

Chapter 9 Newton's Method

1. Introduction
2. Analysis of Newton's Method
3. Levenberg-Marquardt Modification
4. Newton's Method for Nonlinear Least Squares

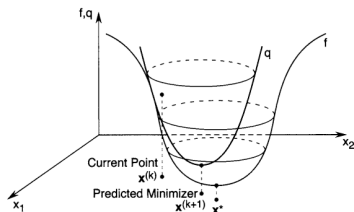


Newton's Method

motivation

- Steepest descent method only uses **1st-order derivative (gradient)** as searching direction.
- If **higher order derivative** is applied to recursion, it may perform better than SD method.

- 1 **Construct a quadratic function** with the same 1st- and 2nd-order derivatives as the objective function at this point, which can be used as an approximation of the objective function.
- 2 **Find the minimizer** of quadratic function as the starting point of the next iteration.
- 3 **Repeat the above procedure** to derive the minimizer of the objective function.



Newton's Method

let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice continuous differentiable

The Taylor expansion of f at $\mathbf{x}^{(k)}$ is

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(k)}) + (\mathbf{x} - \mathbf{x}^{(k)})^\top \mathbf{g}^{(k)} + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(k)})^\top \mathbf{F}(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)}) =: q(\mathbf{x}),$$

where $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$ and $\mathbf{F}(\mathbf{x}^{(k)}) = \nabla^2 f(\mathbf{x}^{(k)})$.

By the optimality conditions, the local minima of $q(\mathbf{x})$ satisfies

$$\mathbf{0} = \nabla q(\mathbf{x}) = \mathbf{g}^{(k)} + \mathbf{F}(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)}).$$

If $\mathbf{F}(\mathbf{x}^{(k)}) \succ 0$, the minimum of q is

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)} \implies \text{Newton's method}$$



Newton's Method

Example (minimize f by Newton's method)

$f(x_1, x_2, x_3, x_4) = (x_1 + 10x_2)^2 + 5(x_3 - x_4)^2 + (x_2 - 2x_3)^4 + 10(x_1 - x_4)^4$.
The initial point is $\mathbf{x}^{(0)} = [3, -1, 0, 1]^\top$. (Ans: see the textbook)

Example (minimize $f(\mathbf{x}) = x_1^2 + \gamma x_2^2$ by Newton's method)

Ans: (see the blackboard), $\mathbf{x}^* = (0, 0)$.

★ Newton's method can be divided into two steps

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)} \iff \begin{cases} \mathbf{d}^{(k)} = -\mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)} \\ \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{d}^{(k)} \end{cases}$$

- Solve $\mathbf{F}(\mathbf{x}^{(k)}) \mathbf{d}^{(k)} = -\mathbf{g}^{(k)}$ to obtain $\mathbf{d}^{(k)}$, which is an n -dimensional linear equations. $\mathbf{d}^{(k)}$ is called Newton direction.
- Determine the next iteration point $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{d}^{(k)}$.



Newton's Method

Require: $\mathbf{x}^{(0)} \in \mathbb{R}^n$ and $\mathbf{H}_0 \succ 0$. set $k := 0$;

1: **repeat**

2: solve linear equations $\mathbf{F}(\mathbf{x}^{(k)})\mathbf{d}^{(k)} = -\mathbf{g}^{(k)}$ for $\mathbf{d}^{(k)} \implies$ **Newton direction**

3: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{d}^{(k)} \implies$ **updating iterative point.**

4: **until** $\|\mathbf{g}^{(k)}\| \leq \epsilon$.

Example (apply Newton's method to minimize f)

$f(\mathbf{x}) = (x_1 - 2)^4 + (x_1 - 2x_2)^2$ with starting point $\mathbf{x}^{(0)} = [0, 3]^\top$.

Ans: S1. $\mathbf{g}^{(0)} = \begin{bmatrix} -44 \\ 24 \end{bmatrix}$ and $\mathbf{F}(\mathbf{x}^{(0)}) = \begin{bmatrix} 50 & -4 \\ -4 & 8 \end{bmatrix}$,

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - [\mathbf{F}(\mathbf{x}^{(0)})]^{-1}\mathbf{g}^{(0)} = \begin{bmatrix} 0 \\ 3 \end{bmatrix} + \begin{bmatrix} 50 & -4 \\ -4 & 8 \end{bmatrix}^{-1} \begin{bmatrix} -44 \\ 24 \end{bmatrix} = \begin{bmatrix} 0.67 \\ 0.33 \end{bmatrix}$$

S2.

