

4.4.2 The Newton Method

Example 4.19 Let us apply the Newton method to find the root of the following function

$$g(x) = \frac{x}{\sqrt{1+x^2}}.$$

Clearly $x^* = 0$.

The Newton method will give:

$$x_{k+1} = x_k - \frac{g(x_k)}{g'(x_k)} = x_k - x_k(1+x_k^2) = -x_k^3.$$

Therefore, the method converges if $|x_0| < 1$, it oscillates if $|x_0| = 1$, and finally, diverges if $|x_0| > 1$.

Assumption 4.20

1. $f \in \mathcal{C}_M^{2,2}(\mathbb{R}^n)$;
2. There is a local minimum \mathbf{x}^* of the function $f(\mathbf{x})$;
3. The Hessian is positive definite at \mathbf{x}^* :

$$\nabla^2 f(\mathbf{x}^*) \succeq \ell \mathbf{I}, \quad \ell > 0;$$

4. Our starting point \mathbf{x}_0 is close enough to \mathbf{x}^* .

Theorem 4.21 Let the function $f(\mathbf{x})$ satisfy the above assumptions. Suppose that the initial starting point \mathbf{x}_0 is close enough to \mathbf{x}^* :

$$\|\mathbf{x}_0 - \mathbf{x}^*\|_2 < \bar{r} := \frac{2\ell}{3M}.$$

Then $\|\mathbf{x}_k - \mathbf{x}^*\|_2 < \bar{r}$ for all k of the Newton method and it converges (Q-)quadratically:

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_2 \leq \frac{M\|\mathbf{x}_k - \mathbf{x}^*\|_2^2}{2(\ell - M\|\mathbf{x}_k - \mathbf{x}^*\|_2)}.$$

Proof:

Let $r_k = \|\mathbf{x}_k - \mathbf{x}^*\|_2$. From Lemma 3.8 and the assumption, we have for $k = 0$,

$$\nabla^2 f(\mathbf{x}_0) \succeq \nabla^2 f(\mathbf{x}^*) - Mr_0 \mathbf{I} \succeq (\ell - Mr_0) \mathbf{I}. \quad (8)$$

Since $r_0 < \bar{r} = \frac{2\ell}{3M} < \frac{\ell}{M}$, we have $\ell - Mr_0 > 0$ and therefore, $\nabla^2 f(\mathbf{x}_0)$ is invertible.

Consider the Newton method for $k = 0$, $\mathbf{x}_1 = \mathbf{x}_0 - [\nabla^2 f(\mathbf{x}_0)]^{-1} \nabla f(\mathbf{x}_0)$.

Then

$$\begin{aligned} \mathbf{x}_1 - \mathbf{x}^* &= \mathbf{x}_0 - \mathbf{x}^* - [\nabla^2 f(\mathbf{x}_0)]^{-1} \nabla f(\mathbf{x}_0) \\ &= \mathbf{x}_0 - \mathbf{x}^* - [\nabla^2 f(\mathbf{x}_0)]^{-1} \int_0^1 \nabla^2 f(\mathbf{x}^* + \tau(\mathbf{x}_0 - \mathbf{x}^*)) (\mathbf{x}_0 - \mathbf{x}^*) d\tau \\ &= [\nabla^2 f(\mathbf{x}_0)]^{-1} \mathbf{G}_0 (\mathbf{x}_0 - \mathbf{x}^*) \end{aligned}$$

where $\mathbf{G}_0 = \int_0^1 [\nabla^2 f(\mathbf{x}_0) - \nabla^2 f(\mathbf{x}^* + \tau(\mathbf{x}_0 - \mathbf{x}^*))] d\tau$.

