

4 NEWTON'S METHOD

4.1 INTRODUCTION

- The method of steepest descent uses **only the first derivatives** (gradients) in selecting a suitable search direction. This strategy is not always the most effective. If higher derivatives are used, the resulting iterative algorithm may perform better than the steepest descent method.
- Newton's method (sometimes called the Newton-Raphson method) uses the **first** and the **second derivatives** and indeed does perform better than the steepest descent method if the initial point is close to the minimizer.
- We can obtain a quadratic approximation to the twice continuously differentiable objection function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ using the Taylor series expansion of f about the current point $\mathbf{x}^{(k)}$, neglecting terms of order three and higher:

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

where, for simplicity, we use the notation $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$. Applying the FONC to q yields

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

- If $\mathbf{F}(\mathbf{x}^{(k)}) > 0$, then q achieves a minimum at

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}.$$

This recursive formula represents **Newton's method**.

Example 4.1.1. Use Newton's method to minimize the Powell function:

$$f(x_1, x_2, x_3, x_4) = (x_1 + 10x_2)^2 + 5(x_3 - x_4)^2 + (x_2 - 2x_3)^4 + 10(x_1 - x_4)^4.$$

Use as the starting point $\mathbf{x}^{(0)} = [3, -1, 0, 1]^\top$. Perform three iterations.

Note that $f(\mathbf{x}^{(0)}) = 215$. We have

$$\nabla f(\mathbf{x}) = \begin{bmatrix} 2(x_1 + 10x_2) + 40(x_1 - x_4)^3 \\ 20(x_1 + 10x_2) + 4(x_2 - 2x_3)^3 \\ 10(x_3 - x_4) - 8(x_2 - 2x_3)^3 \\ -10(x_3 - x_4) - 40(x_1 - x_4)^3 \end{bmatrix},$$

and $\mathbf{F}(\mathbf{x})$ is given by

$$\begin{bmatrix} 2 + 120(x_1 - x_4)^2 & 20 & 0 & -120(x_1 - x_4)^2 \\ 20 & 200 + 12(x_2 - 2x_3)^2 & -24(x_2 - 2x_3)^2 & 0 \\ 0 & -24(x_2 - 2x_3)^2 & 10 + 48(x_2 - 2x_3)^2 & -10 \\ -120(x_1 - x_4)^2 & 0 & -10 & 10 + 120(x_1 - x_4)^2 \end{bmatrix}.$$

Iteration 1.

$$\mathbf{g}^{(0)} = [306, -144, -2, -310]^\top,$$

$$\mathbf{F}(\mathbf{x}^{(0)}) = \begin{bmatrix} 482 & 20 & 0 & -480 \\ 20 & 212 & -24 & 0 \\ 0 & -24 & 58 & -10 \\ -480 & 0 & -10 & 490 \end{bmatrix},$$

$$\mathbf{F}(\mathbf{x}^{(0)})^{-1} = \begin{bmatrix} 0.1126 & -0.0089 & 0.0154 & 0.1106 \\ -0.0089 & 0.0057 & 0.0008 & -0.0087 \\ 0.0154 & 0.0008 & 0.0203 & 0.0155 \\ 0.1106 & -0.0087 & 0.0155 & 0.1107 \end{bmatrix},$$

$$\mathbf{F}(\mathbf{x}^{(0)})^{-1} \mathbf{g}^{(0)} = [1.4127, -0.8413, -0.2540, 0.7460]^\top.$$

Hence,

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} - \mathbf{F}(\mathbf{x}^{(0)})^{-1} \mathbf{g}^{(0)} = [-1.5873, -0.1587, 0.2540, 0.2540]^\top,$$

$$f(\mathbf{x}^{(1)}) = 31.8.$$

Iteration 2.

$$\mathbf{g}^{(1)} = [94.81, -1.179, 2.371, -94.81]^\top,$$

$$\mathbf{F}(\mathbf{x}^{(1)}) = \begin{bmatrix} 215.3 & 20 & 0 & -213.3 \\ 20 & 205.3 & -10.67 & 0 \\ 0 & -10.67 & 31.34 & -10 \\ -213.3 & 0 & -10 & 223.3 \end{bmatrix},$$

$$\mathbf{F}(\mathbf{x}^{(1)})^{-1} \mathbf{g}^{(1)} = [0.5291, -0.0529, 0.0846, 0.0846]^\top.$$

