

Lecture-6

Convergence and order of
convergence

Line Search Methods

$$x_{k+1} \leftarrow x_k + \alpha_k p_k$$

$$p_k \leftarrow -B_k^{-1} \nabla f_k$$

Steepest descent B_k is and identity matrix

Newton B_k is a Hessian matrix

Quasi-Newton B_k is approximation to the Hessian matrix

Line Search Methods

$$x_{k+1} \leftarrow x_k + \alpha_k p_k$$

$$p_i^T A p_j = 0 \quad \forall i \neq j$$

Conjugate gradient

Important Questions

- What are the conditions under which, the method converges?
- What is the rate of convergence?

Conditions of convergence

- Steepest Descent: Wolfe's conditions
- Newton and Quasi-Newton: In addition to Wolfe's conditions, PD Hessian, and bounded condition number
- Conjugate Gradient: subsequence of direction cosines $\cos \mathbf{q}_k$ is bounded away from zero.

Convergence Rate

- Steepest descent: Linear
- Quasi-Newton: Super-linear
- Newton: Quadratic
- Conjugate Gradient: n steps

Convergence of Line Search Methods

- The steepest descent method is globally convergent
- For other algorithms how far p_k can deviate from the steepest descent direction and still gives rise to globally convergent iteration.

Convergence of Line Search Methods (Theorem 3.2)

The angle between p_k and steepest descent direction $-\nabla f_k^T$

$$\cos \mathbf{q}_k = \frac{-\nabla f_k^T p_k}{\|\nabla f_k\| \|p_k\|}$$

We will show (Theorem 3.2):

$$\sum_{k \geq 0} \cos^2 \mathbf{q}_k \|\nabla f_k\|^2 < \infty$$

Convergence of Line Search Methods

$$x_{k+1} = x_k + \mathbf{a}_k p_k \quad \text{Iteration scheme}$$

$$\nabla f(x_k + \mathbf{a}_k p_k)^T p_k \geq c_2 \nabla f_k^T(x_k) p_k, \quad c_2 \in (c_1, 1) \quad \text{Curvature condition}$$

Therefore

$$(\nabla f_{k+1} - \nabla f_k)^T p_k \geq (c_2 - 1) \nabla f_k^T p_k \quad (1)$$

$$\| \nabla f(x) - \nabla f(\tilde{x}) \| \leq L \| x - \tilde{x} \| \quad \text{Lipschitz continuous}$$

$$\begin{aligned} (\nabla f_{k+1} - \nabla f_k)^T p_k &\leq \| (\nabla f_{k+1} - \nabla f_k)^T \| \| p_k \| \\ &\leq \mathbf{a}_k L \| p_k \| \| p_k \| \end{aligned}$$

$$(\nabla f_{k+1} - \nabla f_k)^T p_k \leq \mathbf{a}_k L \| p_k \|^2 \quad (2)$$

Convergence of Line Search Methods

$$(\nabla f_{k+1} - \nabla f_k)^T p_k \leq \mathbf{a}_k L \| p_k \|^2 \quad (2)$$

$$\frac{(\nabla f_{k+1} - \nabla f_k)^T p_k}{L \| p_k \|^2} \leq \mathbf{a}_k$$

$$\mathbf{a}_k \geq \frac{(\nabla f_{k+1} - \nabla f_k)^T p_k}{L \| p_k \|^2}$$

$$(\nabla f_{k+1} - \nabla f_k)^T p_k \geq (c_2 - 1) \nabla f_k^T p_k \quad (1)$$

Combining (1) and (2)

$$\mathbf{a}_k \geq \frac{c_2 - 1}{L} \frac{\nabla f_k^T p_k}{\| p_k \|^2}$$

Convergence of Line Search Methods

$$\mathbf{a}_k \geq \frac{c_2 - 1}{L} \frac{\nabla f_k^T \mathbf{p}_k}{\|\mathbf{p}_k\|^2}$$

$$f(x_k + \mathbf{a}_k \mathbf{p}_k) \leq f(x_k) + c_1 \mathbf{a}_k \nabla f_k^T \mathbf{p}_k \quad \text{Sufficient decrease}$$