k	1	2	3	4	5	6	7
$\mathbf{x}^{(k)}$	[0, 3]	[0.67, 0.33]	[1.11, 0.56]	[1.41, 0.70]	[1.61, 0.80]	[1.74, 0.87]	[1.83, 0.91]
$f(\mathbf{x}^{(k)})$	52	3.13	0.63	0.12	0.02	0.005	0.0009

Newton's Method

Example (apply Newton's method to minimize f)

$$f(\mathbf{x}) = (x_1 - 2)^4 + (x_1 - 2)^2 x_2^2 + (x_2 + 1)^2 \text{ with } \mathbf{x}^{(0)} = [1, 1]^T.$$

			$f(\mathbf{x}_k)$
$\mathbf{x}_0 = (1.0$	$, 1.0$	$)^T$	6.0
$\mathbf{x}_1 = (1.0$	$, -0.5$	$)^T$	1.5
$\mathbf{x}_2 = (1.3913043$	$, -0.69565217)$	$)^T$	4.09×10^{-1}
$\mathbf{x}_3 = (1.7459441$	$, -0.94879809)$	$)^T$	6.49×10^{-2}
$\mathbf{x}_4 = (1.9862783$	$, -1.0482081)$	$)^T$	2.53×10^{-3}
$\mathbf{x}_5 = (1.9987342$	$, -1.0001700)$	$)^T$	1.63×10^{-6}
$\mathbf{x}_6 = (1.9999996$	$, -1.0000016)$	$)^T$	2.75×10^{-12}

note: exact solution is
 $\mathbf{x}^* = [2, -1]^T$.

analysis of Newton's method

- If the Hessian matrix $\mathbf{F}(\mathbf{x}^{(k)}) \neq 0$, Newton direction $\mathbf{d}^{(k)} = -\mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}$ is not necessarily descent direction.
- If the objective function f is quadratic, Newton's method converges to optima in one iteration for any initial point $\mathbf{x}^{(0)}$, i.e., **the convergence rate of Newton's method on minimizing quadratic function is ∞ .**

Newton's Method

Newton's method for nonlinear equations

$$\mathbf{g}(\mathbf{x}) = \begin{bmatrix} g_1(\mathbf{x}) \\ g_2(\mathbf{x}) \\ \vdots \\ g_n(\mathbf{x}) \end{bmatrix} = 0,$$

where $\mathbf{x} \in \mathbb{R}^n$, $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuous differentiable.
Let $\mathbf{F}(\mathbf{x}) \in \mathbb{R}^{n \times n}$ be the Jacobian of \mathbf{g} at \mathbf{x} .
The (i, j) th element of $\mathbf{F}(\mathbf{x})$ is $\frac{\partial g_i}{\partial x_j} \mathbf{x}$.
iterative scheme: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}(\mathbf{x}^{(k)})$

★ if $n = 1$, iterative scheme reduces to: $x^{(k+1)} = x^{(k)} - \frac{g(x^{(k)})}{g'(x^{(k)})}$.

Example (apply Newton's method for nonlinear equation)

$g(x) = e^x - 1$ with starting point $x^{(0)} = -1$.

Ans:

k	0	1	2	3	4	5
$x^{(k)}$	-1.00000	0.71828	0.20587	0.01981	0.00019	0.00000
$g(x^{(k)})$	-0.63212	1.05091	0.22859	0.02000	0.00019	0.00000

★ Typically, 5-10 iterations can converge. Unfortunately, the classical Newton method may also diverge, and often does!



Convergence of Newton's Method

Theorem

If $f \in \mathcal{C}^3$, $\mathbf{x}^* \in \mathbb{R}^n$ satisfies $\nabla f(\mathbf{x}^*) = \mathbf{0}$, and $\mathbf{F}(\mathbf{x}^*)$ is invertible. Then, for all $\mathbf{x}^{(0)}$ sufficiently close to \mathbf{x}^* , Newton's method is well-defined for all k and converges to \mathbf{x}^* with an order of at least 2.

proof. Taylor expansion of gradient ∇f at $\mathbf{x}^{(0)}$:

$$\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^{(0)}) - \mathbf{F}(\mathbf{x}^{(0)})(\mathbf{x} - \mathbf{x}^{(0)}) = O(\|\mathbf{x} - \mathbf{x}^{(0)}\|^2).$$

$\because f \in \mathcal{C}^3$ and $\mathbf{F}(\mathbf{x}^*)$ is invertible.