Then

$$\begin{aligned}
\|\mathbf{G}_0\|_2 &= \left\| \int_0^1 [\nabla^2 \mathbf{f}(\mathbf{x}_0) - \nabla^2 \mathbf{f}(\mathbf{x}^* + \tau(\mathbf{x}_0 - \mathbf{x}^*))] d\tau \right\|_2 \\
&\leq \int_0^1 \|\nabla^2 \mathbf{f}(\mathbf{x}_0) - \nabla^2 \mathbf{f}(\mathbf{x}^* + \tau(\mathbf{x}_0 - \mathbf{x}^*))\|_2 d\tau \\
&\leq \int_0^1 M|1 - \tau|r_0 d\tau = \frac{r_0}{2}M.
\end{aligned}$$

From (8),

$$\|[\nabla^2 \mathbf{f}(\mathbf{x}_0)]^{-1}\|_2 \leq (\ell - Mr_0)^{-1}.$$

Then

$$r_1 \leq \frac{Mr_0^2}{2(\ell - Mr_0)}.$$

Since $r_0 < \bar{r} = \frac{2\ell}{3M}$, $\frac{Mr_0}{2(\ell - Mr_0)} < 1$, and $r_1 < r_0$.

One can see now that the same argument is valid for all k 's. ■

- Comparing this result with the rate of convergence of the steepest descent, we see that the Newton method is much faster.
- Surprisingly, the region of *quadratic convergence* of the Newton method is almost the same as the region of the *linear convergence* of the gradient method.

$$\|\mathbf{x}_0 - \mathbf{x}^*\|_2 < \frac{2\ell}{M} \quad (\text{steepest descent method}) \quad \|\mathbf{x}_0 - \mathbf{x}^*\|_2 < \frac{2\ell}{3M} \quad (\text{Newton method})$$

- This justifies a standard recommendation to use the steepest descent method only at the initial stage of the minimization process in order to get close to a local minimum and then perform the Newton method to refine.

4.4.3 The Conjugate Gradient Methods

The conjugate gradient methods were initially proposed for minimizing convex quadratic functions. Consider the problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

with $f(\mathbf{x}) = \alpha + \langle \mathbf{a}, \mathbf{x} \rangle + \frac{1}{2} \langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle$ and $\mathbf{A} \succ \mathbf{O}$. Since its minimal solution is $\mathbf{x}^* = -\mathbf{A}^{-1}\mathbf{a}$, we can rewrite $f(\mathbf{x})$ as:

$$\begin{aligned}
f(\mathbf{x}) &= \alpha - \langle \mathbf{A}\mathbf{x}^*, \mathbf{x} \rangle + \frac{1}{2} \langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle \\
&= \alpha - \frac{1}{2} \langle \mathbf{A}\mathbf{x}^*, \mathbf{x}^* \rangle + \frac{1}{2} \langle \mathbf{A}(\mathbf{x} - \mathbf{x}^*), \mathbf{x} - \mathbf{x}^* \rangle.
\end{aligned}$$

Thus, $f(\mathbf{x}^*) = \alpha - \frac{1}{2} \langle \mathbf{A}\mathbf{x}^*, \mathbf{x}^* \rangle$ and $\nabla f(\mathbf{x}) = \mathbf{A}(\mathbf{x} - \mathbf{x}^*)$.

Definition 4.22 Given a starting point \mathbf{x}_0 , the linear *Krylov subspaces* is defined as

$$\mathcal{L}_k := \text{span}\{\mathbf{A}(\mathbf{x}_0 - \mathbf{x}^*), \dots, \mathbf{A}^k(\mathbf{x}_0 - \mathbf{x}^*)\}, \quad k \geq 1,$$

where $\text{span}\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_p\}$ is the linear subspace of \mathbb{R}^n spanned by the vectors $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_p \in \mathbb{R}^n$.

We claim temporarily that the sequence of points generated by a *conjugate gradient method* is defined as follows:

$$\mathbf{x}_k := \arg \min \{f(\mathbf{x}) \mid \mathbf{x} \in \mathbf{x}_0 + \mathcal{L}_k\}, \quad k \geq 1.$$

Lemma 4.23 For any $k \geq 1$, $\mathcal{L}_k = \text{span}\{\nabla f(\mathbf{x}_0), \dots, \nabla f(\mathbf{x}_{k-1})\}$.

