{x}_{k+1}), f'(\mathbf{x}_{k+1}) - f'(\mathbf{x}_k) \rangle}{\langle f'(\mathbf{x}_{k+1}) - f'(\mathbf{x}_k), \mathbf{p}_k \rangle}$.
2. *Fletcher-Reeves (1964)*: $\beta_{k+1} = \frac{\|f'(\mathbf{x}_{k+1})\|^2}{\|f'(\mathbf{x}_k)\|^2}$.
3. *Polak-Ribière*: $\beta_{k+1} = \frac{\langle f'(\mathbf{x}_{k+1}), f'(\mathbf{x}_{k+1}) - f'(\mathbf{x}_k) \rangle}{\|f'(\mathbf{x}_k)\|^2}$.
4. *Polak-Ribière plus*: $\beta_{k+1} = \max \left\{ 0, \frac{\langle f'(\mathbf{x}_{k+1}), f'(\mathbf{x}_{k+1}) - f'(\mathbf{x}_k) \rangle}{\|f'(\mathbf{x}_k)\|^2} \right\}$.
5. *Dai-Yuan (1999)*: $\beta_{k+1} = \frac{\|f'(\mathbf{x}_{k+1})\|^2}{\langle f'(\mathbf{x}_{k+1}) - f'(\mathbf{x}_k), \mathbf{p}_k \rangle}$.

1.8 Quasi-Newton methods

The basic idea of the quasi-Newton methods is to approximate the Hessian matrix (or its inverse) which we need to compute in the Newton method. There are of course infinitely many ways to do so, but we choose the ones which satisfy the *secant equation*:

$$\mathbf{H}_{k+1} \mathbf{y}_k = \mathbf{s}_k$$

where $\mathbf{y}_k = f'(\mathbf{x}_{k+1}) - f'(\mathbf{x}_k)$, $\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$.

The general scheme of the quasi-Newton method is as follows.

Quasi-Newton Method	
Step 0:	Let $\mathbf{x}_0 \in \mathbb{R}^n$, $\mathbf{H}_0 := \mathbf{I}$, $k := 0$. Compute $f(\mathbf{x}_0)$, $f'(\mathbf{x}_0)$
Step 1:	Set $\mathbf{p}_k := \mathbf{H}_k f'(\mathbf{x}_k)$
Step 2:	Find $\mathbf{x}_{k+1} := \mathbf{x}_k - h_k \mathbf{p}_k$ by “approximate line search” on the scalar h_k
Step 3:	Compute $f(\mathbf{x}_{k+1})$ and $f'(\mathbf{x}_{k+1})$
Step 4:	Compute \mathbf{H}_{k+1} from \mathbf{H}_k , $k := k + 1$ and go to Step 1

The most popular updates for \mathbf{H}_{k+1} are:

1. *BFGS (Broyden-Fletcher-Goldfarb-Shanno)*

$$\mathbf{H}_{k+1} := \left(\mathbf{I} - \frac{\mathbf{s}_k (\mathbf{y}_k)^T}{(\mathbf{s}_k)^T \mathbf{y}_k} \right) \mathbf{H}_k \left(\mathbf{I} - \frac{\mathbf{y}_k (\mathbf{s}_k)^T}{(\mathbf{s}_k)^T \mathbf{y}_k} \right) + \frac{\mathbf{s}_k (\mathbf{s}_k)^T}{(\mathbf{s}_k)^T \mathbf{y}_k}$$

2. DFP (Davidon-Fletcher-Powell)

$$\mathbf{H}_{k+1} := \mathbf{H}_k + \frac{\mathbf{s}_k(\mathbf{s}_k)^T}{(\mathbf{y}_k)^T \mathbf{s}_k} - \frac{\mathbf{H}_k \mathbf{y}_k (\mathbf{y}_k)^T \mathbf{H}_k}{(\mathbf{y}_k)^T \mathbf{H}_k \mathbf{y}_k}$$

3. Symmetric-Rank-One

$$\mathbf{H}_{k+1} := \mathbf{H}_k + \frac{(\mathbf{s}_k - \mathbf{H}_k \mathbf{y}_k)(\mathbf{s}_k - \mathbf{H}_k \mathbf{y}_k)^T}{(\mathbf{s}_k - \mathbf{H}_k \mathbf{y}_k)^T \mathbf{y}_k}$$

In the same way for the conjugate gradient method, we can show that the quasi-Newton method converges in finite number of iterations for a strictly convex quadratic function. Moreover, under some strict convexity conditions at the neighborhood of the local minimum, it is possible to show that its iterates converge super-linearly [NOCEDAL2006].

