

Introductory Lectures on Optimization

Descent Method (2)

Hui Qian

qianhui@zju.edu.cn

College of Computer Science, Zhejiang University

October 23, 2024

Outline

- 1 Newton Method
 - Basic Scheme
 - Convergence Analysis
 - Variable Metric
- 2 Conjugate Gradients
 - Historical Origin
 - Fundamental Theory
 - CG Algorithm
- 3 Reference

Part I

Newton Descent

Historical Origins

The Newton method is widely known as a technique for **finding a root of a function of one variable**. Let $\phi(t) : \mathbb{R} \rightarrow \mathbb{R}$. Consider the equation

$$\phi(t^*) = 0.$$

The Newton method is based on **linear approximation**. Assume that we get some t close enough to t^* . Note that

$$\phi(\underbrace{t + \Delta t}_{\text{expect that will be } t^*}) = \phi(t) + \phi'(t)\Delta t + o(|\Delta t|).$$

Historical Origins

Therefore the equation $\phi(t + \Delta t) = 0$ can be approximated by the following **linear** equation:

$$\phi(t) + \phi'(t)\Delta t = 0.$$

We can expect that the solution of this equation, the displacement Δt is a good approximation to the **optimal displacement** $\Delta t^* = t^* - t$. Converting this idea in an algorithmic form, we get the process

$$t_{k+1} = t_k - \frac{\phi(t_k)}{\phi'(t_k)}.$$

Historical Origins

This scheme can be naturally extended onto the problem of finding solution to a **system of nonlinear equations**

$$F(\mathbf{x}) = 0$$

where $\mathbf{x} \in \mathbb{R}^n$ and $F(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^n$. In this case we have to define the displacement $\Delta\mathbf{x}$ as a solution to the following system of linear equations:

$$F(\mathbf{x}) + J_F(\mathbf{x})\Delta\mathbf{x} = 0$$

(it is called the **Newton system**). If the **Jacobian** $J_F(\mathbf{x})$ is **nongenerate**, we can compute displacement $\Delta\mathbf{x} = -[J_F(\mathbf{x}_k)]^{-1}F(\mathbf{x}_k)$. The corresponding iterative scheme looks as follows:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [J_F(\mathbf{x}_k)]^{-1}F(\mathbf{x}_k).$$

Basic Scheme

Finally, in view of the first order condition of continuously differentiable function , we can replace the **unconstrained minimization** problem by a problem of **finding roots of the nonlinear system**

$$\nabla f(\mathbf{x}) = 0. \quad (28)$$

Further, for solving (28) we can apply a standard Newton method for systems of nonlinear equations. In this case, the Newton system looks as follows:

$$\nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x})\Delta\mathbf{x} = 0.$$

Hence, the Newton method for optimization problems appears to be in the form

$$\boxed{\mathbf{x}_{k+1} = \mathbf{x}_k - [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k).} \quad (29)$$

Basic Scheme

Note that we can also obtain the process (29), using the idea of quadratic approximation. Consider this approximation, computed with respect to the point \mathbf{x}_k :

$$\phi(\mathbf{x}) = f(\mathbf{x}_k) + \langle \nabla f(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}_k)(\mathbf{x} - \mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle.$$

Assume that $\nabla^2 f(\mathbf{x}_k) \succ 0$. Then we can choose \mathbf{x}_{k+1} as a point of minimum of the quadratic function $\phi(\mathbf{x})$. This means that

$$\nabla \phi(\mathbf{x}_{k+1}) = \nabla f(\mathbf{x}_k) + \nabla^2 f(\mathbf{x}_k)(\mathbf{x}_{k+1} - \mathbf{x}_k) = 0,$$

and we come again to the Newton process (29).

Basic Scheme

Remark.

The convergence of the Newton method in a neighborhood of a strict local minimum is very fast.

Two serious drawbacks.

- Firstly, it can break down if $\nabla^2 f(x_k)$ is degenerate (a square matrix whose determinant is equal to zero).
- Secondly, the Newton process can diverge.

Basic Scheme

Example 7

Let us apply the Newton method for finding a root of the following function of one variable:

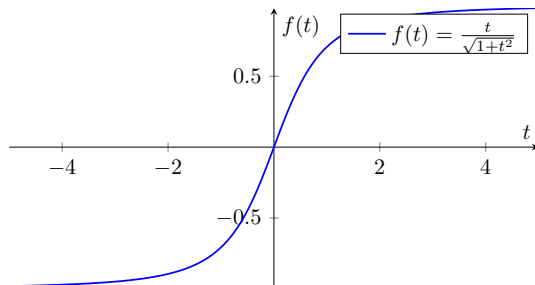
$$\phi(t) = \frac{t}{\sqrt{1+t^2}} \quad \text{and} \quad \phi'(t) = \frac{1}{[1+t^2]^{3/2}}.$$

Clearly, $t^* = 0$. Therefore the Newton process looks as follows:

$$t_{k+1} = t_k - \frac{\phi(t_k)}{\phi'(t_k)} = t_k - \frac{t_k}{\sqrt{1+t_k^2}} \cdot [1+t_k^2]^{3/2} = -t_k^3.$$

Thus, if $|t_0| < 1$, then this method converges and the convergence is extremely fast. The points ± 1 are the oscillation points of this method. If $|t_0| > 1$, then the method diverges.