Proof:

Let us prove by induction hypothesis.

For $k = 1$, the statement is true since $\nabla f(\mathbf{x}_0) = \mathbf{A}(\mathbf{x}_0 - \mathbf{x}^*)$.

Suppose the claim is true for some $k \geq 1$. Then from the definition of the conjugate gradient method,

$$\mathbf{x}_k = \mathbf{x}_0 + \sum_{i=1}^k \lambda_i \mathbf{A}^i(\mathbf{x}_0 - \mathbf{x}^*)$$

with some $\lambda_i \in \mathbb{R}$, $i = 1, \dots, k$. Therefore,

$$\nabla f(\mathbf{x}_k) = \mathbf{A}(\mathbf{x}_0 - \mathbf{x}^*) + \sum_{i=1}^k \lambda_i \mathbf{A}^{i+1}(\mathbf{x}_0 - \mathbf{x}^*) = \mathbf{A}(\mathbf{x}_0 - \mathbf{x}^*) + \sum_{i=1}^{k-1} \lambda_i \mathbf{A}^{i+1}(\mathbf{x}_0 - \mathbf{x}^*) + \lambda_k \mathbf{A}^{k+1}(\mathbf{x}_0 - \mathbf{x}^*).$$

The first two terms of the last expression belongs to \mathcal{L}_k from the definition. And then,

$$\text{span}\{\mathcal{L}_k, \nabla f(\mathbf{x}_k)\} \subseteq \text{span}\{\mathcal{L}_k, \mathbf{A}^{k+1}(\mathbf{x}_0 - \mathbf{x}^*)\} = \mathcal{L}_{k+1}.$$

There are two ways to show that the equality holds. Assume that $\mathbf{A}^{k+1}(\mathbf{x}_0 - \mathbf{x}^*) \in \mathcal{L}_k$. Then it is obvious and $\mathcal{L}_k = \mathcal{L}_{k+1}$. If $\mathbf{A}^{k+1}(\mathbf{x}_0 - \mathbf{x}^*) \notin \mathcal{L}_k$, the equality holds unless $\lambda_k = 0$. However, this possibility implies that $\mathbf{x}_k \in \mathcal{L}_{k-1}$, $\mathbf{x}_{k-1} = \mathbf{x}_k$ and therefore, $\mathcal{L}_{k-1} = \mathcal{L}_k = \mathcal{L}_{k+1}$ again.

An alternative way is to use contradiction. If the equality does not hold, $\nabla f(\mathbf{x}_k) \in \mathcal{L}_k$ implies $\mathbf{A}^{k+1}(\mathbf{x}_0 - \mathbf{x}^*) \in \mathcal{L}_k$, which again implies the equality, or $\lambda_k = 0$, which implies that $\mathbf{x}_k = \mathbf{x}_{k-1}$ (algorithm terminated). ■

Lemma 4.24 For any $k, \ell \geq 0$, $k \neq \ell$, we have $\langle \nabla f(\mathbf{x}_k), \nabla f(\mathbf{x}_\ell) \rangle = 0$.

Proof:

Let $k \geq i$, and consider

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

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

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

Corollary 4.25 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 $\mathcal{L}_k = \text{span}\{\boldsymbol{\delta}_0, \boldsymbol{\delta}_1, \dots, \boldsymbol{\delta}_{k-1}\}$ (Exercise 10).

Lemma 4.26 For any $k, \ell \geq 0$, $k \neq \ell$, $\langle \mathbf{A}\boldsymbol{\delta}_k, \boldsymbol{\delta}_\ell \rangle = 0$.