1.9 Exercises

1. In view of Theorem 1.3.4, find a twice continuously differentiable function on \mathbb{R}^n which satisfies $f'(\mathbf{x}^*) = 0$, $f''(\mathbf{x}^*) \succeq \mathbf{O}$, but \mathbf{x}^* is not a local minimum of $f(\mathbf{x})$.
2. Prove Lemma 1.4.5.
3. Give a geometric interpretation of the following step-size strategies:

Let $0 < c_1 < c_2 < 1$,

- Wolfe condition

$$\begin{aligned} f(\mathbf{x}_k - hf'(\mathbf{x}_k)) &\leq f(\mathbf{x}_k) - c_1 h \|f'(\mathbf{x}_k)\|^2, \\ \langle f'(\mathbf{x}_k - hf'(\mathbf{x}_k)), f'(\mathbf{x}_k) \rangle &\leq c_2 \|f'(\mathbf{x}_k)\|^2. \end{aligned}$$

- Strong Wolfe condition

$$\begin{aligned} f(\mathbf{x}_k - hf'(\mathbf{x}_k)) &\leq f(\mathbf{x}_k) - c_1 h \|f'(\mathbf{x}_k)\|^2, \\ |\langle f'(\mathbf{x}_k - hf'(\mathbf{x}_k)), f'(\mathbf{x}_k) \rangle| &\leq c_2 \|f'(\mathbf{x}_k)\|^2. \end{aligned}$$

4. Consider a sequence $\{\beta_k\}_{k=0}^\infty$ which converges to zero.

The sequence is said to converge *Q-linearly* if there exists a scalar $\rho \in (0, 1)$ such that

$$\frac{\beta_{k+1}}{\beta_k} \leq \rho,$$

for all k sufficiently large. *Q-superlinear* convergence occurs when we have

$$\lim_{k \rightarrow \infty} \frac{\beta_{k+1}}{\beta_k} = 0,$$

while the convergence is *Q-quadratic* if there is a constant C such that

$$\frac{\beta_{k+1}}{\beta_k^2} \leq C$$

for all k sufficiently large. Q -superquadratic convergence is indicated by

$$\lim_{k \rightarrow \infty} \frac{\beta_{k+1}}{\beta_k^2} = 0.$$

(a) Show that the following implications are valid: Q -superquadratic \Rightarrow Q -quadratic \Rightarrow Q -superlinear \Rightarrow Q -linear.

(b) Give examples of sequences which do not imply the opposite directions in the three cases above.

A zero converging sequence $\{\beta_k\}_{k=0}^{\infty}$ is said to converge R -linearly if it is dominated by a Q -linearly converging sequence. That is, if there is a Q -linearly converging sequence $\{\hat{\beta}_k\}_{k=0}^{\infty}$ such that $0 \leq \beta_k \leq \hat{\beta}_k$.

(c) Give a sequence which is R -linearly converging but not Q -linearly converging.

5. Let $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Q} \mathbf{x}$ such that \mathbf{Q} is symmetric, and indefinite. Apply the gradient method with constant step. Show that if the starting point \mathbf{x}_0 belongs to the space spanned by the negative eigenvectors, the sequence generated by the gradient method diverges.
6. In light of Theorem 1.6.3, show that under Assumption 1.6.2, if we want to obtain $\|\mathbf{x}_k - \mathbf{x}^*\| < \varepsilon$, we need an order of $\ln(\ln \varepsilon^{-1})$ iterations for the Newton method.
7. In the Section 1.7, show that $\mathcal{L}_k = \{\boldsymbol{\delta}_0, \boldsymbol{\delta}_1, \dots, \boldsymbol{\delta}_{k-1}\}$.
8. In the same section, arrive at the expression (1.6) for a strictly convex quadratic function.
9. Show that the secant equation is valid for BFGS, DFP and symmetric-rank-one formulae.
10. Given $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ and a non-singular matrix $\mathbf{M} \in \mathbb{R}^{n \times n}$, if $1 + \mathbf{v}^T \mathbf{M}^{-1} \mathbf{u} \neq 0$, then the

$$(\mathbf{M} + \mathbf{u}\mathbf{v}^T)^{-1} = \mathbf{M}^{-1} - \frac{\mathbf{M}^{-1}\mathbf{u}\mathbf{v}^T\mathbf{M}^{-1}}{1 + \mathbf{v}^T\mathbf{M}^{-1}\mathbf{u}}. \quad (\text{Sherman-Morrison formula})$$

Apply this formula to compute the inverses \mathbf{B}_{k+1} of \mathbf{H}_{k+1} for BFGS, DFP and symmetric-rank-one formulae.

11. Apply the quasi-Newton method with BFGS, DFP, and Symmetric-Rank-One updates for the strictly convex function $f(\mathbf{x}) = \alpha + \langle \mathbf{a}, \mathbf{x} \rangle + \frac{1}{2}\langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle$ with $\mathbf{A} \succ \mathbf{O}$.

Chapter 2

Smooth Convex Optimization

2.1 Smooth convex functions

Definition 2.1.1 A continuously differentiable function $f(\mathbf{x})$ is called *convex* on \mathbb{R}^n (notation $\mathcal{F}^1(\mathbb{R}^n)$) if

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle f'(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

if $-f(\mathbf{x})$ is convex, $f(\mathbf{x})$ is called *concave*.

