

Large Scale Optimization: Lecture 7

Jimmy Lin, Vutha Va and David Inouye

The University of Texas at Austin

September 23, 2014

Newton Step

Definition

For $x \in \text{dom } f$, the vector

$$\Delta x_{nt} = -\nabla^2 f(x)^{-1} \nabla f(x) \quad (1)$$

is called *Newton step* for function $f(\cdot)$ at point x .

- Note that $f(\cdot)$ is twice differentiable.
- Because $\nabla^2 f(x)$ is positive definite,

$$\nabla f(x)^T \Delta x_{nt} = -\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) < 0 \quad (2)$$

unless $\nabla f(x) = 0$. So the Newton step is a descent direction unless x is already optimal.

Interpretation I

Minimizer of second-order Approximation

Second-order Taylor approximation \hat{f} of f at x is

$$\hat{f}(x+v) = f(x) + \nabla f(x)^T v + \frac{1}{2} v^T \nabla^2 f(x) v \quad (3)$$

where RHS is minimized at direction

$$v = -\nabla^2 f(x)^{-1} \nabla f(x) = \Delta x_{nt} \quad (4)$$

Combined with update rule, we have Newton update

$$x^+ = x - t \nabla^2 f(x)^{-1} \nabla f(x) \quad (5)$$

where t is fixed step size.

Interpretation I

Minimizer of second-order Approximation

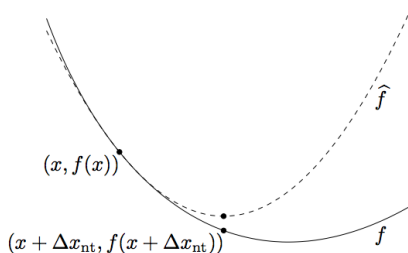


Figure 1 : The function f (shown solid) and its second-order approximation \hat{f} at x (dashed). The Newton step Δx_{nt} is what must be added to x to give the minimizer of \hat{f} . [Figure and caption from *Convex Optimization* by Boyd and Vandenberghe.]

Interpretation II

Steepest descent direction in Hessian norm

Newton step can also be interpreted as steepest descent direction when the norm is defined as

$$\|u\|_{\nabla^2 f(x)} \triangleq \sqrt{u^T \nabla^2 f(x) u} \quad (6)$$

Interpretation II

Steepest descent direction in Hessian norm

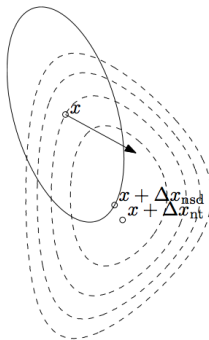


Figure 2 : The dashed lines are level curves of a convex function. The ellipsoid shown (with solid line) is $\{x + v \mid v^T \nabla^2 f(x) v \leq 1\}$. The arrow shows $-\nabla f(x)$, the gradient descent direction. The Newton step Δx_{nt} is the steepest descent direction in the norm $\|\cdot\|_{\nabla^2 f(x)}$. [Figure and caption from *Convex Optimization* by Boyd and Vandenberghe.]

Interpretation III

Solution of linearized optimality condition

Newton step can also be interpreted as linear approximation over gradient $\nabla f(x)$ around x .

$$\nabla f(x + v) \approx \nabla f(x) + \nabla^2 f(x)v \quad (7)$$

Set RHS to zero gives Newton step Δx_{nt}

$$v = -\nabla^2 f(x)^{-1} \nabla f(x) = \Delta x_{nt} \quad (8)$$

So the Newton step Δx_{nt} is what must be added to x so that the linearized optimality condition holds.

Again, this suggests that when x is near x^* , the update $x + \Delta x_{nt}$ should be a very good approximation of x^* .

Interpretation III

Solution of linearized optimality condition

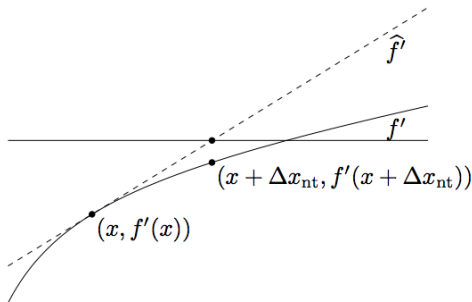


Figure 3 : The solid curve is the derivative f' of the function f shown in Figure 2. \hat{f}' is the linear approximation of f' at x . The Newton step Δx_{nt} is the difference between the root of \hat{f}' and the point x . [Figure and caption from *Convex Optimization* by Boyd and Vandenberghe.]

Affine Invariance of Newton step

Lemma

Newton step is affine invariant.

