

# Statistical Optimization: Lecture 10

Newton's method

Zijian Guo

Zhejiang University  
Center for Data Science

April 13, 2026

# 1-dimensional Newton method

---

Newton's method was originally introduced for solving

$$f(\theta) = 0.$$

Starting from  $\theta_0 \in \mathbb{R}$ , it computes

$$\theta_{t+1} = \theta_t - \frac{f(\theta_t)}{f'(\theta_t)}, \quad t \geq 0.$$

Geometric view:  $\theta_{t+1}$  is the intersection of the tangent line at  $(\theta_t, f(\theta_t))$  with the horizontal axis, i.e. it solves

$$f(\theta_t) + f'(\theta_t)(\theta - \theta_t) = 0.$$

**Remark.** We can apply this idea to find zero gradient.

# 1-dimensional Newton method

---

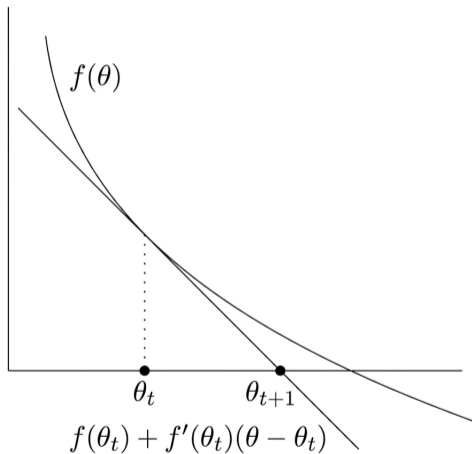


Figure: One step of Newton's method

# Newton's method for optimization

---

Suppose we want to minimize a twice differentiable function

$$f: \mathbb{R} \rightarrow \mathbb{R}.$$

In one dimension, this can be reduced to finding a zero of  $f'$ . Applying Newton's method to  $f'$ , we obtain

$$\theta_{t+1} := \theta_t - \frac{f'(\theta_t)}{f''(\theta_t)} = \theta_t - (f''(\theta_t))^{-1} f'(\theta_t), \quad t \geq 0. \quad (1)$$

There is no reason to restrict to  $d = 1$ . For

$$f: \mathbb{R}^d \rightarrow \mathbb{R},$$

Newton's method iterates

$$\theta_{t+1} := \theta_t - (\nabla^2 f(\theta_t))^{-1} \nabla f(\theta_t), \quad t \geq 0. \quad (2)$$

## Second-order intuition and gradient descent

---

To minimize a smooth function  $f$ , we approximate it locally near  $\theta_t$  by a **second-order Taylor expansion**:

$$f(\theta_t + \rho_t) \approx f(\theta_t) + \nabla f(\theta_t)^\top \rho_t + \frac{1}{2} \rho_t^\top \nabla^2 f(\theta_t) \rho_t.$$

Minimizing this quadratic model (when  $\nabla^2 f(\theta_t)$  is positive definite) gives

$$\rho_t = -(\nabla^2 f(\theta_t))^{-1} \nabla f(\theta_t).$$

If we replace the Hessian by a scaled identity matrix,

$$\nabla^2 f(\theta_t) \approx \frac{1}{\gamma_t} I,$$

then minimizing this simplified model yields

$$\rho_t = -\gamma_t \nabla f(\theta_t).$$

which is exactly the gradient descent update.

## Newton step minimizes the quadratic model

---

### Theorem

Let  $f$  be convex and twice differentiable at  $\theta_t$ , and assume

$$\nabla^2 f(\theta_t) \succ 0$$

is invertible. Then the Newton iterate  $\theta_{t+1}$  satisfies

$$\theta_{t+1} = \arg \min_{\theta \in \mathbb{R}^d} \left[ f(\theta_t) + \nabla f(\theta_t)^\top (\theta - \theta_t) + \frac{1}{2} (\theta - \theta_t)^\top \nabla^2 f(\theta_t) (\theta - \theta_t) \right].$$

## Proof of the theorem

---

**Proof.** Consider the quadratic function

$$m_t(\theta) := f(\theta_t) + \nabla f(\theta_t)^\top (\theta - \theta_t) + \frac{1}{2}(\theta - \theta_t)^\top \nabla^2 f(\theta_t)(\theta - \theta_t).$$

Since  $\nabla^2 f(\theta_t) \succ 0$ , the function  $m_t$  is strictly convex and has a unique minimizer. Setting its gradient to zero gives

$$\nabla m_t(\theta) = \nabla f(\theta_t) + \nabla^2 f(\theta_t)(\theta - \theta_t) = 0.$$

Hence

$$\theta = \theta_t - (\nabla^2 f(\theta_t))^{-1} \nabla f(\theta_t) = \theta_{t+1}.$$

Therefore  $\theta_{t+1}$  is the unique minimizer of  $m_t$ . □

# Gradient step minimizes the simplified quadratic model

---

## Theorem

Let  $f$  be differentiable at  $\theta_t$ , then the gradient descent iterate  $\theta_{t+1}$  satisfies

$$\theta_{t+1} = \arg \min_{\theta \in \mathbb{R}^d} \left[ f(\theta_t) + \nabla f(\theta_t)^\top (\theta - \theta_t) + \frac{1}{2\gamma_t} (\theta - \theta_t)^\top (\theta - \theta_t) \right].$$

