

CS 2601 Linear and Convex Optimization

7. Newton's method

Bo Jiang

John Hopcroft Center for Computer Science
Shanghai Jiao Tong University

Fall 2022

Outline

- Newton's method and properties
- Analysis of Newton's method
- Damped Newton's method

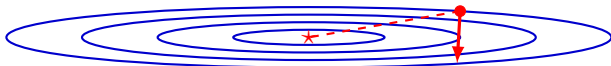
Better descent direction

Gradient descent uses first-order information (i.e. gradient),

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

Locally $-\nabla f(\mathbf{x}_k)$ is the max-rate descending direction, but globally it may not be the “right” direction.

Example. For $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Q}\mathbf{x}$ with $\mathbf{Q} = \text{diag}\{0.01, 1\}$, minimum is $\mathbf{x}^* = \mathbf{0}$.



The negative gradient is

$$-\nabla f(\mathbf{x}) = -\mathbf{Q}\mathbf{x} = -(0.01x_1, x_2)^T$$

quite different from the “right” descent direction $\mathbf{d} = -\mathbf{x}$. Note

$$\mathbf{d} = -\mathbf{Q}^{-1}\nabla f(\mathbf{x}) = -[\nabla^2 f(\mathbf{x})]^{-1}\nabla f(\mathbf{x}) \quad \text{加入二阶信息}$$

With second-order information (i.e. Hessian), we hope to do better.

Newton's method

Gradient step $\mathbf{x}_{k+1} = \mathbf{x}_k - t_k \nabla f(\mathbf{x}_k)$ can be interpreted as minimizing a quadratic approximation of f at \mathbf{x}_k , 具体内容写在黑板上

$$f(\mathbf{x}) \approx \hat{f}_{\text{gd}}(\mathbf{x}) \triangleq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x} - \mathbf{x}_k) + \frac{1}{2t_k} \|\mathbf{x} - \mathbf{x}_k\|^2$$

Newton's method minimizes the second-order Taylor approximation,

$$f(\mathbf{x}) \approx \hat{f}_{\text{nt}}(\mathbf{x}) \triangleq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x} - \mathbf{x}_k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_k)^T \nabla^2 f(\mathbf{x}_k) (\mathbf{x} - \mathbf{x}_k)$$

Newton step. Assuming $\nabla^2 f(\mathbf{x}_k) \succ \mathbf{0}$, 需要可逆

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

Note. If f is quadratic, then $f = \hat{f}_{\text{nt}}$, and Newton's method gets to the minimum in a single step starting from any \mathbf{x}_0 .

Newton's method (cont'd)

- 1: initialization $\mathbf{x} \leftarrow \mathbf{x}_0 \in \mathbb{R}^n$
- 2: **while** $\|\nabla f(\mathbf{x})\| > \delta$ **do**
- 3: $\mathbf{x} \leftarrow \mathbf{x} - [\nabla^2 f(\mathbf{x})]^{-1} \nabla f(\mathbf{x})$
- 4: **end while**
- 5: **return** \mathbf{x}

Note. As in the case of gradient descent, other stopping criteria can be used. [BV] uses $\nabla f(\mathbf{x})[\nabla^2 f(\mathbf{x})]^{-1} \nabla f(\mathbf{x}) > \delta$.

The Newton step is a special case of $\mathbf{x}_{k+1} = \mathbf{x}_k + t_k \mathbf{d}_k$ with

- **Newton direction** $\mathbf{d}_k = -[\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$
- **constant step size** $t_k = 1$

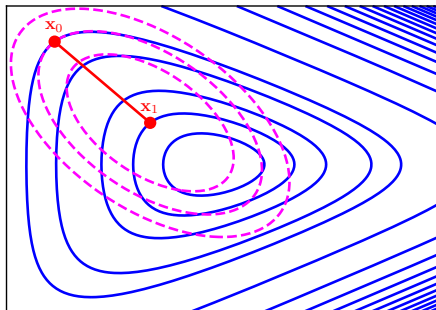
For $\nabla^2 f(\mathbf{x}_k) \succ \mathbf{O}$, the Newton direction is a descent direction

$$\nabla f(\mathbf{x}_k)^T \mathbf{d}_k = -\nabla f(\mathbf{x}_k)^T [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k) < 0 \quad \text{if } \nabla f(\mathbf{x}_k) \neq \mathbf{0}$$

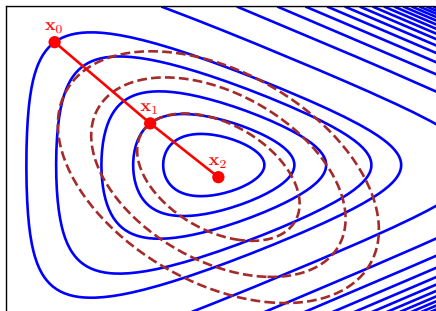
Newton's method (cont'd)

洋红色

The magenta curves are the level curves of the quadratic approximation of f at x_0



The brown curves are the level curves of the quadratic approximation of f at x_1 .