Hence,

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} - \mathbf{F}(\mathbf{x}^{(1)})^{-1} \mathbf{g}^{(1)} = [1.0582, -0.1058, 0.1694, 0.1694]^\top,$$

$$f(\mathbf{x}^{(2)}) = 6.28.$$

Iteration 3.

$$\mathbf{g}^{(2)} = [28.09, -0.3475, 0.7031, -28.08]^\top,$$

$$\mathbf{F}(\mathbf{x}^{(2)}) = \begin{bmatrix} 96.80 & 20 & 0 & -94.80 \\ 20 & 202.4 & -4.744 & 0 \\ 0 & -4.744 & 19.49 & -10 \\ -94.80 & 0 & -10 & 104.80 \end{bmatrix},$$

$$\mathbf{x}^{(3)} = \mathbf{x}^{(2)} - \mathbf{F}(\mathbf{x}^{(2)})^{-1} \mathbf{g}^{(2)} = [0.7037, -0.0704, 0.1121, 0.1111]^T,$$

$$f(\mathbf{x}^{(3)}) = 1.24$$

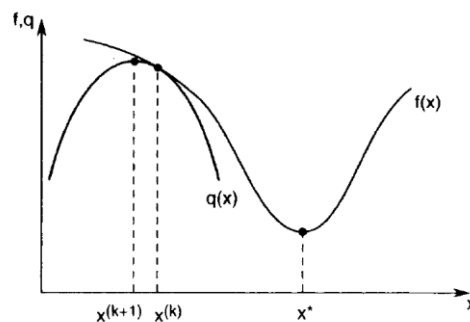
- Observe that the k th iteration of Newton's method can be written in two steps as
 1. Solve $\mathbf{F}(\mathbf{x}^{(k)})\mathbf{d}^{(k)} = -\mathbf{g}^{(k)}$ for $\mathbf{d}^{(k)}$.
 2. Set $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{d}^{(k)}$.
- As in the one-variable case, Newton's method can also be viewed as a technique for iteratively solving the equation

$$\mathbf{g}(\mathbf{x}) = \mathbf{0},$$

where $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{g}: \mathbb{R}^n \rightarrow \mathbb{R}^n$.

4.2 ANALYSIS OF NEWTON'S METHOD

- As in the one-variable case there is no guarantee that Newton's algorithm heads in the direction of decreasing values of the objective function if $\mathbf{F}(\mathbf{x}^{(k)})$ is not positive definite.



- Moreover, even if $\mathbf{F}(\mathbf{x}^{(k)}) > 0$, Newton's method may not be a descent method; that is, it is possible that $f(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k)})$. For example, this may occur if our starting point $\mathbf{x}^{(0)}$ is far away from the solution.
- **Theorem 4.2.1** Suppose that $f \in C^3$ and $\mathbf{x}^* \in \mathbb{R}^n$ is a point such that $\nabla f(\mathbf{x}^*) = \mathbf{0}$ 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 convergence at least 2.

- **Warning:** In the Theorem 4.2.1, we did not state that \mathbf{x}^* is a local minimizer. For example, if \mathbf{x}^* is a local maximizer, then provided that $f \in C^3$ and $\mathbf{F}(\mathbf{x}^*)$ is invertible, Newton's method converges to \mathbf{x}^* if we start close enough to it.

Example 4.2.1. Suppose f is a quadratic function such that

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{x}^\top \mathbf{b},$$

where $\mathbf{Q} = \mathbf{Q}^\top$ is invertible. Show that Newton's method reaches the point \mathbf{x}^* such that $\nabla f(\mathbf{x}^*) = 0$ in just one step starting from any initial point $\mathbf{x}^{(0)}$. Determine the order of convergence.

Note that $\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x}) = \mathbf{Q}\mathbf{x} - \mathbf{b}$ and $\mathbf{F}(\mathbf{x}) = \mathbf{Q}$.

Hence, given any initial point $\mathbf{x}^{(0)}$, by Newton's algorithm

$$\begin{aligned} \mathbf{x}^{(1)} &= \mathbf{x}^{(0)} - \mathbf{F}(\mathbf{x}^{(0)})^{-1} \mathbf{g}^{(0)} \\ &= \mathbf{x}^{(0)} - \mathbf{Q}^{-1}(\mathbf{Q}\mathbf{x}^{(0)} - \mathbf{b}) \\ &= \mathbf{Q}^{-1} \mathbf{b} \\ &= \mathbf{x}^* \end{aligned}$$

In this case, if for all $p > 0$,

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^p} = 0.$$

Therefore, the order of convergence is ∞ .