Proof:

Left for exercise. ■

The vectors $\{\delta_i\}$ 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:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - h_k \nabla \mathbf{f}(\mathbf{x}_k) + \sum_{j=0}^{k-1} \lambda_j \delta_j$$

Using the previous properties, we arrive that (see Exercise 11)

$$\lambda_j = 0, \quad (j = 0, 1, \dots, k-2), \quad \lambda_{k-1} = \frac{h_k \|\nabla \mathbf{f}(\mathbf{x}_k)\|_2^2}{\langle \nabla \mathbf{f}(\mathbf{x}_k) - \nabla \mathbf{f}(\mathbf{x}_{k-1}), \delta_{k-1} \rangle}. \quad (9)$$

Thus

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

where

$$\mathbf{p}_k = \nabla \mathbf{f}(\mathbf{x}_k) - \frac{\|\nabla \mathbf{f}(\mathbf{x}_k)\|_2^2 \mathbf{p}_{k-1}}{\langle \nabla \mathbf{f}(\mathbf{x}_k) - \nabla \mathbf{f}(\mathbf{x}_{k-1}), \mathbf{p}_{k-1} \rangle}.$$

Finally, we can present the Conjugate Gradient Method

Conjugate Gradient Method	
Step 0:	Let $\mathbf{x}_0 \in \mathbb{R}^n$, compute $f(\mathbf{x}_0)$, $\nabla \mathbf{f}(\mathbf{x}_0)$ and set $\mathbf{p}_0 := \nabla \mathbf{f}(\mathbf{x}_0)$, $k := 0$
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 $\nabla \mathbf{f}(\mathbf{x}_{k+1})$
Step 3:	Compute the coefficient β_{k+1}
Step 4:	Set $\mathbf{p}_{k+1} := \nabla \mathbf{f}(\mathbf{x}_{k+1}) - \beta_{k+1} \mathbf{p}_k$, $k := k + 1$ and go to Step 1

The most popular choices for the coefficient β_k are:

1. *Hestenes-Stiefel (1952)*: $\beta_{k+1} = \frac{\langle \nabla \mathbf{f}(\mathbf{x}_{k+1}), \nabla \mathbf{f}(\mathbf{x}_{k+1}) - \nabla \mathbf{f}(\mathbf{x}_k) \rangle}{\langle \nabla \mathbf{f}(\mathbf{x}_{k+1}) - \nabla \mathbf{f}(\mathbf{x}_k), \mathbf{p}_k \rangle}.$
2. *Fletcher-Reeves (1964)*: $\beta_{k+1} = \frac{\|\nabla \mathbf{f}(\mathbf{x}_{k+1})\|_2^2}{\|\nabla \mathbf{f}(\mathbf{x}_k)\|_2^2}.$
3. *Polak-Ribière (1969)*: $\beta_{k+1} = \frac{\langle \nabla \mathbf{f}(\mathbf{x}_{k+1}), \nabla \mathbf{f}(\mathbf{x}_{k+1}) - \nabla \mathbf{f}(\mathbf{x}_k) \rangle}{\|\nabla \mathbf{f}(\mathbf{x}_k)\|_2^2}.$
4. *Polak-Ribière plus*: $\beta_{k+1} = \max \left\{ 0, \frac{\langle \nabla \mathbf{f}(\mathbf{x}_{k+1}), \nabla \mathbf{f}(\mathbf{x}_{k+1}) - \nabla \mathbf{f}(\mathbf{x}_k) \rangle}{\|\nabla \mathbf{f}(\mathbf{x}_k)\|_2^2} \right\}.$
5. *Dai-Yuan (1999)*: $\beta_{k+1} = \frac{\|\nabla \mathbf{f}(\mathbf{x}_{k+1})\|_2^2}{\langle \nabla \mathbf{f}(\mathbf{x}_{k+1}) - \nabla \mathbf{f}(\mathbf{x}_k), \mathbf{p}_k \rangle}.$

Among them, Hestenes-Stiefel and Polak-Ribière are empirically preferred.