Basic Scheme



$$\phi(t) = \frac{t}{\sqrt{1+t^2}}.$$

Damped Newton Method

In order to avoid a possible **divergence**, in practice we can apply a **damped Newton method**:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - h_k [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k),$$

where, $h_k > 0$ is a step-size parameter.

Remark.

- At the initial stage of the method we can use the same step size strategies as for the gradient scheme.
- At the final stage it is reasonable to choose $h_k = 1$.

Local Convergence of The Newton Method.

Let us study the local convergence of the Newton method. Consider the problem

- 1 $f \in C_M^{2,2}(\mathbb{R}^n)$,
- 2 There exists a local minimum of function f with positive definite Hessian:

$$\nabla^2 f(\mathbf{x}^*) \succeq lI_n, \quad l > 0. \quad (30)$$

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

Consider the process:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k).$$

Local Convergence of The Newton Method.

Using the same reasoning as for the gradient method, we obtain the following representation:

$$\begin{aligned}
 \boxed{\mathbf{x}_{k+1} - \mathbf{x}^*} &= \mathbf{x}_k - \mathbf{x}^* - [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k) \quad \text{by definition} \\
 &= \mathbf{x}_k - \mathbf{x}^* - [\nabla^2 f(\mathbf{x}_k)]^{-1} \int_0^1 \nabla^2 f(\mathbf{x}^* + \tau(\mathbf{x}_k - \mathbf{x}^*)) (\mathbf{x}_k - \mathbf{x}^*) d\tau \\
 &\quad \text{since we have } \nabla f(\mathbf{x}^*) = 0 \\
 &= [\nabla^2 f(\mathbf{x}_k)]^{-1} G_k \cdot (\boxed{\mathbf{x}_k - \mathbf{x}^*}), \tag{31}
 \end{aligned}$$

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

We consider the contracting mappings, thus both G_k and $[\nabla^2 f(\mathbf{x}_k)]^{-1}$ should be explored.

Bound for G_k

Denote $r_k = \|\mathbf{x}_k - \mathbf{x}^*\|$. Then

$$\begin{aligned}\|G_k\| &= \left\| \int_0^1 [\nabla^2 f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}^* + \tau(\mathbf{x}_k - \mathbf{x}^*))] d\tau \right\| \\ &\leq \int_0^1 \|\nabla^2 f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}^* + \tau(\mathbf{x}_k - \mathbf{x}^*))\| d\tau \\ &\leq \int_0^1 M(1 - \tau)r_k d\tau = \frac{r_k}{2}M. \quad C_M^{2,2}(\mathbb{R}^n) \text{ Lipschitz condition.}\end{aligned}$$

Bound for $[\nabla^2 f(x_k)]^{-1}$

In view of Corollary 1.2.2 of Nesterov [2003], we have for $f \in C_M^{2,2}(\mathbb{R}^n)$ and $\|\mathbf{y} - \mathbf{x}\| = r$,

$$\nabla^2 f(\mathbf{x}) - MrI_n \preceq \nabla^2 f(\mathbf{y}) \preceq \nabla^2 f(\mathbf{x}) + MrI_n,$$

where I_n is the identity matrix on \mathbb{R}^n .

In view of Corollary (1.2.2) and (30), we have

$$\nabla^2 f(\mathbf{x}_k) \succeq \nabla^2 f(\mathbf{x}^*) - Mr_k I_n \succeq (l - Mr_k)I_n.$$

Therefore, if $r_k < \frac{l}{M}$, then $\nabla^2 f(\mathbf{x}_k)$ is positive definite and

$$\|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| \leq (l - Mr_k)^{-1}.$$

Local Convergence of The Newton Method.

$$r_{k+1} = [\nabla^2 f(\mathbf{x}_k)]^{-1} G_k \cdot r_k,$$

$$\|G_k\| \leq \frac{r_k}{2} M,$$

$$\|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| \leq (l - Mr_k)^{-1}.$$

Hence, for r_k small enough ($r_k < \frac{2l}{3M}$), we have

$$r_{k+1} \leq \frac{M}{2(l - Mr_k)} r_k^2 \quad (< r_k).$$

The rate of convergence of this type is called **quadratic**.

Local Convergence of The Newton Method.

Thus, we have proved the following theorem.

Theorem 8

Let function $f(x)$ satisfy our assumptions. Suppose that the initial starting point x_0 is close enough to x^* :

$$\|x_0 - x^*\| < \bar{r} = \frac{2l}{3M}.$$

Then, $\|x_k - x^*\| < \bar{r}$ for all k , and the Newton method converges quadratically:

$$\|x_{k+1} - x^*\| \leq \frac{M \|x_k - x^*\|^2}{2(l - M \|x_k - x^*\|)}.$$

Other Approximation derived Methods: Variable Metric

Revisit GD: Let us fix some $\bar{x} \in \mathbb{R}^n$. Consider the following approximation of the function:

$$\phi_1(\mathbf{x}) = f(\bar{\mathbf{x}}) + \langle \nabla f(\bar{\mathbf{x}}), \mathbf{x} - \bar{\mathbf{x}} \rangle + \frac{1}{2h} \|\mathbf{x} - \bar{\mathbf{x}}\|^2,$$

where the parameter h is positive.