- The Newton's method may not be a descent method; that is, it is possible that $f(\mathbf{x}^{(k+1)}) \geq f(\mathbf{x}^{(k)})$. Fortunately, it is possible to modify the algorithm such that the descent property holds.
- **Theorem 4.2.2** Let $\{\mathbf{x}^{(k)}\}$ be the sequence generated by Newton's method for minimizing a given objective function $f(\mathbf{x})$. If the Hessian $\mathbf{F}(\mathbf{x}^{(k)}) > 0$ and $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$, then the search direction

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

from $\mathbf{x}^{(k)}$ to $\mathbf{x}^{(k+1)}$ is a descent direction for f in the sense that there exists an $\bar{\alpha} > 0$ such that for all $\alpha \in (0, \bar{\alpha})$,

$$f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) < f(\mathbf{x}^{(k)}).$$

Proof Let

$$\phi(\alpha) = f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}).$$

Then, using the chain rule, we obtain

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

Hence,

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

because $\mathbf{F}(\mathbf{x}^{(k)})^{-1} > 0$ and $\mathbf{g}^{(k)} \neq \mathbf{0}$. Thus, there exists an $\bar{\alpha} > 0$ so that for all $\alpha \in (0, \bar{\alpha})$, $\phi(\alpha) < \phi(0)$. This implies that for all $\alpha \in (0, \bar{\alpha})$,

$$f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) < f(\mathbf{x}^{(k)}).$$

- Theorem 4.2.2 motivates the following modification of Newton's method:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)},$$

where

$$\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)});$$

that is, at each iteration, we perform a line search in the direction $-\mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}$. By Theorem 4.2.2 we conclude that the modified Newton's method has the descent property; that is,

$$f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$$

whenever $\mathbf{g}^{(k)} \neq \mathbf{0}$.

- A drawback of Newton's method is that evaluation of $\mathbf{F}(\mathbf{x}^{(k)})$ for large n can be computationally expensive. Furthermore, we have to solve the set of n linear equations $\mathbf{F}(\mathbf{x}^{(k)}) \mathbf{d}^{(k)} = -\mathbf{g}^{(k)}$.
- Another source of potential problems in Newton's method arises from the Hessian matrix not being positive definite.

4.3 LEVENBERG-MARQUARDT MODIFICATION

- If the Hessian matrix $\mathbf{F}(\mathbf{x}^{(k)})$ is not positive definite, then the search direction $\mathbf{d}^{(k)} = -\mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}$ may not point in a descent direction. A simple technique to ensure that the search direction is a descent direction is to introduce the **Levenberg-Marquardt modification** of Newton's algorithm:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k (\mathbf{F}(\mathbf{x}^{(k)}) + \mu_k \mathbf{I})^{-1} \mathbf{g}^{(k)},$$

where $\mu_k \geq 0$.

- The idea underlying the Levenberg-Marquardt modification is as follows. Consider a symmetric matrix \mathbf{F} , which may not be positive definite. Let $\lambda_1, \dots, \lambda_n$

be the eigenvalues of \mathbf{F} with corresponding eigenvectors $\mathbf{v}_1, \dots, \mathbf{v}_n$. The eigenvalues $\lambda_1, \dots, \lambda_n$ are real, but may not all be positive. Next, consider the matrix $\mathbf{G} = \mathbf{F} + \mu \mathbf{I}$, where $\mu \geq 0$. Note that the eigenvalues of \mathbf{G} are $\lambda_1 + \mu, \dots, \lambda_n + \mu$. Indeed,

$$\mathbf{G}\mathbf{v}_i = (\mathbf{F} + \mu \mathbf{I})\mathbf{v}_i = \mathbf{F}\mathbf{v}_i + \mu \mathbf{I}\mathbf{v}_i = \lambda_i \mathbf{v}_i + \mu \mathbf{v}_i = (\lambda_i + \mu)\mathbf{v}_i,$$

which shows that for all $i = 1, \dots, n$, \mathbf{v}_i is also an eigenvector of \mathbf{G} with eigenvalue $\lambda_i + \mu$. Therefore, if μ is sufficiently large, then all the eigenvalues of \mathbf{G} are positive and \mathbf{G} is positive definite.

- The Levenberg-Marquardt modification of Newton's algorithm can be made to approach the behavior of the pure Newton's method by letting $\mu_k \rightarrow 0$. On the other hand, by letting $\mu_k \rightarrow \infty$, the algorithm approaches a pure gradient method with small step size.

