Then

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

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

$$\begin{split} \|\boldsymbol{G}_{0}\|_{2} &= \left\| \int_{0}^{1} [\boldsymbol{\nabla}^{2} \boldsymbol{f}(\boldsymbol{x}_{0}) - \boldsymbol{\nabla}^{2} \boldsymbol{f}(\boldsymbol{x}^{*} + \tau(\boldsymbol{x}_{0} - \boldsymbol{x}^{*}))] d\tau \right\|_{2} \\ &\leq \int_{0}^{1} \|\boldsymbol{\nabla}^{2} \boldsymbol{f}(\boldsymbol{x}_{0}) - \boldsymbol{\nabla}^{2} \boldsymbol{f}(\boldsymbol{x}^{*} + \tau(\boldsymbol{x}_{0} - \boldsymbol{x}^{*}))\|_{2} d\tau \\ &\leq \int_{0}^{1} M |1 - \tau| r_{0} d\tau = \frac{r_{0}}{2} M. \end{split}$$

From (8),

$$\|[\nabla^2 f(x_0)]^{-1}\|_2 \le (\ell - Mr_0)^{-1}.$$

Then

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

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.

$$\|\boldsymbol{x}_0 - \boldsymbol{x}^*\|_2 < rac{2\ell}{M}$$
 (steepest descent method) $\|\boldsymbol{x}_0 - \boldsymbol{x}^*\|_2 < rac{2\ell}{3M}$ (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

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

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

$$f(\boldsymbol{x}) = \alpha - \langle \boldsymbol{A}\boldsymbol{x}^*, \boldsymbol{x} \rangle + \frac{1}{2} \langle \boldsymbol{A}\boldsymbol{x}, \boldsymbol{x} \rangle$$

= $\alpha - \frac{1}{2} \langle \boldsymbol{A}\boldsymbol{x}^*, \boldsymbol{x}^* \rangle + \frac{1}{2} \langle \boldsymbol{A}(\boldsymbol{x} - \boldsymbol{x}^*), \boldsymbol{x} - \boldsymbol{x}^* \rangle.$

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

Definition 4.22 Given a starting point x_0 , the linear *Krylov subspaces* is defined as

$$\mathcal{L}_k := \text{span}\{A(x_0 - x^*), \dots, A^k(x_0 - x^*)\}, \quad k \ge 1,$$

where span $\{a_1, a_2, \ldots, a_p\}$ is the linear subspace of \mathbb{R}^n spanned by the vectors $a_1, a_2, \ldots, a_p \in \mathbb{R}^n$.

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

$$\boldsymbol{x}_k := rg\min\{f(\boldsymbol{x}) \mid \boldsymbol{x} \in \boldsymbol{x}_0 + \mathcal{L}_k\}, \ k \ge 1.$$

Lemma 4.23 For any $k \geq 1$, $\mathcal{L}_k = \operatorname{span}\{\nabla f(x_0), \ldots, \nabla f(x_{k-1})\}.$

Proof:

Let us prove by induction hypothesis.

For k = 1, the statement is true since $\nabla f(x_0) = A(x_0 - x^*)$.

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

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

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

$$\nabla f(x_k) = A(x_0 - x^*) + \sum_{i=1}^k \lambda_i A^{i+1}(x_0 - x^*) = A(x_0 - x^*) + \sum_{i=1}^{k-1} \lambda_i A^{i+1}(x_0 - x^*) + \lambda_k A^{k+1}(x_0 - x^*).$$

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

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

There are two ways to show that the equality holds. Assume that $A^{k+1}(x_0 - x^*) \in \mathcal{L}_k$. Then it is obvious and $\mathcal{L}_k = \mathcal{L}_{k+1}$. If $A^{k+1}(x_0 - x^*) \notin \mathcal{L}_k$, the equality holds unless $\lambda_k = 0$. However, this possibility implies that $x_k \in \mathcal{L}_{k-1}$, $x_{k-1} = x_k$ 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, $\nabla f(x_k) \in \mathcal{L}_k$ implies $A^{k+1}(x_0 - x^*) \in \mathcal{L}_k$, which again implies the equality, or $\lambda_k = 0$, which implies that $x_k = x_{k-1}$ (algorithm terminated).

Lemma 4.24 For any $k, \ell \ge 0, k \ne \ell$, we have $\langle \nabla f(x_k), \nabla f(x_\ell) \rangle = 0$.

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

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

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

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

Corollary 4.25 The sequence generated by the conjugate gradient method for the convex quadratic function is finite.