## Theorem (Quadratic Convergence)

Let  $f : \text{dom}(f) \rightarrow \mathbb{R}$  be twice differentiable with a critical point  $\theta^*$ . Suppose there is a ball  $\Theta \subseteq \text{dom}(f)$  centered at  $\theta^*$  such that:

**(i) Bounded inverse Hessians:**  $\exists \mu > 0$  with

$$\|(\nabla^2 f(\theta))^{-1}\| \leq \frac{1}{\mu}, \quad \forall \theta \in \Theta.$$

**(ii) Lipschitz continuous Hessians:**  $\exists B \geq 0$  with

$$\|\nabla^2 f(\theta) - \nabla^2 f(\eta)\| \leq B\|\theta - \eta\|, \quad \forall \theta, \eta \in \Theta.$$

(Here  $\|\cdot\|$  is the spectral norm.) Then for  $\theta_t \in \Theta$  and  $\theta_{t+1}$  given by Newton's method,

$$\|\theta_{t+1} - \theta^*\| \leq \frac{B}{2\mu} \|\theta_t - \theta^*\|^2.$$

**Proof.** By Newton's method,

$$\theta_{t+1} - \theta^* = \theta_t - \theta^* - (\nabla^2 f(\theta_t))^{-1} \nabla f(\theta_t).$$

Since  $\nabla f(\theta^*) = 0$ , we expand  $\nabla f(\theta^*)$  around  $\theta_t$ :

$$\nabla f(\theta^*) = \nabla f(\theta_t) + \int_0^1 \nabla^2 f(\theta_t + s(\theta^* - \theta_t)) (\theta^* - \theta_t) ds.$$

Therefore,

$$0 = \nabla f(\theta_t) + \nabla^2 f(\theta_t) (\theta^* - \theta_t) + \int_0^1 \left( \nabla^2 f(\theta_t + s(\theta^* - \theta_t)) - \nabla^2 f(\theta_t) \right) (\theta^* - \theta_t) ds.$$

Rearranging gives

$$\nabla f(\theta_t) = \nabla^2 f(\theta_t) (\theta_t - \theta^*) - \int_0^1 \left( \nabla^2 f(\theta_t + s(\theta^* - \theta_t)) - \nabla^2 f(\theta_t) \right) (\theta^* - \theta_t) ds.$$

Substituting this identity into the Newton step, we obtain

$$\theta_{t+1} - \theta^* = (\nabla^2 f(\theta_t))^{-1} \int_0^1 \left( \nabla^2 f(\theta_t + s(\theta^* - \theta_t)) - \nabla^2 f(\theta_t) \right) (\theta^* - \theta_t) ds.$$

Taking norms, and using the spectral norm, gives

$$\|\theta_{t+1} - \theta^*\| \leq \|(\nabla^2 f(\theta_t))^{-1}\| \int_0^1 \left\| \left( \nabla^2 f(\theta_t + s(\theta^* - \theta_t)) - \nabla^2 f(\theta_t) \right) (\theta^* - \theta_t) \right\| ds.$$

Hence

$$\|\theta_{t+1} - \theta^*\| \leq \|(\nabla^2 f(\theta_t))^{-1}\| \|\theta^* - \theta_t\| \int_0^1 \left\| \nabla^2 f(\theta_t + s(\theta^* - \theta_t)) - \nabla^2 f(\theta_t) \right\| ds.$$

By assumptions (i) and (ii),

$$\|\theta_{t+1} - \theta^*\| \leq \frac{1}{\mu} \|\theta^* - \theta_t\| \int_0^1 B \|s(\theta^* - \theta_t)\| ds \leq \frac{B}{2\mu} \|\theta_t - \theta^*\|^2.$$

Therefore,

$$\|\theta_{t+1} - \theta^*\| \leq \frac{B}{\mu} \|\theta^* - \theta_t\|^2 \int_0^1 s ds = \frac{B}{2\mu} \|\theta_t - \theta^*\|^2.$$

# Convergence of Newton's method

## Corollary (Convergence of Newton's method)

With the assumptions of Theorem, if  $\theta_0 \in \Theta$  satisfies

$$\|\theta_0 - \theta^*\| \leq \frac{\mu}{B},$$

then Newton's method yields

$$\|\theta_T - \theta^*\| \leq \frac{\mu}{B} \left(\frac{1}{2}\right)^{2^T - 1}, \quad T \geq 0.$$

**Proof.** Let  $e_t := \|\theta_t - \theta^*\|$ . Theorem gives  $e_{t+1} \leq \frac{B}{2\mu} e_t^2$ . Assuming  $e_t \leq \frac{\mu}{B} \left(\frac{1}{2}\right)^{2^t - 1}$ ,

$$e_{t+1} \leq \frac{B}{2\mu} \left(\frac{\mu}{B}\right)^2 \left(\frac{1}{2}\right)^{2(2^t - 1)} = \frac{\mu}{B} \left(\frac{1}{2}\right)^{2^{t+1} - 1}.$$