4.4 NEWTON'S METHOD FOR NONLINEAR LEAST SQUARES

- We now examine a particular class of optimization problems and the use of Newton's method for solving them. Consider the following problem:

$$\text{minimize } \sum_{i=1}^m (r_i(\mathbf{x}))^2,$$

where $r_i: \mathbb{R}^n \rightarrow \mathbb{R}$, $i = 1, \dots, m$, are given functions. This particular problem is called a **nonlinear least-squares problem**.

Example 4.4.1 Suppose that we are given m measurements of a process at m points in time (here, $m = 21$). Let t_1, \dots, t_m denote the measurement times and y_1, \dots, y_m the measurement values. Note that $t_1 = 0$ while $t_{21} = 10$. We wish to fit a sinusoid to the measurement data. The equation of the sinusoid is

$$y = A \sin(\omega t + \phi)$$

with appropriate choices of the parameters A , ω , and ϕ . Formulate the data-fitting problem and derive the Newton's method.

To formulate the data-fitting problem, we construct the objective function

$$\sum_{i=1}^m (y_i - A \sin(\omega t_i + \phi))^2,$$

representing the sum of the squared errors between the measurement values and the function values at the corresponding points in time. Let $\mathbf{x} = [A, \omega, \phi]^T$ represent the vector of decision variables. We therefore obtain a nonlinear least-squares problem with

$$r_i(\mathbf{x}) = y_i - A \sin(\omega t_i + \phi).$$

Defining $\mathbf{r} = [r_1, \dots, r_m]^\top$, we write the objective function as $f(\mathbf{x}) = \mathbf{r}(\mathbf{x})^\top \mathbf{r}(\mathbf{x})$. To apply Newton's method, we need to compute the gradient and the Hessian of f . The j th component of $\nabla f(\mathbf{x})$ is

$$(\nabla f(\mathbf{x}))_j = \frac{\partial f}{\partial x_j}(\mathbf{x}) = 2 \sum_{i=1}^m r_i(\mathbf{x}) \frac{\partial r_i}{\partial x_j}(\mathbf{x}).$$

Then, the gradient of f can be represented as

$$\nabla f(\mathbf{x}) = 2\mathbf{J}(\mathbf{x})^\top \mathbf{r}(\mathbf{x})$$

where

$$\mathbf{J}(\mathbf{x}) = \begin{bmatrix} \frac{\partial r_1}{\partial x_1}(\mathbf{x}) & \frac{\partial r_1}{\partial x_2}(\mathbf{x}) & \frac{\partial r_1}{\partial x_3}(\mathbf{x}) \\ \vdots & \vdots & \vdots \\ \frac{\partial r_m}{\partial x_1}(\mathbf{x}) & \frac{\partial r_m}{\partial x_2}(\mathbf{x}) & \frac{\partial r_m}{\partial x_3}(\mathbf{x}) \end{bmatrix}.$$

Next, we compute the Hessian matrix of f . The (k, j) th component of the Hessian is given by

$$\begin{aligned} \frac{\partial^2 f}{\partial x_k \partial x_j}(\mathbf{x}) &= \frac{\partial}{\partial x_k} \left(\frac{\partial f}{\partial x_j}(\mathbf{x}) \right) \\ &= \frac{\partial}{\partial x_k} \left(2 \sum_{i=1}^m r_i(\mathbf{x}) \frac{\partial r_i}{\partial x_j}(\mathbf{x}) \right) \\ &= 2 \sum_{i=1}^m \left(\frac{\partial r_i}{\partial x_k}(\mathbf{x}) \frac{\partial r_i}{\partial x_j}(\mathbf{x}) + r_i(\mathbf{x}) \frac{\partial^2 r_i}{\partial x_k \partial x_j}(\mathbf{x}) \right) \end{aligned}$$

Letting $\mathbf{S}(\mathbf{x})$ be the matrix whose (k, j) th component is

$$\sum_{i=1}^m r_i(\mathbf{x}) \frac{\partial^2 r_i}{\partial x_k \partial x_j}(\mathbf{x}),$$

we write the Hessian matrix as

$$\mathbf{F}(\mathbf{x}) = 2(\mathbf{J}(\mathbf{x})^\top \mathbf{J}(\mathbf{x}) + \mathbf{S}(\mathbf{x})).$$

Therefore, Newton's method applied to the nonlinear least-squares problem is given by

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