Affine invariance

Given f , and invertible $A \in \mathbb{R}^{n \times n}$, let $g(\mathbf{y}) = f(A\mathbf{y})$. By the chain rule,

$$\nabla g(\mathbf{y}) = A^T \nabla f(A\mathbf{y}), \quad \nabla^2 g(\mathbf{y}) = A^T \nabla^2 f(A\mathbf{y}) A$$

If we run Newton's method on f and g with $\mathbf{x}_0 = A\mathbf{y}_0$, then

$$\begin{aligned} \mathbf{y}_1 &= \mathbf{y}_0 - [\nabla^2 g(\mathbf{y}_0)]^{-1} \nabla g(\mathbf{y}_0) \\ &= \mathbf{y}_0 - [A^T \nabla^2 f(\mathbf{x}_0) A]^{-1} A^T \nabla f(\mathbf{x}_0) \\ &= \mathbf{y}_0 - A^{-1} [\nabla^2 f(\mathbf{x}_0)]^{-1} \nabla f(\mathbf{x}_0) \\ &= A^{-1} [\mathbf{x}_0 - [\nabla^2 f(\mathbf{x}_0)]^{-1} \nabla f(\mathbf{x}_0)] \\ &= A^{-1} \mathbf{x}_1 \end{aligned}$$

By induction $\mathbf{x}_k = A\mathbf{y}_k$. Same progress independent of scaling by A .

For gradient descent, if $AA^T \neq I$, then in general,

$$\mathbf{x}_1 = \mathbf{x}_0 - t \nabla f(\mathbf{x}_0) \neq A\mathbf{y}_1 = A(\mathbf{y}_0 - t A^T \nabla f(\mathbf{y}_0)) = \mathbf{x}_0 - t A A^T \nabla f(\mathbf{x}_0)$$

Connection to root finding

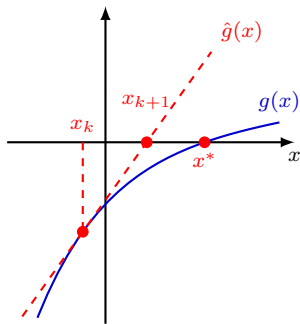
Newton's method is originally an algorithm for solving $g(x) = 0$.

By the first-order Taylor expansion,

$$g(x) \approx \hat{g}(x) \triangleq g(x_k) + g'(x_k)(x - x_k)$$

Use the root of $\hat{g}(x)$ as the next approximation

$$x_{k+1} = x_k - \frac{g(x_k)}{g'(x_k)}$$



Example (computing \sqrt{C}). \sqrt{C} is a root of $g(x) = x^2 - C$. Newton's method yields

$$x_{k+1} = x_k - \frac{x_k^2 - C}{2x_k} = \frac{1}{2} \left(x_k + \frac{C}{x_k} \right)$$

For $x_0 > 0$, x_k converges to \sqrt{C} .

Connection to root finding (cont'd)

Back to the optimization problem,

$$x^* = \operatorname{argmin}_x f(x) \iff f'(x^*) = 0$$

Letting $g = f'$ in Newton's root finding algorithm,

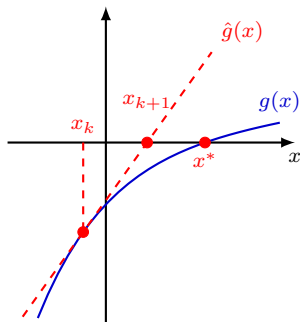
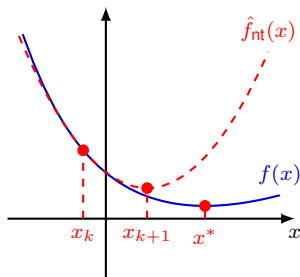
$$x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)} = x_k - [f''(x_k)]^{-1} f'(x_k)$$

In n -dimension, $f' \rightarrow \nabla f, f'' \rightarrow \nabla^2 f$.

$$\mathbf{x}^* = \operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}) \iff \nabla f(\mathbf{x}^*) = \mathbf{0}$$

Newton's algorithm becomes

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



Example

$$f(x_1, x_2) = e^{x_1+3x_2-0.1} + e^{x_1-3x_2-0.1} + e^{-x_1-0.1}$$

Newton step at $\mathbf{x}_0 = (-2, 1)^T$.