The first-order optimality condition provides us with the following equation for \mathbf{x}_1^* , the unconstrained minimum of the function $\phi_1(\mathbf{x})$:

$$\nabla \phi_1(\mathbf{x}_1^*) = \nabla f(\bar{\mathbf{x}}) + \frac{1}{h}(\mathbf{x}_1^* - \bar{\mathbf{x}}) = 0.$$

Thus, $\mathbf{x}_1^* = \bar{\mathbf{x}} - h \nabla f(\bar{\mathbf{x}})$. That is exactly the iterate of the **gradient descent method**.

Other Approximation derived Methods: Variable Metric

Revisit Newton: Further, consider a quadratic approximation of function $f(\mathbf{x})$:

$$\phi_2(\mathbf{x}) = f(\bar{\mathbf{x}}) + \langle \nabla f(\bar{\mathbf{x}}), \mathbf{x} - \bar{\mathbf{x}} \rangle + \frac{1}{2} \langle \nabla^2 f(\bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}}), \mathbf{x} - \bar{\mathbf{x}} \rangle.$$

We have already seen that the minimum of this function is

$$\mathbf{x}_2^* = \bar{\mathbf{x}} - [\nabla^2 f(\bar{\mathbf{x}})]^{-1} \nabla f(\bar{\mathbf{x}}),$$

and that is exactly the iterate of the Newton method.

逆矩阵的计算复杂度大约是 $O(n^3)$ 。迭代次数虽少，但是每次迭代开销大。

Other Approximation derived Methods: Variable Metric

Lightweight Newton: Thus, we can try to use some approximations of function $f(\mathbf{x})$, which are better than $\phi_1(\mathbf{x})$ and which are less expensive than $\phi_2(\mathbf{x})$.

Let G be a positive definite $n \times n$ matrix. Denote

$$\phi_G(\mathbf{x}) = f(\bar{\mathbf{x}}) + \langle \nabla f(\bar{\mathbf{x}}), \mathbf{x} - \bar{\mathbf{x}} \rangle + \frac{1}{2} \langle G(\mathbf{x} - \bar{\mathbf{x}}), \mathbf{x} - \bar{\mathbf{x}} \rangle.$$

Its minimum we obtain from the above equation is

$$\mathbf{x}_G^* = \bar{\mathbf{x}} - G^{-1} \nabla f(\bar{\mathbf{x}}). \quad (32)$$

Other Approximation derived Methods: Variable Metric

Variable Metric: The first-order methods, which form a sequence of matrices

$$\{G_k\} : G_k \rightarrow \nabla^2 f(\mathbf{x}^*)$$

or

$$\{H_k\} : H_k \equiv G_k^{-1} \rightarrow [\nabla^2 f(\mathbf{x}^*)]^{-1},$$

are called the **variable metric** methods (sometime, use **Quasi-Newton** instead).

In these methods **only** the gradients are involved in the process of generating the sequences $\{G_k\}$ or $\{H_k\}$.

Other Approximation derived Methods: Variable Metric

New Inner Product: Note that the **gradient** and the **Hessian** of a nonlinear function $f(\mathbf{x})$ are defined with respect to a standard Euclidean inner product on \mathbb{R}^n :

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^n \mathbf{x}^{(i)} \mathbf{y}^{(i)}, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \quad \|\mathbf{x}\| = \langle \mathbf{x}, \mathbf{x} \rangle^{1/2}.$$

Indeed, the definition of the gradient is

$$f(\mathbf{x} + h) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), h \rangle + o(\|h\|),$$

and from this equation we derive its coordinate representation:

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}^{(1)}}, \dots, \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}^{(n)}} \right)^\top.$$

Other Approximation derived Methods: Variable Metric

New Inner Product: Consider a symmetric positive definite $n \times n$ -matrix A . For $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ denote

$$\langle \mathbf{x}, \mathbf{y} \rangle_A = \langle A\mathbf{x}, \mathbf{y} \rangle, \quad \|\mathbf{x}\|_A = \langle A\mathbf{x}, \mathbf{x} \rangle^{1/2}.$$

The function $\|\mathbf{x}\|_A$ is a new norm on \mathbb{R}^n . Note that topologically this new metric is equivalent to the old one:

$$\lambda_n(A)^{1/2} \|\mathbf{x}\| \leq \|\mathbf{x}\|_A \leq \lambda_1(A)^{1/2} \|\mathbf{x}\|,$$

where $\lambda_n(A)$ and $\lambda_1(A)$ are the smallest and largest eigenvalues of the matrix A^*A respectively.