$\therefore \exists \varepsilon > 0$, $c_1 > 0$, and $c_2 > 0$ such that: for all $\mathbf{x}^{(0)} \in \mathcal{B}(\mathbf{x}^*, \varepsilon)$,
 $\mathbf{x} \in \mathcal{B}(\mathbf{x}^*, \varepsilon)$,

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^{(0)}) - \mathbf{F}(\mathbf{x}^{(0)})(\mathbf{x} - \mathbf{x}^{(0)})\| \leq c_1 \|\mathbf{x} - \mathbf{x}^{(0)}\|^2. \quad (1)$$

$\because \mathbf{F}(\mathbf{x})$ is invertible and $\|\mathbf{F}(\mathbf{x})^{-1}\| \leq c_2$, by replacing \mathbf{x} in (1) with \mathbf{x}^* and using $\nabla f(\mathbf{x}^*) = \mathbf{0}$, we have

$$\|\mathbf{F}(\mathbf{x}^{(0)})(\mathbf{x}^{(0)} - \mathbf{x}^*) - \nabla f(\mathbf{x}^{(0)})\| \leq c_1 \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2.$$



Convergence of Newton's Method

$$\begin{aligned}\therefore \|\mathbf{x}^{(1)} - \mathbf{x}^*\| &\stackrel{\text{why?}}{=} \|\mathbf{x}^{(0)} - \mathbf{x}^* - \mathbf{F}(\mathbf{x}^{(0)})^{-1} \nabla f(\mathbf{x}^{(0)})\| \\ &= \|\mathbf{F}(\mathbf{x}^{(0)})^{-1} [\mathbf{F}(\mathbf{x}^{(0)}) (\mathbf{x}^{(0)} - \mathbf{x}^*) - \nabla f(\mathbf{x}^{(0)})]\| \\ &\leq \|\mathbf{F}(\mathbf{x}^{(0)})^{-1}\| \cdot \|\mathbf{F}(\mathbf{x}^{(0)}) (\mathbf{x}^{(0)} - \mathbf{x}^*) - \nabla f(\mathbf{x}^{(0)})\| \\ &\leq \|\mathbf{F}(\mathbf{x}^{(0)})^{-1}\| \cdot c_1 \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2 \\ &\leq c_1 c_2 \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2.\end{aligned}$$

By choosing an $\mathbf{x}^{(0)}$ satisfying $\|\mathbf{x}^{(0)} - \mathbf{x}^*\| \leq \frac{\alpha}{c_1 c_2}$ with $\alpha \in (0, 1)$, we have

$$\|\mathbf{x}^{(1)} - \mathbf{x}^*\| \leq \alpha \|\mathbf{x}^{(0)} - \mathbf{x}^*\|.$$

$\therefore \lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\| = 0 \implies$ the sequence $\{\mathbf{x}^{(k)}\}$ converges to \mathbf{x}^* .

$\therefore \|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq c_1 c_2 \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2$, the order of convergence is at least 2.

★ If the initial point is close to the minima (maxima), Newton's method will have a good convergence. However, if it is far from the minima (maxima), Newton's method is not necessarily convergent.



Convergence of Newton's Method

Theorem (descent property of Newton's method)

Let $\mathbf{x}^{(k)}$ be the sequence generated by Newton's method for minimizing f . If $\mathbf{F}(\mathbf{x}^{(k)}) \succ \mathbf{0}$, and $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$, then the Newton direction $\mathbf{d}^{(k)} = -\mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)} = \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}$ is a descent direction in the sense that there exists a $\bar{\alpha} > 0$ such that $\alpha \in (0, \bar{\alpha})$, $f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) < f(\mathbf{x}^{(k)})$.

proof. let $\phi(\alpha) = f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) \implies \phi'(\alpha) = \nabla f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})^\top \mathbf{d}^{(k)}$.

$$\because \mathbf{F}(\mathbf{x}^{(k)})^{-1} \succ \mathbf{0} \text{ and } \mathbf{g}^{(k)} \neq \mathbf{0},$$

$$\because \phi'(0) = \nabla f(\mathbf{x}^{(k)})^\top \mathbf{d}^{(k)} = -\mathbf{g}^{(k)\top} \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)} < 0.$$

$$\because \exists \bar{\alpha} > 0 \text{ such that } \phi(\alpha) < \phi(0) \text{ for all } \alpha \in (0, \bar{\alpha}),$$

$$\because f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) < f(\mathbf{x}^{(k)}) \text{ for all } \alpha \in (0, \bar{\alpha}).$$

performance of Newton's method

- very fast when it converges (how fast? quadratic).
- may diverge (or worse, $\mathbf{F}(\mathbf{x}^{(k)})$ is nonsingular) when $\mathbf{x}^{(0)}$ is far from a nonsingular local optimum.