- gradient

$$\nabla f(\mathbf{x}_0) = e^{-0.1} \begin{pmatrix} e^{x_1+3x_2} + e^{x_1-3x_2} - e^{-x_1} \\ 3e^{x_1+3x_2} - 3e^{x_1-3x_2} \end{pmatrix} \Big|_{\mathbf{x}=\mathbf{x}_0} = \begin{pmatrix} -4.22019458 \\ 7.36051909 \end{pmatrix}$$

- Hessian

$$\begin{aligned} \nabla^2 f(\mathbf{x}_0) &= e^{-0.1} \begin{pmatrix} e^{x_1+3x_2} + e^{x_1-3x_2} + e^{-x_1} & 3e^{x_1+3x_2} - 3e^{x_1-3x_2} \\ 3e^{x_1+3x_2} - 3e^{x_1-3x_2} & 9e^{x_1+3x_2} + 9e^{x_1-3x_2} \end{pmatrix} \Big|_{\mathbf{x}=\mathbf{x}_0} \\ &= \begin{pmatrix} 9.1515943 & 7.36051909 \\ 7.36051909 & 22.19129872 \end{pmatrix} \end{aligned}$$

- Newton direction

$$\nabla^2 f(\mathbf{x}_0)\mathbf{d} = -\nabla f(\mathbf{x}_0) \implies \mathbf{d} = (0.99274936, -0.66096491)^T$$

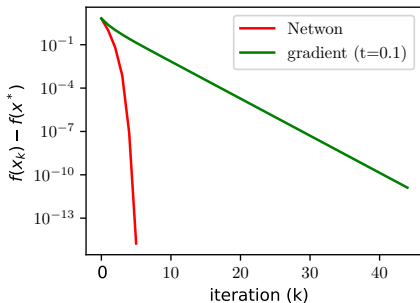
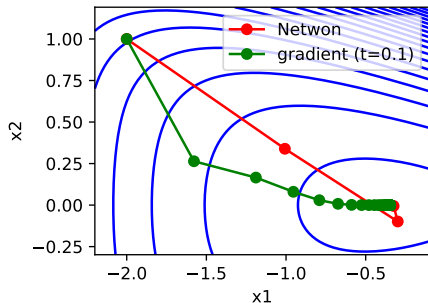
- Newton step

$$\mathbf{x}_1 = \mathbf{x}_0 + \mathbf{d} = (-1.00725064, 0.33903509)^T$$

Example (cont'd)

$$f(x_1, x_2) = e^{x_1+3x_2-0.1} + e^{x_1-3x_2-0.1} + e^{-x_1-0.1}$$

Solution using Newton's method and gradient descent with constant step size 0.1. Initial point $x_0 = (-2, 1)^T$.



- Newton's method takes a more “direct” path
- Newton's method ~~requires much fewer iterations~~, but each iteration is more expensive

Outline

- Newton's method and properties
- Analysis of Newton's method
- Damped Newton's method

Convergence of Newton's method

Example. Consider the minimization of $f(x) = \sqrt{1+x^2}$.

$$f'(x) = \frac{x}{\sqrt{1+x^2}}, \quad f''(x) = \frac{1}{(1+x^2)^{3/2}}$$

The Newton direction is

$$d_k = -\frac{f'(x_k)}{f''(x_k)} = -x_k - x_k^3$$

The Newton step is

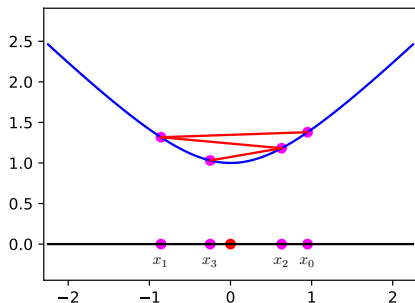
$$x_{k+1} = x_k + d_k = -x_k^3$$

Note $x_k \rightarrow x^* = 0$ iff $|x_0| < 1$.