For example, let $g(y) = f(Ay)$, y^+ be Newton update on function $g(\cdot)$, and x^+ be Newton update on function $f(\cdot)$. Then if $x = Ay$, we have $x^+ = Ay^+$.

Remark

Affine Invariance indicates that Newton Method is invulnerable to the selection of coordinate system.

Remark

Gradient Descent Method is not affine invariant. This means that bad coordinate choice may limit the power of Gradient Descent Method.

Proof of Affine Invariance of Newton step

Let $x = Ay$ and $g(y) = f(Ay)$, then we have

$$\nabla_y^2 g(y) = \nabla_y^2 f(Ay) = A^T \nabla_x^2 f(x) A \quad (9)$$

$$\nabla_y g(y) = \nabla_y f(Ay) = A^T \nabla_x f(x) \quad (10)$$

Newton update y^+ for $g(\cdot)$ can be expanded as

$$\begin{aligned} y^+ &= y - t(\nabla_y^2 g(y))^{-1} \nabla_y g(y) \\ &= y - t(A^T \nabla_x^2 f(x) A)^{-1} A^T \nabla_x f(x) \\ &= y - t A^{-1} \nabla_x^2 f(x)^{-1} \nabla_x f(x) \end{aligned} \quad (11)$$

Multiply both sides by A ,

$$\begin{aligned} Ay^+ &= Ay - A \cdot t A^{-1} \nabla_x^2 f(x)^{-1} \nabla_x f(x) \\ &= x - t \nabla_x^2 f(x)^{-1} \nabla_x f(x) \\ &= x^+ \end{aligned} \quad (12)$$

Convergence Analysis: Assumption

We assume $f(\cdot)$ satisfies the following:

- $f(\cdot)$ is strongly convex, such that

$$ml \preceq \nabla^2 f(x) \preceq MI, \quad \forall x \quad (13)$$

- $\nabla^2 f(x)$ is L -Lipschitz with constant $L > 0$, such that

$$\|\nabla^2 f(y) - \nabla^2 f(x)\|_2 \leq L\|x - y\|_2, \quad \forall x, y \quad (14)$$

Note that induced matrix norm $\|\cdot\|_2$ equals to the largest singular value of the matrix.

Convergence Analysis: Theorem

Theorem (Part I)

There exist η , γ , where $0 < \eta \leq \frac{m^2}{L}$, $\gamma = \frac{\alpha\beta m}{M^2} \eta^2$ such that Newton method with BTLS has two phases:

(a) *Global or Damped Phase*: If $\|\nabla f(x)\|_2 \geq \eta$, then

$$f(x^+) - f(x) \leq -\gamma \quad (15)$$

Inequality (15) has three implications:

- Every Newton step with BTLS gets closer to global optima by at least γ .
- Damped phase has at most $\frac{f(x^{(0)}) - f^*}{\gamma}$ iterations.
- The damped phase is also at least linearly convergent if not better, i.e. $f(x^+) - f^* \leq c(f(x) - f^*)$

Convergence Analysis: Theorem

Theorem (Part II)

(b) *Local or Quadratic Phase: If $\|\nabla f(x)\|_2 < \eta$, then BTLS will give $t = 1$ and we have*

$$\frac{L}{2m^2} \|\nabla f(x^+)\|_2 \leq \left(\frac{L}{2m^2} \|\nabla f(x)\|_2 \right)^2 \quad (16)$$

Implication:

- To achieve an accuracy of ε , only $O(\log \log \varepsilon)$ iterations are needed once in the quadratic phase, which is known as quadratic convergence. [detail next slide]

Implication of Part (b)

If we have $a_1(< 1), a_2, \dots$ such that $a_{i+1} \leq a_i^2$, then

$$a_\ell \leq (a_{\ell-1})^2 \leq (a_{\ell-2})^{2^2} \leq (a_{\ell-3})^{2^3} \leq \dots \leq (a_1)^{2^{\ell-1}}$$

for $a_\ell \leq \varepsilon$ need,

$$(a_1)^{2^{\ell-1}} \leq \varepsilon$$

$$2^{\ell-1} \log a_1 \leq \log \varepsilon$$

$$\ell \geq \log \log \frac{1}{\varepsilon} + \text{constant}$$

- In this case to make it more precise, set

$$a_\ell = \frac{L}{2m^2} \|\nabla f(x^{(k+\ell-1)})\|_2 \text{ and note that } a_1 < \frac{L}{2m^2} \eta < \frac{1}{2}$$

Convergence Analysis: Part (a)

For readability of the proof, we will divide the proof into lemmas. Use the following lemma for part (a).

Lemma (Global BTLS)

$t = \frac{m}{M}$ satisfies the exit condition of BTLS.