Other Approximation derived Methods: Variable Metric

New Inner Product: However, the gradient and the Hessian, computed with respect to the new inner product are changing:

$$\begin{aligned} f(\mathbf{x} + h) &= f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), h \rangle + \frac{1}{2} \langle \nabla^2 f(\mathbf{x}) h, h \rangle + o(\|h\|) \\ &= f(\mathbf{x}) + \langle A^{-1} \nabla f(\mathbf{x}), h \rangle_A + \frac{1}{2} \langle A^{-1} \nabla^2 f(\mathbf{x}) h, h \rangle_A + o(\|h\|_A). \end{aligned}$$

Hence, $\nabla_A f(\mathbf{x}) = A^{-1} \nabla f(\mathbf{x})$ is the new **gradient** and $\nabla_A^2 f(\mathbf{x}) = A^{-1} \nabla^2 f(\mathbf{x})$ is the new **Hessian**. Thus, the direction used in the Newton method $[\nabla^2 f(\mathbf{x})]^{-1} \nabla f(\mathbf{x})$ can be seen as a gradient computed with respect to the metric defined by $A = \nabla^2 f(\mathbf{x})$. Note that the Hessian of $f(\mathbf{x})$ at \mathbf{x} computed with respect to $A = \nabla^2 f(\mathbf{x})$ is I_n .

Other Approximation derived Methods: Variable Metric

Example 9

Consider quadratic function

$$f(\mathbf{x}) = \alpha + \langle \mathbf{a}, \mathbf{x} \rangle + \frac{1}{2} \langle A\mathbf{x}, \mathbf{x} \rangle,$$

where $A = A^\top \succeq 0$. Note that $\nabla f(\mathbf{x}) = A\mathbf{x} + \mathbf{a}$, $\nabla^2 f(\mathbf{x}) = A$ and

$$\nabla f(\mathbf{x}^*) = A\mathbf{x}^* + \mathbf{a} = 0$$

for $\mathbf{x}^* = -A^{-1}\mathbf{a}$.

Other Approximation derived Methods: Variable Metric

Example 9 (Continued.) Let us compute the Newton direction at some $\mathbf{x} \in \mathbb{R}^n$:

$$d_N(\mathbf{x}) = [\nabla^2 f(\mathbf{x})]^{-1} \nabla f(\mathbf{x}) = A^{-1}(A\mathbf{x} + \mathbf{a}) = \mathbf{x} + A^{-1}\mathbf{a}.$$

Therefore for any $\mathbf{x} \in \mathbb{R}^n$ we have $\mathbf{x} - d_N(\mathbf{x}) = -A^{-1}\mathbf{a} = \mathbf{x}^*$. Thus, for a quadratic function the Newton method converges in **one step**. Note also that

$$f(\mathbf{x}) = \alpha + \langle A^{-1}\mathbf{a}, \mathbf{x} \rangle_A + \frac{1}{2} \|\mathbf{x}\|_A^2,$$

$$\nabla f_A(\mathbf{x}) = A^{-1} \nabla f(\mathbf{x}) = d_N(\mathbf{x}),$$

$$\nabla^2 f_A(\mathbf{x}) = A^{-1} \nabla^2 f(\mathbf{x}) = I_n.$$

更一般的情况下, $A^{-1} = H$, A 可以不是 $\nabla^2 f$ 。

Other Approximation derived Methods: Variable Metric

Let us write down a general scheme of the variable metric methods.

Variable metric method
<p>0. Choose $\mathbf{x}_0 \in \mathbb{R}^n$. Set $H_0 = I_0$. Compute $f(\mathbf{x}_0)$ and $\nabla f(\mathbf{x}_0)$.</p>
<p>1. kth iteration ($k \geq 0$).</p> <ul style="list-style-type: none">a Set $p_k = H_k \nabla f(\mathbf{x}_k)$.b Find $\mathbf{x}_{k+1} = \mathbf{x}_k - h_k p_k$.c Compute $f(\mathbf{x}_{k+1})$ and $\nabla f(\mathbf{x}_{k+1})$.d Update the matrix $H_k : H_k \rightarrow H_{k+1}$.

Other Approximation derived Methods: Variable Metric

Quasi-Newton rule: The variable metric schemes differ one from another only in implementation of Step 1d), which updates matrix H_k . For that, they use new information, accumulated at Step 1c), namely the gradient $\nabla f(\mathbf{x}_{k+1})$. The idea is justified by the following property of a quadratic function. Let

$$f(\mathbf{x}) = \alpha + \langle \alpha, \mathbf{x} \rangle + \frac{1}{2} \langle A\mathbf{x}, \mathbf{x} \rangle, \quad \nabla f(\mathbf{x}) = A\mathbf{x} + \alpha.$$

Then, for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ we have $\nabla f(\mathbf{x}) - \nabla f(\mathbf{y}) = A(\mathbf{x} - \mathbf{y})$. This identity explains the origin of the so-called **Quasi-Newton rule**:

$$H_{k+1} (\nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)) = \mathbf{x}_{k+1} - \mathbf{x}_k.$$

Other Approximation derived Methods: Variable Metric

Example 10

Denote

$$\Delta H_k = H_{k+1} - H_k, \quad \gamma_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k), \quad \delta_k = \mathbf{x}_{k+1} - \mathbf{x}_k.$$

Then the quasi-Newton relation is satisfied by the following rules.

