

Goal: In this lecture we talk about Newton's method and its convergence properties.

1 Introduction to Newton's Method

The convergence of gradient descent depends inherently on the condition number; a change of coordinates can improve the condition number, and hence we should do a change of coordinates at each step. There are indeed several motivations for the Newton step. Before we discuss these, we give the basic definition of the Newton step.

Definition 1. For f having positive definite Hessian $\nabla^2 f(x) \succ 0$, the Newton updating rule with step size t is defined as,

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

where, Δx_{nt} is called the Newton step.

Note that since $\nabla^2 f(x) \succ 0$ we have,

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

always holds for $\nabla f(x) \neq 0$, so this is a descent direction.

There are various ways to interpret this choice of updating rule.

Minimizer of Quadratic Approximation

Consider a quadratic approximation of f around x ,

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

This quadratic function is minimized at $v^* = -\nabla^2 f(x)^{-1} \nabla f(x)$. Note that if f is quadratic, this approximation is exact and $x + v^*$ is the exact minimizer of f .

Linear Approximation of Gradient around x

Consider a linear approximation of $\nabla f(x+v)$,

$$\nabla f(x+v) \simeq \nabla f(x) + \nabla^2 f(x)v. \quad (3)$$

The Newton updating rule is obtained by setting the right hand side to 0, which is an approximation to the optimality condition $\nabla f(x^*) = 0$.

2 Convergence of Newton's Method

We make two major assumptions in this analysis:

1. f is strongly convex, such that $mI \preceq \nabla^2 f(x) \preceq MI$.
2. $\nabla^2 f(x)$ is Lipschitz continuous with constant $L > 0$, i.e.

$$\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq L\|x - y\|_2 \quad \text{for all } x, y. \quad (4)$$

The first assumption is the standard assumption we made for gradient descent, namely, that our function f is smooth (in the sense that it has an upper bound on curvature), and it is strongly convex (it has a lower bound on curvature). The assumption also requires that f be twice differentiable, for the definition to make sense.

Beyond the assumptions for gradient descent, we also need some control on the smoothness of the Hessian. This is given in the second assumption. There, note that the norm on the left is the spectral norm, defined as the largest singular value (which coincides with the largest eigenvalue, for positive semidefinite matrices). L can be interpreted as a bound on the third derivative of f . The smaller L is, the better f can be approximated by a quadratic function. Since each step of Newton's method minimizes a quadratic approximation of f , the performance of Newton's method works best for functions with small L .

Lemma 1. *Consider a twice differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Further assume that $\nabla^2 f(x)$ is Lipschitz continuous with constant $L > 0$. Then, we have*

$$f(y) \leq f(x) + \nabla f(x)^\top (y - x) + \frac{1}{2} \nabla f(x)^\top \nabla^2 f(x) \nabla f(x) + \frac{L}{6} \|y - x\|^3$$

Proof. Exercise. □

Lemma 2. *Consider a twice differentiable function $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Further assume that $\nabla^2 f(x)$ is Lipschitz continuous with constant $L > 0$. Then, we have*

$$\|\nabla f(y) - \nabla f(x) - \nabla^2 f(x)(y - x)\| \leq \frac{L}{2} \|y - x\|^2$$

Proof. According to the fundamental theorem of calculus we have

$$\begin{aligned} \nabla f(y) &= \nabla f(x) + \int_0^1 \nabla^2 f((1 - \alpha)x + \alpha y)(y - x) d\alpha \\ &= \nabla f(x) + \nabla^2 f(x)(y - x) + \int_0^1 (\nabla^2 f((1 - \alpha)x + \alpha y) - \nabla^2 f(x))(y - x) d\alpha \end{aligned}$$

and therefore we have

$$\nabla f(y) - \nabla f(x) - \nabla^2 f(x)(y - x) = \int_0^1 (\nabla^2 f((1 - \alpha)x + \alpha y) - \nabla^2 f(x))(y - x) d\alpha$$

By computing the norm of both sides and using L -smoothness of the Hessians we have

$$\begin{aligned} \|\nabla f(y) - \nabla f(x) - \nabla^2 f(x)(y - x)\| &= \left\| \int_0^1 (\nabla^2 f((1 - \alpha)x + \alpha y) - \nabla^2 f(x))(y - x) d\alpha \right\| \\ &\leq \int_0^1 \|\nabla^2 f((1 - \alpha)x + \alpha y) - \nabla^2 f(x)\| \|y - x\| d\alpha \\ &\leq \int_0^1 L\alpha \|y - x\|^2 d\alpha = \frac{L}{2} \|y - x\|^2 \end{aligned}$$

□

For notational convenience we denote $g = \nabla f(x)$ and $H = \nabla^2 f(x)$ from this point on. Our main result of this lecture is Theorem 1. We devote the rest of the lecture trying to understand and prove this theorem. Before stating the theorem, let us first recall Backtracking Line Search (BTLS). With BTLS, first α and β are chosen such that $0 < \alpha < \frac{1}{2}$ and $0 < \beta < 1$, starting with $t = 1$, repeat

```

while true
  if  $f(x + t\Delta x) \leq f(x) + \alpha t g^T \Delta x$ 
    exit
  else
     $t \leftarrow \beta t$ 
  end
end

```

We are now ready to state and prove the theorem.