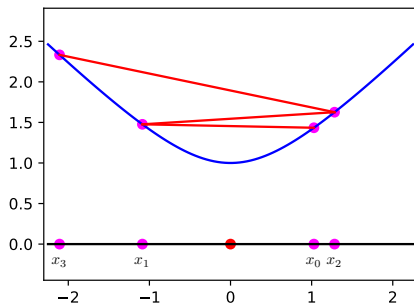
When $|x_0| > 1$, x_k diverges, and

$$f(x_{k+1}) > f(x_k)$$

Convergence of Newton's method (cont'd)



$$x_0 = 0.95$$



$$x_0 = 1.02$$

In general, Newton's method does not guarantee global convergence. When it does converge, the convergence is usually very fast.

Convergence analysis: 1D case

Theorem. If f is m -strongly convex, f'' is M -Lipschitz continuous, and x^* is a minimum of f , then the sequence $\{x_k\}$ produced by Newton's method satisfies

$$|x_{k+1} - x^*| \leq \frac{M}{2m} |x_k - x^*|^2$$

Notes. Let $\xi_k = \frac{M}{2m} |x_k - x^*|$. The above inequality becomes $\xi_{k+1} \leq \xi_k^2$.

- If $\xi_k = 10^{-p}$, then $\xi_{k+1} \leq 10^{-2p}$, the number of significant digits doubles in each iteration!
- If $\xi_0 < 1$ i.e. $|x_0 - x^*| < \frac{2m}{M}$, then $\xi_k \leq \xi_0^{2^k}$ converges to 0 extremely fast. The number of iterations to ensure $\xi_k \leq \epsilon$ is $k \geq \log_2 \log_{\frac{1}{\xi_0}} \frac{1}{\epsilon}$.

For $\epsilon = 10^{-p}$, $k \geq \log_2 p + \log_2 \log_{\frac{1}{\xi_0}} 10$, only logarithmic in the number of digits. Very few iterations are required!

- This theorem is a **local** convergence result. **Fast convergence if x_0 is close enough to x^* , i.e. $|x_0 - x^*| < \frac{2m}{M}$. No guarantee if $|x_0 - x^*|$ is large.**

和初始点选取很有关系！否则可能不会收敛

Proof: 1D case

f'' is M -Lipschitz continuous

$$\begin{aligned} & |x_{k+1} - x^*| \\ &= |x_k - x^* - [f''(x_k)]^{-1} f'(x_k)| \\ &= |f''(x_k)|^{-1} \cdot |f'(x^*) - f'(x_k) - f''(x_k)(x^* - x_k)| \\ &= \frac{|x_k - x^*|}{|f''(x_k)|} \cdot \left| \int_0^1 [f''(x_k + t(x^* - x_k)) - f''(x_k)] dt \right| \\ &\leq \frac{|x_k - x^*|}{|f''(x_k)|} \cdot \int_0^1 |f''(x_k + t(x^* - x_k)) - f''(x_k)| dt \\ &\leq \frac{|x_k - x^*|}{|f''(x_k)|} \cdot \int_0^1 M t |x_k - x^*| dt \\ &= \frac{M}{2|f''(x_k)|} |x_k - x^*|^2 \\ &\leq \frac{M}{2m} |x_k - x^*|^2 \end{aligned}$$

Newton step

$$f'(x^*) = 0$$

Newton-Leibniz

$$\left| \int f \right| \leq \int |f|$$

M -Lipschitz of f''

m -strong convexity

Matrix norm Any function satisfies following conditions can be called matrix norm

The set of $m \times n$ matrices $\mathbb{R}^{m \times n}$ is a mn -dimensional vector space

A **matrix norm** on $\mathbb{R}^{m \times n}$ is a function $\|\cdot\| : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$ s.t.

1. $\|A\| \geq 0, \forall A \in \mathbb{R}^{m \times n}$
2. $\|A\| = 0$ iff $A = \mathbf{O}$
3. $\|cA\| = |c| \cdot \|A\|, \forall c \in \mathbb{R}, A \in \mathbb{R}^{m \times n}$ (**positive homogeneity**)
4. $\|A + B\| \leq \|A\| + \|B\|, \forall A, B \in \mathbb{R}^{m \times n}$ (**triangle inequality**)

Example. The **Frobenius norm** on $\mathbb{R}^{m \times n}$ is the 2-norm on \mathbb{R}^{mn} .