1 Rank-one correcton scheme

$$\Delta H_k = \frac{(\delta_k - H_k \gamma_k)(\delta_k - H_k \gamma_k)^\top}{\langle \delta_k - H_k \gamma_k, \gamma_k \rangle}.$$

2 Davidon-Fletcher-Powell scheme (DFP))

$$\Delta H_k = \frac{\delta_k \delta_k^\top}{\langle \gamma_k, \delta_k \rangle} - \frac{H_k \gamma_k \gamma_k^\top H_k}{\langle H_k \gamma_k, \gamma_k \rangle}.$$

Other Approximation derived Methods: Variable Metric

Example 10

(Continued.)

3 Broyden-Fletcher-Goldfarb-Shanno scheme (BFGS)

$$\Delta H_k = \frac{H_k \gamma_k \delta_k^\top + \delta_k \gamma_k^\top H_k}{\langle H_k \gamma_k, \gamma_k \rangle} - \beta_k \frac{H_k \gamma_k \gamma_k^\top H_k}{\langle H_k \gamma_k, \gamma_k \rangle},$$

where $\beta_k = 1 + \langle \gamma_k, \delta_k \rangle / \langle H_k \gamma_k, \gamma_k \rangle$.

Clearly, there are many other possibilities. From the computational point of view, BFGS is considered as the most stable scheme.

Remarks on Variable Metric

- 1 Note that for quadratic functions the variable metric methods usually terminate in n iterations.
- 2 In a neighborhood of strict minimum they have a **superlinear** rate of convergence: for any $\mathbf{x}_0 \in \mathbb{R}^n$ there exists a number N such that for all $k \geq N$ we have

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \text{const} \cdot \|\mathbf{x}_k - \mathbf{x}^*\| \cdot \|\mathbf{x}_{k-n} - \mathbf{x}^*\|$$

(the proofs are very long and technical).

- 3 As far as global convergence is concerned, these methods are not better than the gradient method (at least, from the theoretical point of view).

Remarks on Variable Metric

- 4 Note that in the variable metric schemes it is necessary to store and update a symmetric $n \times n$ -matrix. Thus, each iteration needs $O(n^2)$ auxiliary arithmetic operations. During many years this feature was considered as one of **the main drawbacks** of the variable metric methods. That stimulated the interest in so-called **conjugate gradients** schemes, which have much lower complexity of each iteration.

Part II

Conjugate Gradient

Methods Based on Krylov Subspace

In the previous sections, all \mathbf{x}_k are spanned from $\text{Lin}\{\mathbf{x}_0, \nabla f(\mathbf{x}_i), i = 1, \dots, k-1\}$. Essentially, pick a point from Lin . Is it possible to consider finding solutions in a subspace?

Krylov Subspace: To solve a large-scale sparse system of linear equations

$$A\mathbf{x} = \mathbf{b}, \quad A \in \mathbb{R}^{n \times n}, \mathbf{b} \in \mathbb{R}^n,$$

Krylov subspace can be used as a target subspace in many methods ([Arnoldi](#), [Lanczos](#), [Conjugate gradient](#), [GMRES](#)).

Remark. This type of method is also regarded as a projection method, that is, to find the projection of the true solution in a certain subspace (which can be an orthogonal projection or an oblique projection).

Historical Origin of CG

The **conjugate gradient** method were initially proposed for the numerical solution of particular systems of linear equations (Stiefel [1952]). It also can be used to minimize quadratic functions like

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad (33)$$

where $f(\mathbf{x}) = \alpha + \langle \mathbf{a}, \mathbf{x} \rangle + \frac{1}{2} \langle A\mathbf{x}, \mathbf{x} \rangle$, and $A = A^\top \succ 0$, since

$$A\mathbf{x} = -\mathbf{a}.$$

We have already seen that the solution to this problem is $\boxed{\mathbf{x}^* = -A^{-1}\mathbf{a}}$.

Historical Origin of CG

Therefore, our objective function can be written in the following form

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

Thus,

$$f^* = \alpha - \frac{1}{2} \langle A\mathbf{x}^*, \mathbf{x}^* \rangle \text{ and}$$

$$\boxed{\nabla f(\mathbf{x}) = A(\mathbf{x} - \mathbf{x}^*)}.$$

Historical Origin of CG

Suppose we are given a starting point $\mathbf{x}_0 \in \mathbb{R}^n$. Consider the following linear **Krylov subspace**

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

where A^k is the k -th power of matrix A . A sequence of points $\{\mathbf{x}_k\}$ is generated by the **Conjugate Gradient** Method in accordance with the following rule.

$$\boxed{\mathbf{x}_k = \operatorname{argmin}\{f(\mathbf{x}) \mid \mathbf{x} \in \mathbf{x}_0 + \mathcal{L}_k\}, \quad k \geq 1.} \quad (34)$$

This definition looks quite artificial. However, later we will see that this method can be written in a pure “algorithmic” form. We need representation (34) only for theoretical analysis.