Theorem 1. *Consider the update of Newton's method using BTLS with parameters $0 < \alpha < 0.5$ and $0 < \beta < 1$. Further, consider the definitions $\eta := \min\{1, 3(1 - 2\alpha)\} \frac{m^2}{L}$ and $\gamma := \alpha\beta\eta^2 \frac{m}{M^2}$. Then we have*

(a). *Damped Newton Phase: If $\|g\|_2 \geq \eta$ then $f(x^+) - f(x) \leq -\gamma$.*

(b). *Quadratic Phase: If $\|g\|_2 < \eta$ then BTLS $t = 1$ and*

$$\frac{L}{2m^2} \|\nabla f(x^+)\|_2 \leq \left(\frac{L}{2m^2} \|\nabla f(x)\|_2 \right)^2. \quad (5)$$

2.1 Convergence Proof

For readability of the proof, we divide the proof into lemmas to emphasize the flow of the proof and not get lost in the details of the derivation.

Lemma 3. *With the assumptions in part (a), $t = \frac{m}{M}$ satisfies BTLS exit condition, i.e., $f(x + t\Delta x_{\text{nt}}) \leq f(x) + \alpha t g^T \Delta x_{\text{nt}}$.*

Proof.

$$f(x^+) = f(x - tH^{-1}g) \quad (6)$$

$$\leq f(x) - t g^T H^{-1} g + \frac{M}{2} t^2 g^T H^{-1} H^{-1} g \quad (7)$$

Note that¹,

$$g^T H^{-1} H^{-1} g = g^T H^{-1/2} H^{-1} H^{-1/2} g \leq \frac{1}{m} g^T H^{-1} g \quad (8)$$

¹Recall the definition of square root of a matrix. If A is positive definite then we can write $A = U\Lambda U^T$, where U is unitary and Λ is diagonal, and $A^{1/2} = U\Lambda^{1/2}U^T$.

Thus,

$$f(x^+) \leq f(x) - tg^T H^{-1}g + \frac{M}{2m}t^2 g^T H^{-1}g \quad (9)$$

Setting $t = \frac{m}{M}$,

$$f(x^+) \leq f(x) - \frac{m}{2M}g^T H^{-1}g \quad (10)$$

This satisfies the exit condition for $t = \frac{m}{M}$, $f(x^+) \leq f(x) - \alpha \frac{m}{M}g^T H^{-1}g$ since $\alpha < 1/2$. \square

This result shows that if we are in phase (a) then $t \geq \beta m/M$.

Lemma 4. *With the assumptions in part (b), $t = 1$ satisfies BTLS exit condition.*

Proof. In this proof we find $\alpha < \frac{1}{2}$ such that $t = 1$ satisfies BTLS exit condition. Our goal is to find α such that,

$$f(x + \Delta x_{\text{nt}}) \leq f(x) + \alpha g^T \Delta x_{\text{nt}}. \quad (11)$$

For notational convenience we denote

$$\lambda(x) = (\Delta x_{\text{nt}}^T H \Delta x_{\text{nt}})^{1/2} = (g^T H^{-1}g)^{1/2} \quad (12)$$

which is known as the Newton decrement at x . The second equality follows because $\Delta x_{\text{nt}} = -H^{-1}g$. By Lemma 1 we know that if we follow the update of Newton's method with stepsize $t = 1$, i.e., $x^+ = x - H^{-1}g$, we have

$$f(x^+) \leq f(x) - g^T H^{-1}g + \frac{1}{2}g^T H^{-1}g + \frac{L}{6}\|H^{-1}g\|^3$$

which can be upper bounded by

$$f(x^+) \leq f(x) - g^T H^{-1}g + \frac{1}{2}g^T H^{-1}g + \frac{L}{6m^{3/2}}\|H^{-1/2}g\|^3$$

and simplified as

$$\begin{aligned} f(x + \Delta x_{\text{nt}}) &\leq f(x) - \frac{1}{2}\lambda^2(x) + \frac{L}{6m^{3/2}}\lambda^3(x) \\ &= f(x) - \lambda^2(x) \left(\frac{1}{2} - \frac{L\lambda(x)}{6m^{3/2}} \right) \\ &= f(x) + g^T \Delta x_{\text{nt}} \left(\frac{1}{2} - \frac{L\lambda(x)}{6m^{3/2}} \right) \end{aligned}$$

Note that since $\eta \leq 3(1 - 2\alpha)\frac{m^2}{L}$ we can show that

$$\lambda(x) = (g^T H^{-1}g)^{1/2} \leq \frac{1}{m^{1/2}}\|g\|_2 < \frac{1}{m^{1/2}}\eta \leq 3(1 - 2\alpha)\frac{m^{3/2}}{L}.$$