We will first prove the lemma and then the main result, which is

Theorem (Part (a))

If $\|\nabla f(x)\|_2 \geq \eta$, then $f(x^+) - f(x) \leq -\gamma$, where $\gamma = \frac{\alpha\beta m}{M^2}\eta^2$

■ Notation: For convenience, $g \triangleq \nabla f(x)$ and $H \triangleq \nabla^2 f(x)$

Convergence Analysis: Part (a)

Proof of Global BTLS Lemma

$$\begin{aligned} f(x^+) &= f(x - tH^{-1}g) \\ &\leq f(x) - tg^TH^{-1}g + \frac{M}{2}t^2g^TH^{-1}H^{-1}g \end{aligned} \quad (17)$$

$$\leq f(x) - tg^TH^{-1}g + \frac{M}{2m}t^2g^TH^{-1}g \quad (18)$$

$$= f(x) - \frac{m}{2M}g^TH^{-1}g \quad \text{by setting } t = \frac{m}{M}$$

$$\leq f(x) - \alpha \frac{m}{M}g^TH^{-1}g \quad \text{since } \alpha < \frac{1}{2}$$

Hence, $t = \frac{m}{M}$ satisfies the BTLS exit condition.

Note (17) \Rightarrow (18) follows because,

$$g^TH^{-1}H^{-1}g = g^TH^{-1/2}H^{-1}H^{-1/2}g \leq \frac{1}{m}g^TH^{-1}g$$

Convergence Analysis: Main Proof of Part (a)

$$t \leq \beta \frac{m}{M} \quad (\text{Global BTLS Lemma})$$

$$f(x^+) \leq f(x) - \alpha \left(\beta \frac{m}{M} \right) g^T H^{-1} g \quad (\text{BTLS condition})$$

$$\leq f(x) - \alpha \left(\beta \frac{m}{M} \right) \left(\frac{1}{M} \|g\|_2^2 \right) \quad (H^{-1} \preceq I/m)$$

$$= f(x) - \alpha \beta \frac{m}{M^2} \|g\|_2^2$$

$$\leq f(x) - \underbrace{\alpha \beta \frac{m}{M^2} \eta^2}_{\gamma} \quad (19)$$

$$\implies f(x^+) - f(x) = -\gamma \quad (20)$$

Convergence Analysis: Part (b)

We will use the following lemma for this part.

Lemma (Local BTLS)

With the assumptions in (b), $t = 1$ satisfies the exit condition of BTLS.

Proof: See lecture note and text book.

Theorem (Part (b))

If $\|\nabla f(x)\|_2 < \eta$, then $\frac{L}{2m^2} \|\nabla f(x^+)\|_2 \leq \left(\frac{L}{2m^2} \|\nabla f(x)\|_2 \right)^2$

Convergence Analysis: Part (b)

Let $x^+ = x - H^{-1}g$ (BTLS Quad. Lemma)

$$\|\nabla f(x^+)\|_2 = \|\nabla f(x - H^{-1}g) - g + HH^{-1}g\|_2 \quad (\text{Add zero})$$

$$= \left\| \int_0^1 \nabla^2 f(x - tH^{-1}g)(-H^{-1}g) + HH^{-1}g dt \right\|_2$$

(Fund. Theorem of Calculus)

$$= \left\| \int_0^1 (\nabla^2 f(x - tH^{-1}g) - H)(-H^{-1}g) dt \right\|_2$$

(Rearrange)

$$\leq \int_0^1 \|(\nabla^2 f(x - tH^{-1}g) - H)\|_2 \|(-H^{-1}g)\|_2 dt$$

(Triangle inequality of norms)

Convergence Analysis: Part (b) Cont.

$$\begin{aligned}\|\nabla f(x^+)\|_2 &\leq \int_0^1 \|(\nabla^2 f(x - tH^{-1}g) - H)\|_2 \|H^{-1}g\|_2 dt \\ &\leq \int_0^1 L \| -tH^{-1}g \|_2 \|H^{-1}g\|_2 dt \\ &\quad \text{(Lipschitz Continuity of Hessian)} \\ &= L \|H^{-1}g\|_2^2 \int_0^1 t dt = \frac{L}{2} \|H^{-1}g\|_2^2 \\ &\leq \frac{L}{2m^2} \|g\|_2^2 \quad \text{(Strong convexity } (H^{-1} \preceq I/m)) \\ \implies \frac{L}{2m^2} \|\nabla f(x^+)\|_2 &\leq \left(\frac{L}{2m^2} \|g\|_2 \right)^2\end{aligned}\tag{21}$$

Summary

- Mainly cover convergence analysis of Newton method
- Newton method is affine invariant
- There are two phases
 - Damped phase: linear convergence
 - Quadratic phase: quadratic convergence