Historical Origin of CG

Calculating the optimum on $\mathbf{x}_0 + \mathcal{L}_k$ means that in the k -th iteration, k -th weights are to be calculated. This is not what we want. In practice, we hope to have the following calculation process.

- This is a descent algorithm. That is, for $k = 1, \dots$,

$$\begin{aligned}\mathbf{x}_k &= \mathbf{x}_{k-1} + \alpha_{k-1} \mathbf{p}_{k-1}, \Rightarrow (A\mathbf{x}_k + \mathbf{a}) = (A\mathbf{x}_{k-1} + \mathbf{a}) + \alpha_{k-1} A\mathbf{p}_{k-1}, \\ &\Rightarrow \boxed{\nabla f(\mathbf{x}_k) = \nabla f(\mathbf{x}_{k-1}) + \alpha_{k-1} A\mathbf{p}_{k-1}},\end{aligned}\quad (35)$$

where $\alpha_{k-1} > 0$, for $k = 1, \dots$

- The direction is updated according to the gradient and last direction as follows. With $\beta_k > 0$, $\mathbf{p}_0 = \nabla f(\mathbf{x}_0)$, and for $k = 1, \dots$,

$$\boxed{\mathbf{p}_k = \nabla f(\mathbf{x}_k) - \beta_{k-1} \mathbf{p}_{k-1}}. \quad (36)$$

Fundamental Theory

Lemma 11

In view of (35) and (36), for any $k \geq 1$,

$$\begin{aligned} \{\mathcal{L}_k &\triangleq \text{Lin}\{A(\mathbf{x}_0 - \mathbf{x}^*), A^2(\mathbf{x}_0 - \mathbf{x}^*), \dots, A^k(\mathbf{x}_0 - \mathbf{x}^*)\}\} \\ &\equiv \{\mathcal{R}_k \triangleq \text{Lin}\{\nabla f(\mathbf{x}_0), \nabla f(\mathbf{x}_1), \dots, \nabla f(\mathbf{x}_{k-1})\}\} \end{aligned} \quad (37)$$

and

$$\begin{aligned} \{\mathcal{L}_k &\triangleq \text{Lin}\{A(\mathbf{x}_0 - \mathbf{x}^*), A^2(\mathbf{x}_0 - \mathbf{x}^*), \dots, A^k(\mathbf{x}_0 - \mathbf{x}^*)\}\} \\ &\equiv \{\mathcal{P}_k \triangleq \text{Lin}\{\mathbf{p}_0, \mathbf{p}_1, \dots, \mathbf{p}_{k-1}\}\}. \end{aligned} \quad (38)$$

Fundamental Theory

Proof. The proof is by induction. (37) and (38) hold trivially for $k = 1$. since

$$\mathcal{L}_1 = \mathcal{R}_1 = \mathcal{P}_1, \text{ since } \nabla f(\mathbf{x}_0) = A(\mathbf{x}_0 - \mathbf{x}^*) = \mathbf{p}_0.$$

We here suppose that

$$\mathcal{L}_k = \mathcal{R}_k \text{ and } \mathcal{L}_k = \mathcal{P}_k. \quad (39)$$

Let's check the $\nabla f(\mathbf{x}_k)$ in \mathcal{R}_{k+1} and $A^{k+1}(\mathbf{x}_0 - \mathbf{x}^*)$ in \mathcal{L}_{k+1} .

■ For $\nabla f(\mathbf{x}_k)$: In view of (35), we have

$$\begin{aligned} \nabla f(\mathbf{x}_k) &= \nabla f(\mathbf{x}_{k-1}) + \alpha_{k-1} A \mathbf{p}_{k-1} \\ &\in \mathcal{R}_k \cup A \mathcal{P}_k \\ &= \mathcal{L}_k \cup A \mathcal{L}_k = \mathcal{L}_{k+1}. \end{aligned}$$

Fundamental Theory

Proof. (Continued.)

- For $A^{k+1}(\mathbf{x}_0 - \mathbf{x}^*)$: In view of (39), we have

$$\begin{aligned} A^{k+1}(\mathbf{x}_0 - \mathbf{x}^*) &= A(A^k(\mathbf{x}_0 - \mathbf{x}^*)) \in A\mathcal{L}_k = A\mathcal{P}_k \\ &= \text{Lin}\{A\mathbf{p}_0, A\mathbf{p}_1, \dots, A\mathbf{p}_{k-1}\} = \text{Lin}\{\nabla f(\mathbf{x}_0), \nabla f(\mathbf{x}_1), \dots, \nabla f(\mathbf{x}_k)\}. \end{aligned}$$

The last equality comes from the fact that in view of (35), we have

$$A\mathbf{p}_0 = \frac{1}{\alpha_0}(\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_0)), \dots, A\mathbf{p}_{k-1} = \frac{1}{\alpha_{k-1}}(\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k-1})).$$

Thus, $A^{k+1}(\mathbf{x}_0 - \mathbf{x}^*) \in \mathcal{R}_{k+1}$. Thus, we arrive at $\mathcal{R}_{k+1} = \mathcal{L}_{k+1}$.