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2} \quad \text{for } A = (a_{ij}) \in \mathbb{R}^{m \times n}$$

Operator norm

定义了一种映射

A matrix $A \in \mathbb{R}^{m \times n}$ defines a linear transformation from \mathbb{R}^n to \mathbb{R}^m

$$A : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$\mathbf{x} \mapsto A\mathbf{x}$$

Given two vector norms $\|\cdot\|_a$ and $\|\cdot\|_b$ on \mathbb{R}^n and \mathbb{R}^m , respectively, the operator norm or induced norm of A is defined by

$$\|A\|_{a,b} = \max_{\mathbf{x}:\mathbf{x} \neq \mathbf{0}} \frac{\|A\mathbf{x}\|_b}{\|\mathbf{x}\|_a} = \max_{\mathbf{x}:\|\mathbf{x}\|_a=1} \|A\mathbf{x}\|_b = \max_{\mathbf{x}:\|\mathbf{x}\|_a \leq 1} \|A\mathbf{x}\|_b$$

可以被证明

Exercise. Show the three definitions are equivalent.

The induced norm has the following important property.

Proposition (compatibility of norms).

$$\|A\mathbf{x}\|_b \leq \|A\|_{a,b} \|\mathbf{x}\|_a$$

Spectral norm

When the norms on \mathbb{R}^n and \mathbb{R}^m are both 2-norms, the induced norm on $\mathbb{R}^{n \times m}$ is simply called the **2-norm** or **spectral norm**, denoted by $\|\cdot\|_2$.

Proposition.

$$\|A\|_2 = \sqrt{\lambda_{\max}(A^T A)},$$

where $\lambda_{\max}(A^T A)$ is the maximum eigenvalue of $A^T A$.

Proof. Let $\|x\|_2 = 1$. By slide 32 of §2,

$$\|Ax\|_2^2 = x^T A^T A x \leq \lambda_{\max}(A^T A) \|x\|_2^2 = \lambda_{\max}(A^T A), \quad \forall x \in \mathbb{R}^n$$

with equality iff x is an eigenvector of $A^T A$ associated with $\lambda_{\max}(A^T A)$.

Corollary. If A is symmetric,

$$\|A\|_2 = \max\{|\lambda_{\max}(A)|, |\lambda_{\min}(A)|\}$$

If $A \succeq O$, then $\|A\|_2 = \lambda_{\max}(A)$.

Examples

Example.

$$\mathbf{A} = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$$

To find the 2-norm,

$$\mathbf{A}^T \mathbf{A} = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} = \begin{pmatrix} 10 & 14 \\ 14 & 20 \end{pmatrix}$$

$$\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})} = \sqrt{15 + \sqrt{221}} \approx 5.465$$

Example.

$$\mathbf{A} = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix} \succeq \mathbf{0}$$

$$\|\mathbf{A}\|_2 = \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})} = \sqrt{\lambda_{\max}(\mathbf{A}^2)} = \sqrt{\lambda_{\max}^2(\mathbf{A})} = \lambda_{\max}(\mathbf{A}) = 5$$

Convergence analysis

$\nabla^2 f$ is M -Lipschitz continuous if

$$\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\|_2 \leq M\|\mathbf{x} - \mathbf{y}\|_2, \quad \forall \mathbf{x}, \mathbf{y}$$

Theorem. If f is m -strongly convex, $\nabla^2 f$ is M -Lipschitz continuous, and \mathbf{x}^* is a minimum of f , then the sequence $\{\mathbf{x}_k\}$ produced by Newton's method satisfies

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq \frac{M}{2m} \|\mathbf{x}_k - \mathbf{x}^*\|^2$$

Note. The same remarks on slide 14 apply here with $|x_k - x^*|$ replaced by $\|\mathbf{x}_k - \mathbf{x}^*\|$. In particular, if $\|\mathbf{x}_0 - \mathbf{x}^*\| < \frac{2m}{M}$, then

$$\|\mathbf{x}_k - \mathbf{x}^*\| \leq \frac{2m}{M} \left(\frac{M}{2m} \|\mathbf{x}_0 - \mathbf{x}^*\| \right)^{2^k}$$

The proof is also very similar with only minor modifications.