Therefore

$$f_{k+1} \leq f_k - c_1 \frac{1 - c_2}{L} \frac{(\nabla f_k^T \mathbf{p}_k)^2}{\|\mathbf{p}_k\|^2}$$

$$f_{k+1} \leq f_k - c \cos^2 \mathbf{q}_k \|\nabla f_k\|^2, \quad c = c_1(1 - c_2)/L$$

$$f_{k+1} \leq f_0 - c \sum_{j=0}^k \cos^2 \mathbf{q}_j \|\nabla f_j\|^2$$

Convergence of Line Search Methods

$$f_{k+1} \leq f_0 - c \sum_{j=0}^k \cos^2 \mathbf{q}_j \|\nabla f_j\|^2$$

Since f is bounded below, we have $f_0 - f_{k+1}$ is less than some positive constant for all k

Taking the limits:

$$\sum_{k \geq 0} \cos^2 \mathbf{q}_k \|\nabla f_k\|^2 < \infty$$

Convergence of Line Search Methods

$$\sum_{k \geq 0} \cos^2 \mathbf{q}_k \|\nabla f_k\|^2 < \infty \quad \cos^2 \mathbf{q}_k \|\nabla f_k\|^2 \rightarrow 0$$

$$\cos \mathbf{q}_k \geq \mathbf{d} > 0 \quad \text{If angle is bounded away From } 90^\circ$$

$$\lim_{k \rightarrow \infty} \|\nabla f_k\|^2 = 0$$

We can be sure that gradient norms converges to zero, provided that the search directions are never too close to orthogonality with the gradient

Therefore, the steepest descent produces a gradient sequence that converges to zero, provided that it uses a line search satisfying Wolf's conditions.

We can not guarantee that the method converges to a minimizer, but only that it is attracted by stationary points.

Newton-Like

$$x_{k+1} = x_k + \mathbf{a}_k p_k \quad p_k = -B_k^{-1} \nabla f_k$$

Assume Hessian is a PD with a uniformly bounded condition number

$$\|B_k\| \|B_k^{-1}\| \leq M, \forall k$$

Using
$$\cos \mathbf{q}_k = \frac{-\nabla f_k^T p_k}{\|\nabla f_k\| \|p_k\|}$$

Show that (Homework)

$$\cos \mathbf{q}_k \geq \frac{1}{M}$$

Newton-Like

$$\cos \mathbf{q}_k \geq \frac{1}{M}$$

Using $\sum_{k \geq 0} \cos^2 \mathbf{q}_k \|\nabla f_k\|^2 < \infty$ Theorem 3.2

$$\cos^2 \mathbf{q}_k \|\nabla f_k\|^2 \rightarrow 0$$

Therefore

$$\cos^2 \mathbf{q}_k \geq 0$$

$$\lim_{k \rightarrow \infty} \|\nabla f_k\|^2 = 0$$

Therefore:

We have shown that :
 Newton and Quasi Newton
 are globally convergent
 if Hessians have bounded condition
 numbers and are PD, and if the step
 lengths satisfy Wolf's conditions

Conjugate Gradient

Only subsequence of the gradient norms converges to zero,
 rather than the whole sequence.

$$\liminf_{k \rightarrow \infty} \|\nabla f_k\|^2 = 0$$

Sketch of proof by contradiction:

$$\|\nabla f_k\| \geq \mathbf{g}$$

Then $\cos^2 \mathbf{q}_k \|\nabla f_k\|^2 \rightarrow 0$

Implies $\cos \mathbf{q}_k \rightarrow 0$

Therefore it is enough to show that a subsequence $\{\cos \mathbf{q}_k\}$ is
 bounded away from zero.

General Class of Algorithms

- Algorithm
 - Every iteration produces a decrease in the objective function
 - Every m *the*-th iteration is a steepest descent step, with the step length chosen to satisfy the Wolf's conditions.
- Then
 - Since $\cos \alpha_k = 1$ for steepest descent, then following holds

$$\lim_{k \rightarrow \infty} \|\nabla f_k\|^2 = 0$$

Convergence Rate of Steepest Descent: Quadratic Function

$$\|x_{k+1} - x^*\|_Q^2 \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right)^2 \|x_k - x^*\|_Q^2 \quad \text{Theorem 3.3}$$

As the condition number increases the contours of the quadratic become more elongated, the zigzags of line search becomes more pronounced.

Theorem 3.4: Steepest Descent

$$f(x_{k+1}) - f(x^*) \leq \left(\frac{I_n - I_1}{I_n + I_1} \right)^2 (f(x_k) - f(x^*))$$

where $0 \leq I_1 \leq I_2 \leq \dots \leq I_n$ are eigenvalues of Hessian

If the condition number is 800, and $f(x_1) = 1$ and $f(x^*) = 0$,
After 1000 iterations the value of function will be .08.