Fundamental Theory

CG: a two-term recurrence

$$\begin{aligned}\mathbf{x}_k &= \mathbf{x}_{k-1} + \alpha_{k-1} \mathbf{p}_{k-1}, \\ \mathbf{p}_k &= \nabla f(\mathbf{x}_k) - \beta_{k-1} \mathbf{p}_{k-1}.\end{aligned}$$

Thus, we have

$$\begin{aligned}\mathbf{x}_{k+1} &= \mathbf{x}_k + \alpha_k \mathbf{p}_k \\ &= \mathbf{x}_k + \alpha_k (\nabla f(\mathbf{x}_k) - \beta_{k-1} \mathbf{p}_{k-1}) \\ &= \mathbf{x}_k + \alpha_k \nabla f(\mathbf{x}_k) - \underbrace{\alpha_k \beta_{k-1} \mathbf{p}_{k-1}}_{\in \mathcal{L}_k}.\end{aligned}\tag{40}$$

Fundamental Theory

Lemma 12 (Lemma 1.3.2)

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

Proof. Let $k > i$. Consider the function

$$\phi(\lambda) = f \left(\mathbf{x}_0 + \sum_{j=1}^k \lambda^{(j)} \nabla f(\mathbf{x}_{j-1}) \right), \quad \lambda \in \mathbb{R}^k.$$

In view of Lemma 11, for some λ_* we have $\mathbf{x}_k = \mathbf{x}_0 + \sum_{j=1}^k \lambda_*^{(j)} \nabla f(\mathbf{x}_{j-1})$.

Fundamental Theory

Lemma 12

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

Proof. (Continued)

However, by definition, \mathbf{x}_k is the minimum point of $f(\mathbf{x})$ on $\mathbf{x}_0 + \mathcal{L}_k$. Therefore $\nabla \phi(\lambda_*) = 0$.

It remains to compute the components of the gradient:

$$0 = \frac{\partial \phi(\lambda_*)}{\partial \lambda^{(i)}} = \left\langle \frac{\partial \phi(\lambda_*)}{\partial \mathbf{x}_k}, \frac{\partial \mathbf{x}_k}{\partial \lambda^{(i)}} \right\rangle = \langle \nabla f(\mathbf{x}_k), \nabla f(\mathbf{x}_i) \rangle.$$



Fundamental Theory

Corollary 13

The sequence generated by the Conjugate Gradient Method for problem (33) is finite.

证明.

Indeed, the number of nonzero orthogonal directions in \mathbb{R}^n cannot exceed n . □

Corollary 14

For any $p \in \mathcal{L}_k$, $k \geq 1$, we have $\langle \nabla f(x_k), p \rangle = 0$.

$\mathcal{L}_k = \mathcal{P}_k$, $\nabla f(x_k) \in \mathcal{L}_{k+1}$ and $\nabla f(x_k) \notin \mathcal{L}_k$.

Fundamental Theory

The last auxiliary result explains the name of the method. Let $\delta_i = \mathbf{x}_{i+1} - \mathbf{x}_i$. It is clear that $\mathcal{L}_k = \text{Lin}\{\delta_0, \dots, \delta_{k-1}\}$. (since we have $\delta_i = c\mathbf{p}_i$)

Lemma 15

For any $k \neq i$, we have $\langle A\delta_k, \delta_i \rangle = 0$.

(Such directions are called conjugate with respect to A)

Proof. Without loss of generality, we can assume that $k > i$. Then (Note that $\nabla f(\mathbf{x}) = A(\mathbf{x} - \mathbf{x}^*)$)

$$\langle A\delta_k, \delta_i \rangle = \langle A(\mathbf{x}_{k+1} - \mathbf{x}_k), \delta_i \rangle = \underbrace{\langle \nabla f(\mathbf{x}_{k+1}), \delta_i \rangle}_{k+2 \text{ Kry}} - \underbrace{\langle \nabla f(\mathbf{x}_k), \delta_i \rangle}_{k+1 \text{ Kry}} = 0,$$

since $\delta_i = \mathbf{x}_{i+1} - \mathbf{x}_i \in \mathcal{L}_{i+1} \subseteq \mathcal{L}_k$. □

CG Algorithm

Let us show how we can write down the **Conjugate Gradient** Method in a more algorithmic form. Since $\mathcal{L}_{k+1} = \text{Lin}\{\delta_0, \dots, \delta_k\}$, according to (40), we can represent \mathbf{x}_{k+1} as follows:

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

In our notation, this is

$$\delta_k = -h_k \nabla f(\mathbf{x}_k) + \sum_{j=0}^{k-1} \lambda^{(j)} \delta_j. \quad (41)$$

$$\begin{aligned} \mathcal{L}_{k+1} &= \text{Lin}\{\delta_0, \dots, \delta_k\} = \text{Lin}\{\nabla f(\mathbf{x}_0), \dots, \nabla f(\mathbf{x}_k)\} \\ &= \text{Lin}\{\delta_0, \dots, \delta_{k-1}, \nabla f(\mathbf{x}_k)\} \end{aligned}$$