Theorem 2.1.2 If $f \in \mathcal{F}^1(\mathbb{R}^n)$ and $f'(\mathbf{x}^*) = 0$, then \mathbf{x}^* is the *global minimum* of $f(\mathbf{x})$ on \mathbb{R}^n .

Proof: Left for exercise. ■

Lemma 2.1.3 If $f \in \mathcal{F}^1(\mathbb{R}^m)$, $\mathbf{b} \in \mathbb{R}^m$, and $\mathbf{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, then

$$\phi(\mathbf{x}) = f(\mathbf{A}\mathbf{x} + \mathbf{b}) \in \mathcal{F}^1(\mathbb{R}^n).$$

Proof: Left for exercise. ■

Example 2.1.4 The following functions are differentiable and convex:

1. $f(x) = e^x$
2. $f(x) = |x|^p, \quad p > 1$
3. $f(x) = \frac{x^2}{1+|x|}$
4. $f(x) = |x| - \ln(1 + |x|)$
5. $f(\mathbf{x}) = \sum_{i=1}^m e^{\alpha_i + \langle \mathbf{a}_i, \mathbf{x} \rangle}$
6. $f(\mathbf{x}) = \sum_{i=1}^m |\langle \mathbf{a}_i, \mathbf{x} \rangle - b_i|^p, \quad p > 1$

Theorem 2.1.5 Let f be a continuously differentiable function. The following conditions are equivalent:

1. $f \in \mathcal{F}^1(\mathbb{R}^n)$.

2. $f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y}), \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \forall \alpha \in [0, 1].$
3. $\langle f'(\mathbf{x}) - f'(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq 0, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$

Proof: Left for exercise. ■

Theorem 2.1.6 Let f be a twice continuously differentiable function. Then $f \in \mathcal{F}^2(\mathbb{R}^n)$ if and only if

$$f''(\mathbf{x}) \succeq \mathbf{O}, \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

Proof: Let $f \in \mathcal{F}^2(\mathbb{R}^n)$, and denote $\mathbf{x}_\tau = \mathbf{x} + \tau \mathbf{s}$, $\tau > 0$. Then, from the previous result

$$\begin{aligned} 0 &\leq \frac{1}{\tau^2} \langle f'(\mathbf{x}_\tau) - f'(\mathbf{x}), \mathbf{x}_\tau - \mathbf{x} \rangle = \frac{1}{\tau} \langle f'(\mathbf{x}_\tau) - f'(\mathbf{x}), \mathbf{s} \rangle \\ &= \frac{1}{\tau} \int_0^\tau \langle f''(\mathbf{x} + \lambda \mathbf{s}) \mathbf{s}, \mathbf{s} \rangle d\lambda \\ &= \frac{F(\tau) - F(0)}{\tau} \end{aligned}$$

where $F(\tau) = \int_0^\tau \langle f''(\mathbf{x} + \lambda \mathbf{s}) \mathbf{s}, \mathbf{s} \rangle d\lambda$. Therefore, tending τ to 0, we get $0 \leq F'(0) = \langle f''(\mathbf{x}) \mathbf{s}, \mathbf{s} \rangle$, and we have the result.

Conversely, $\forall \mathbf{x} \in \mathbb{R}^n$,

$$\begin{aligned} f(\mathbf{y}) &= f(\mathbf{x}) + \langle f'(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \int_0^1 \int_0^\tau \langle f''(\mathbf{x} + \lambda(\mathbf{y} - \mathbf{x}))(\mathbf{y} - \mathbf{x}), \mathbf{y} - \mathbf{x} \rangle d\lambda d\tau \\ &\geq f(\mathbf{x}) + \langle f'(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle. \end{aligned}$$
■

Corollary 2.1.7 Let f be a two times continuously differentiable function. $f \in \mathcal{F}_L^{2,1}(\mathbb{R}^n)$ if and only if $\mathbf{O} \preceq f''(\mathbf{x}) \preceq L\mathbf{I}$, $\forall \mathbf{x} \in \mathbb{R}^n$.

Theorem 2.1.8 Let f be a continuously differentiable function in \mathbb{R}^n , $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, and $\alpha \in [0, 1]$. Then the following conditions are equivalent:

1. $f \in \mathcal{F}_L^{1,1}(\mathbb{R}^n)$.
2. $0 \leq f(\mathbf{y}) - f(\mathbf{x}) - \langle f'(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \leq \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2$.
3. $f(\mathbf{x}) + \langle f'(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2L} \|f'(\mathbf{x}) - f'(\mathbf{y})\|^2 \leq f(\mathbf{y})$.
4. $0 \leq \frac{1}{L} \|f'(\mathbf{x}) - f'(\mathbf{y})\|^2 \leq \langle f'(\mathbf{x}) - f'(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$.
5. $0 \leq \langle f'(\mathbf{x}) - f'(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \leq L \|\mathbf{x} - \mathbf{y}\|^2$.
6. $f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) + \frac{\alpha(1-\alpha)}{2L} \|f'(\mathbf{x}) - f'(\mathbf{y})\|^2 \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y})$.
7. $0 \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y}) - f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \leq \alpha(1 - \alpha) \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2$.