Theorem 3.5

Suppose f is three times continuously differentiable. Consider iteration $x_{k+1} = x_k + \alpha_k p_k$, where p_k is a descent direction, α_k satisfies Wolfe's conditions, with $\alpha_k > 0$. If the sequence converges to a point x^* such that $\nabla f(x^*) = 0$ and B_k is pd, and if the search direction satisfies

$$\lim_{k \rightarrow \infty} \frac{\| \nabla f_k + \nabla^2 f_k p_k \|}{\| p_k \|} = 0$$

$$\lim_{k \rightarrow \infty} \frac{\| (B_k - \nabla^2 f(x^*)) p_k \|}{\| p_k \|} = 0$$

Then

- (i) if p_k is admissible for all $k > k_0$
- (ii) If $\| p_k \|$ is bounded away from zero for all $k > k_0$, then $\{x_k\}$ converges to x^* superlinearly.

Theorem 3.6

Suppose f is three times continuously differentiable. Consider iteration $x_{k+1} = x_k + p_k$, where p_k is given by Quasi-Newton direction. Assume the sequence x_k converges to a point x^* such that B_k and $\nabla^2 f(x^*)$ is pd, the sequence converges superlinearly iff the following condition holds.

$$\lim_{k \rightarrow \infty} \frac{\| (B_k - \nabla^2 f(x^*)) p_k \|}{\| p_k \|} = 0$$

Order Notations

Given two non-negative infinite sequences

$$\mathbf{h}_k = O(\mathbf{n}_k)$$

$$\text{if } |\mathbf{h}_k| \leq C |\mathbf{n}_k|, \text{ for } C > 0, \forall k$$

$$\mathbf{h}_k = o(\mathbf{n}_k)$$

$$\text{if } \lim_{k \rightarrow \infty} \frac{\mathbf{h}_k}{\mathbf{n}_k} = 0$$

Sketch of a Proof

$$\begin{aligned}
 p_k - p_k^N &= \nabla^2 f_k^{-1} (\nabla^2 f_k p_k + \nabla f_k) \\
 &= \nabla^2 f_k^{-1} (\nabla^2 f_k - B_k) p_k \\
 &= O(\| (\nabla^2 f_k - B_k) p_k \|) \\
 &= o(\| p_k \|)
 \end{aligned}$$

$$h_k = O(n_k)$$

$$\text{if } \|h_k\| \leq C \|n_k\|, \text{ for } C > 0$$

$$h_k = o(n_k)$$

$$\text{if } \lim_{k \rightarrow \infty} \frac{h_k}{n_k} = 0$$

$$\lim_{k \rightarrow 0} \frac{\| (B_k - \nabla^2 f(x^*)) p_k \|}{\| p_k \|} = 0$$

$$\lim_{k \rightarrow 0} \frac{\| \nabla f_k + \nabla^2 f_k p_k \|}{\| p_k \|} = 0$$

$$\begin{aligned}
 \| x_k + p_k - x^* \| &\leq \| x_k + p_k^N - x^* \| + \| p_k - p_k^N \| \\
 &= O(\| x_k - x^* \|^2) + o(\| p_k \|)
 \end{aligned}$$

$$\| x_k + p_k - x^* \| \leq o(\| x_k - x^* \|)$$

Theorem 3.7

Suppose that f is twice differentiable and that Hessian is Lipschitz continuous. Consider the iteration $p_k^N = -\nabla^2 f_k^{-1} \nabla f_k$ where is given by

$$p_k^N = -\nabla^2 f_k^{-1} \nabla f_k$$

Then:

1. If the starting point x_0 is sufficiently close to x^* , the sequence converges to x^* .
2. The rate of convergence is quadratic
3. The sequence of gradient norms $\| \nabla f_k \|$ converges quadratically to zero.

Coordinate Descent Method

Cycle through n coordinate directions e_1, e_2, \dots, e_n using each in turn as a search direction.

Fix all other variables except one, and minimize the function.

It is an inefficient method, it can iterate infinitely without ever approaching a point, where the gradient vanishes.

The gradient may become more and more perpendicular to search Directions, making $\cos \theta$ approach to zero, but not the gradient.