CG Algorithm

Let us compute the coefficients in this representation. Multiplying (41) by A and δ_i , $0 \leq i \leq k-1$, and using lemma 15, we obtain

$$\begin{aligned}
 0 &= \langle A\delta_k, \delta_i \rangle = -h_k \langle A\nabla f(\mathbf{x}_k), \delta_i \rangle + \sum_{j=0}^{k-1} \lambda^{(j)} \langle A\delta_j, \delta_i \rangle \\
 &= -h_k \langle A\nabla f(\mathbf{x}_k), \delta_i \rangle + \lambda^{(i)} \langle A\delta_i, \delta_i \rangle \\
 &= \underbrace{-h_k \langle \nabla f(\mathbf{x}_k), \nabla f(\mathbf{x}_{i+1}) - \nabla f(\mathbf{x}_i) \rangle}_{i < k-1 \Rightarrow \text{above}=0} + \underbrace{\lambda^{(i)} \langle A\delta_i, \delta_i \rangle}_{i < k-1 \Rightarrow \lambda^{(i)}=0 \text{ since } \langle A\delta_i, \delta_i \rangle \text{ is positive.}}.
 \end{aligned}$$

CG Algorithm

Hence, in view of Lemma 12, $\lambda^{(i)} = 0$, for $i < k - 1$. For $\boxed{i = k - 1}$, we have

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

Thus, $\mathbf{x}_{k+1} = \mathbf{x}_k - h_k (\nabla f(\mathbf{x}_k) - \lambda^{(k-1)}\delta_{k-1}) = \mathbf{x}_k - h_k \mathbf{p}_k$, where

$$\begin{aligned} \mathbf{p}_k &= \nabla f(\mathbf{x}_k) - \frac{\|\nabla f(\mathbf{x}_k)\|^2 \delta_{k-1}}{\langle \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k-1}), \delta_{k-1} \rangle} \\ &= \nabla f(\mathbf{x}_k) - \frac{\|\nabla f(\mathbf{x}_k)\|^2 \mathbf{p}_{k-1}}{\langle \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_{k-1}), \mathbf{p}_{k-1} \rangle}, \end{aligned}$$

since $\delta_{k-1} = \mathbf{x}_k - \mathbf{x}_{k-1} = -h_{k-1}\mathbf{p}_{k-1}$ by definition of the directions $\{\mathbf{p}_k\}$.

CG Algorithm

Conjugate Gradient Method

0. Let $\mathbf{x}_0 \in \mathbb{R}^n$. Compute $f(\mathbf{x}_0)$, $\nabla f(\mathbf{x}_0)$. 设 $\mathbf{p}_0 = \nabla f(\mathbf{x}_0)$.

1. k -th iteration ($k \geq 0$)

- a Find $\mathbf{x}_{k+1} = \mathbf{x}_k + h_k \mathbf{p}_k$ (by "exact" line search).
- b Compute $f(\mathbf{x}_{k+1})$ and $\nabla f(\mathbf{x}_{k+1})$.
- c Compute the coefficient β_k .
- d Define $\mathbf{p}_{k+1} = \nabla f(\mathbf{x}_{k+1}) - \beta_k \mathbf{p}_k$.

CG Algorithm

The specification of the coefficient β_k .

1 Dai-Yuan

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

2 Fletcher-Rieves:

$$\beta_k = -\frac{\|\nabla f(\mathbf{x}_{(k+1)})\|^2}{\|\nabla f(\mathbf{x}_{(k)})\|^2}.$$

3 Polak-Ribbiere:

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

CG Algorithm

- Recall that in the quadratic case, the Conjugate Gradient Method terminates in n iterations (or less). Algorithmically, this means that $\mathbf{p}_n = 0$.
- In the general nonlinear case, this is not true. However, after n iterations, this direction loses its interpretation. Therefore, in all practical schemes, there exists a restarting strategy. This ensures the global convergence of the process.

In a neighborhood of a strict minimum, the conjugate gradient schemes demonstrate a local n -step quadratic convergence:

$$\|\mathbf{x}_{n+1} - \mathbf{x}^*\| \leq \text{const} \cdot \|\mathbf{x}_0 - \mathbf{x}^*\|^2$$

Note that this local convergence is slower than that of the variable metric methods. However, the conjugate gradient methods have the advantage of cheap iteration. As far as the global convergence is concerned, these schemes, in general, are not better than the simplest Gradient Method.

References I

Yurii Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*, volume 87. Springer Science & Business Media, 2003.

Eduard Stiefel. Methods of conjugate gradients for solving linear systems. *J. Res. Nat. Bur. Standards*, 49:409–435, 1952.

Lloyd N Trefethen and David Bau III. *Numerical linear algebra*, volume 50. Siam, 1997.

Stephen Wright, Jorge Nocedal, et al. Numerical optimization. *Springer Science*, 35(67-68): 7, 1999.

Stephen J Wright and Benjamin Recht. *Optimization for data analysis*. Cambridge University Press, 2022.

Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.

Yurii Nesterov. *Lectures on convex optimization*, volume 137. Springer, 2018.

Thank You!

Email: qianhui@zju.edu.cn