Question: how to modify Newton's method so that it converges globally?

Modification of Newton's Method

Newton's method with stepsize

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}, \text{ where } \alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}).$$

★ Modified Newton method has the descent property $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$.

Levenberg-Marquardt (L-M) modification

- If the Hessian $\mathbf{F}(\mathbf{x}^{(k)})$ is not necessarily positive definite, then $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k [\mathbf{F}(\mathbf{x}^{(k)}) + \mu_k \mathbf{I}]^{-1} \mathbf{g}^{(k)}$, where $\mu_k \geq 0$.
- If μ_k is large enough, e.g., $\mu_k > -\lambda_{\min}(\mathbf{F})$, then $\mathbf{G} = \mathbf{F} + \mu \mathbf{I} \succ 0$. Thus, the searching direction $\mathbf{d}^{(k)} = -(\mathbf{F}(\mathbf{x}^{(k)}) + \mu_k \mathbf{I})^{-1} \mathbf{g}^{(k)}$ is a descent direction.

Behavior of Parameter μ_k

- ① if $\mu_k \rightarrow 0$, L-M modification can be made to approach the pure Newton's method.
- ② if $\mu_k \rightarrow \infty$, L-M modification approaches the pure gradient descent method with a small step size.



Gauss-Newton Newton's Method

nonlinear least squares (NLS): e.g., applications to curve-fitting (see textbook)

$\min \sum_{i=1}^m [r_i(\mathbf{x})]^2$, where $r_i : \mathbb{R}^n \rightarrow \mathbb{R}$ ($i = 1, \dots, m$) are differentiable functions.

application of Newton's method to NLS

let $\mathbf{r} = [r_1, \dots, r_m]^\top$ be vector-valued function from $\mathbb{R}^n \rightarrow \mathbb{R}^m$. The objective of NLS amounts to $f(\mathbf{x}) = \mathbf{r}(\mathbf{x})^\top \mathbf{r}(\mathbf{x})$.

- $\nabla f(\mathbf{x}) = 2\mathbf{J}(\mathbf{x})^\top \mathbf{r}(\mathbf{x})$ with $\mathbf{J}(\mathbf{x}) = \begin{bmatrix} \frac{\partial r_1}{\partial x_1}(\mathbf{x}) & \dots & \frac{\partial r_1}{\partial x_n}(\mathbf{x}) \\ \vdots & & \vdots \\ \frac{\partial r_m}{\partial x_1}(\mathbf{x}) & \dots & \frac{\partial r_m}{\partial x_n}(\mathbf{x}) \end{bmatrix}$ as Jacobian of \mathbf{r} .
- $\nabla f(\mathbf{x})$ can also be rewritten by $(\nabla f(\mathbf{x}))_j = \frac{\partial f}{\partial x_j}(\mathbf{x}) = 2 \sum_{i=1}^m r_i(\mathbf{x}) \frac{\partial r_i}{\partial x_j}(\mathbf{x})$.
- Hessian of f is $\mathbf{F}(\mathbf{x}) = 2[\mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}) + \mathbf{S}(\mathbf{x})]$, where $\mathbf{S}(\mathbf{x})$ is a matrix whose (h, j) th element is $\sum_{i=1}^m r_i(\mathbf{x}) \frac{\partial^2 r_i}{\partial x_h \partial x_j}(\mathbf{x})$.

Gauss-Newton Newton's Method

recursion of Newton's method for solving NLS

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - [\mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}) + \mathbf{S}(\mathbf{x})]^{-1} \mathbf{J}(\mathbf{x})^\top \mathbf{r}(\mathbf{x}).$$

- ★ As the matrix $\mathbf{S}(\mathbf{x})$ usually contains the 2nd derivative of \mathbf{r} , where the elements are all small.

Gauss-Newton method for NLS, i.e., neglecting $\mathbf{S}(\mathbf{x})$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - [\mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x})]^{-1} \mathbf{J}(\mathbf{x})^\top \mathbf{r}(\mathbf{x}).$$

- ★ A pitfall of Gauss-Newton method is $\mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x})$ may not be positive definite.
- ★ By using Levenberg-Marquardt modification, Gauss-Newton method can be modified as

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - [\mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}) + \mu_k \mathbf{I}]^{-1} \mathbf{J}(\mathbf{x})^\top \mathbf{r}(\mathbf{x}).$$



Homework

Exercise in the textbook: 9.4.