By regrouping the terms we can show that

$$\alpha < \frac{1}{2} - \frac{L\lambda(x)}{6m^{3/2}}$$

and hence we have

$$f(x + \Delta x_{\text{nt}}) \leq f(x) + \alpha g^T \Delta x_{\text{nt}}$$

then $t = 1$ satisfies BTLS exit condition. \square

Proof of Theorem 1

Theorem 1 part (a). Using Lemma 3, with $t \geq \beta \frac{m}{M}$ and substituting $\Delta x_{\text{nt}} = -H^{-1}g$, we have

$$f(x^+) \leq f(x) - \alpha\beta \frac{m}{M} g^T H^{-1} g. \quad (13)$$

By strong convexity $H \preceq MI$, so $H^{-1} \succeq \frac{1}{M}I$, and we have

$$g^T H^{-1} g \geq \frac{1}{M} \|g\|_2^2, \quad (14)$$

therefore,

$$f(x^+) \leq f(x) - \alpha\beta \frac{m}{M} \frac{1}{M} \|g\|_2^2 \quad (15)$$

$$f(x^+) - f(x) \leq - \underbrace{\alpha\beta \frac{m}{M^2} \eta^2}_{\gamma} \quad (16)$$

where the last inequality follows because $\|g\|_2 \geq \eta$. \square

Theorem 1 part (b). Using Lemma 4, we can set $t = 1$. By using Lemma 2 we can write that

$$\|\nabla f(x^+) - \nabla f(x) - \nabla^2 f(x)(x^+ - x)\| \leq \frac{L}{2} \|x^+ - x\|^2$$

which can be simplified as

$$\|\nabla f(x^+)\| \leq \frac{L}{2} \|H^{-1}g\|^2$$

Also, $H \succeq mI$ so $H^{-1} \preceq \frac{1}{m}I$ and,

$$\|H^{-1}g\|_2^2 \leq \frac{1}{m^2} \|g\|_2^2. \quad (17)$$

Substituting this,

$$\|\nabla f(x^+)\|_2 \leq \frac{L}{2m^2} \|g\|_2^2, \quad (18)$$

and multiplying both sides by $\frac{L}{2m^2}$ we obtain the result stated in the theorem. \square

Implication of (a)

In the damped Newton phase, f decreases by at least γ at each iteration, there the total of iterations in this phase cannot exceed,

$$\frac{f(x^{(0)}) - p^*}{\gamma}$$

since otherwise $f(x)$ would be less than p^* , which contradicts the optimality of p^* . In other words, the quadratic phase starts after at most $\frac{f(x^{(0)}) - p^*}{\gamma}$ iterations.

Implication of (b)

Let k be the first iteration in which $\|g\| < \eta$. And let $\ell \geq 0$ be the number of iterations after k . For simplicity, let us define:

$$a_\ell = \frac{L}{2m^2} \|\nabla f(x^{(k+\ell-1)})\|_2 \quad (19)$$

First, let us establish a bound on a_1 . In the quadratic phase, since $\|\nabla f(x^{(k)})\|_2 < \eta$ and $\eta < \frac{m^2}{L}$ by assumption, we have

$$\frac{L}{2m^2} \|\nabla f(x^{(k)})\|_2 < \frac{L}{2m^2} \eta < \frac{1}{2}.$$

Thus, $a_1 < \frac{1}{2}$. Now from (b) of the theorem, we also have that $a_{\ell+1} \leq a_\ell^2$. Therefore, we have the following sequence:

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

$$\implies a_\ell \leq \left(\frac{1}{2}\right)^{2^{\ell-1}}$$

$$\begin{aligned} \implies \frac{L}{2m^2} \|\nabla f(x^{(\ell)})\|_2 &\leq \left(\frac{1}{2}\right)^{2^{\ell-1}} \\ \implies \|\nabla f(x^{(\ell)})\|_2 &\leq \frac{2m^2}{L} \left(\frac{1}{2}\right)^{2^{\ell-1}} \end{aligned}$$

For strongly convex functions, we also have

$$f(x^{(\ell)}) - p^\star \leq \frac{1}{2m} \|\nabla f(x^{(\ell)})\|_2^2 \quad (20)$$

$$\leq \frac{1}{2m} \left(\frac{2m^2}{L} \left(\frac{1}{2}\right)^{2^{\ell-1}} \right)^2 \quad (21)$$

$$\leq \frac{2m^3}{L^2} \left(\frac{1}{2}\right)^{2^{\ell-1}} \quad (22)$$

thus, $f(x) \rightarrow p^\star$ quadratically.

Therefore, if we want $a_\ell \leq \epsilon$, we only need the following number of iterations:

$$\begin{aligned} (a_1)^{2^{\ell-1}} &\leq \epsilon \\ 2^{\ell-1} \log a_1 &\leq \log \epsilon \\ 2^{\ell-1} &\geq \text{constant} \times \log \epsilon \\ \ell - 1 &\geq \log \log \epsilon + \text{constant}. \end{aligned}$$