Proof

$$\begin{aligned} & \|\mathbf{x}_{k+1} - \mathbf{x}^*\| \\ &= \|\mathbf{x}_k - \mathbf{x}^* - [\nabla^2 f(\mathbf{x}_k)]^{-1} [\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}^*)]\| \end{aligned} \quad (1)$$

$$\leq \|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| \cdot \|\nabla f(\mathbf{x}^*) - \nabla f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k)(\mathbf{x}^* - \mathbf{x}_k)\| \quad (2)$$

$$= \|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| \cdot \left\| \int_0^1 [\nabla^2 f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k)) - \nabla^2 f(\mathbf{x}_k)](\mathbf{x}^* - \mathbf{x}_k) dt \right\| \quad (3)$$

$$\leq \|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| \int_0^1 \|[\nabla^2 f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k)) - \nabla^2 f(\mathbf{x}_k)](\mathbf{x}^* - \mathbf{x}_k)\| dt \quad (4)$$

$$\leq \|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| \int_0^1 \|\nabla^2 f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k)) - \nabla^2 f(\mathbf{x}_k)\| \cdot \|\mathbf{x}^* - \mathbf{x}_k\| dt \quad (5)$$

$$\leq \|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| \int_0^1 Mt \|\mathbf{x}^* - \mathbf{x}_k\|^2 dt \quad (6)$$

$$= \|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| \cdot \frac{M}{2} \|\mathbf{x}^* - \mathbf{x}_k\|^2 \quad (7)$$

$$\leq \frac{M}{2m} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \quad (8)$$

Proof (cont'd)

1. Step (1) uses the Newton updating rule

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

and the optimality condition $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

2. Step (2) applies the compatibility of norms on slide 17 to

$$[\nabla^2 f(\mathbf{x}_k)]^{-1} [\nabla f(\mathbf{x}^*) - \nabla f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k)(\mathbf{x}^* - \mathbf{x}_k)]$$

3. Step (3) applies the Newton-Leibniz formula to the function $\mathbf{h}(t) = \nabla f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k))$,

$$\nabla f(\mathbf{x}^*) - \nabla f(\mathbf{x}_k) = \mathbf{h}(1) - \mathbf{h}(0) = \int_0^1 \mathbf{h}'(t) dt$$

where $\mathbf{h}'(t)$ is given by the chain rule,

$$\mathbf{h}'(t) = \nabla^2 f(\mathbf{x}_k + t(\mathbf{x}^* - \mathbf{x}_k))(\mathbf{x}^* - \mathbf{x}_k)$$

Proof (cont'd)

4. Step (4) uses the following inequality

$$\left\| \int \mathbf{f}(t) dt \right\| \leq \int \|\mathbf{f}(t)\| dt$$

Proof. Let $\mathbf{z} = \int \mathbf{f}(t) dt$.

$$\|\mathbf{z}\|^2 = \mathbf{z}^T \int \mathbf{f}(t) dt \stackrel{(a)}{=} \int \mathbf{z}^T \mathbf{f}(t) dt \stackrel{(b)}{\leq} \int \|\mathbf{z}\| \cdot \|\mathbf{f}(t)\| dt = \|\mathbf{z}\| \int \|\mathbf{f}(t)\| dt,$$

where (a) uses linearity of integration and (b) Cauchy-Schwarz.

5. Step (5) again applies the compatibility of norms on slide 17
6. Step (6) uses the Lipschitz continuity of $\nabla^2 f$
7. Step (7) performs the integration over t
8. Step (8) uses the m -strong convexity of f

$$\|[\nabla^2 f(\mathbf{x}_k)]^{-1}\| = \lambda_{\max}([\nabla^2 f(\mathbf{x}_k)]^{-1}) = \frac{1}{\lambda_{\min}(\nabla^2 f(\mathbf{x}_k))} \leq \frac{1}{m}$$

Outline

- Newton's method and properties
- Analysis of Newton's method
- Damped Newton's method

Damped Newton's method

The Newton direction $-\left[\nabla^2 f(\mathbf{x})\right]^{-1} \nabla f(\mathbf{x})$ is a descent direction, but with step size 1, Newton's method does not guarantee $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$.

To ensure $f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$, **damped Newton's method** does backtracking line search along the Newton direction.

Damped Newton's method

converge globally

- 1: initialization $\mathbf{x} \leftarrow \mathbf{x}_0 \in \mathbb{R}^n$
- 2: **while** $\|\nabla f(\mathbf{x})\| > \delta$ **do**
- 3: $\mathbf{d} \leftarrow -\left[\nabla^2 f(\mathbf{x})\right]^{-1} \nabla f(\mathbf{x})$ ▷ solve $\nabla^2 f(\mathbf{x}) \mathbf{d} = -\nabla f(\mathbf{x})$
- 4: $t \leftarrow 1$
- 5: **while** $f(\mathbf{x} + t\mathbf{d}) > f(\mathbf{x}) + \alpha t \nabla f(\mathbf{x})^T \mathbf{d}$ **do**
- 6: $t \leftarrow \beta t$
- 7: **end while**
- 8: $\mathbf{x} \leftarrow \mathbf{x} + t\mathbf{d}$
- 9: **end while**
- 10: **return** \mathbf{x}

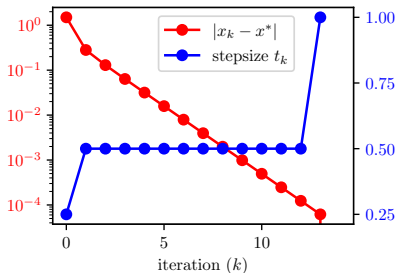
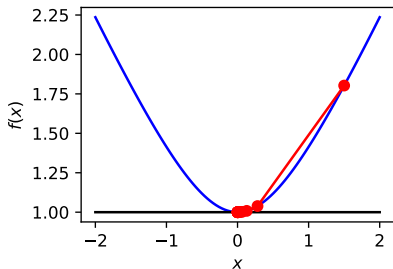
where $\alpha, \beta \in (0, 1)$

Example

$$f(x) = \sqrt{1+x^2}$$

Recall pure Newton's method converges iff $|x_0| < 1$.

Damped Newton's method **converges globally**, e.g. for $x_0 = 1.5$.



一旦到了下一个阶段，就会进入
纯牛顿法，会很快，第一阶段比较慢

这里stepsize变为1了，其实
基本分为two phase

Convergence analysis

Theorem. Assume f is m -strongly convex and L -smooth, $\nabla^2 f$ is M -Lipschitz, and \mathbf{x}^* is a minimum of f . Damped Newton's method satisfies the following error bounds

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \begin{cases} f(\mathbf{x}_0) - f(\mathbf{x}^*) - \gamma k, & \text{if } k \leq k_0 \\ \frac{2m^3}{M^2} \left(\frac{1}{2}\right)^{2^{k-k_0+1}}, & \text{if } k > k_0 \end{cases}$$

where $\gamma = 2\alpha\bar{\alpha}\beta\eta^2m/L^2$, $\eta = \min\{1, 3(1 - 2\alpha)\}m^2/M$, and k_0 is the number of steps until $\|\nabla f(\mathbf{x}_{k_0+1})\| \leq \eta$. 其实这就告诉我们此时离optima还很远

Notes.

- Damped Newton's method guarantees **global** convergence.
- To get $f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \epsilon$, we need at most

$$\text{第一阶段} \quad \frac{f(\mathbf{x}_0) - f(\mathbf{x}^*)}{\gamma} + \log_2 \log_2 \frac{\epsilon_0}{\epsilon} \quad \text{第二阶段}$$

where $\epsilon_0 = \frac{2m^3}{M^2}$. It can be slow if γ is small.

Convergence analysis (cont'd)

Detailed analysis shows that the convergence follows two stages

- **Damped Newton phase.** When $\|\nabla f(\mathbf{x}_k)\| > \eta$, backtracking selects a step size $t_k \leq 1$, and

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leq -\gamma$$

Summing over k from 0 to $k_0 - 1$,

$$f(\mathbf{x}^*) - f(\mathbf{x}_0) \leq f(\mathbf{x}_{k_0}) - f(\mathbf{x}_0) \leq -k_0\gamma \implies k_0 \leq \frac{f(\mathbf{x}_0) - f(\mathbf{x}^*)}{\gamma}$$

- **Pure Newton phase.** When $\|\nabla f(\mathbf{x}_k)\| \leq \eta$, backtracking selects step size $t_k = 1$, and

$$\|\nabla f(\mathbf{x}_{k+1})\| \leq \frac{M}{2m^2} \|\nabla f(\mathbf{x}_k)\|^2 \leq \frac{1}{2} \|\nabla f(\mathbf{x}_k)\| \leq \eta$$

By induction, $\|\nabla f(\mathbf{x}_k)\| \leq \eta$ for all $k \geq k_0$, so we will stay in the pure Newton phase with $t_k = 1$.