

Chapter 9 Newton's Method

An Introduction to Optimization

Spring, 2014

Wei-Ta Chu

Introduction

- ▶ The steepest descent method uses only first derivatives in selecting a suitable search direction.
- ▶ Newton's method (sometimes called Newton-Raphson method) uses first and second derivatives and indeed performs better.
- ▶ Given a starting point, construct a quadratic approximation to the objective function that matches the first and second derivative values at that point. We then minimize the approximate (quadratic function) instead of the original objective function. The minimizer of the approximate function is used as the starting point in the next step and repeat the procedure iteratively.

$$f(b) = f(a) + \frac{h}{1!}f^{(1)}(a) + \frac{h^2}{2!}f^{(2)}(a) + \dots + \frac{h^{m-1}}{(m-1)!}f^{(m-1)}(a) + R_m$$

Introduction

- ▶ We can obtain a quadratic approximation to the twice continuously differentiable function $f : R^n \rightarrow R$ using the Taylor series expansion of f about the current point $\mathbf{x}^{(k)}$, neglecting terms of order three and higher.

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(k)}) + (\mathbf{x} - \mathbf{x}^{(k)})^T \mathbf{g}^{(k)} + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(k)})^T F(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)}) \triangleq q(\mathbf{x})$$

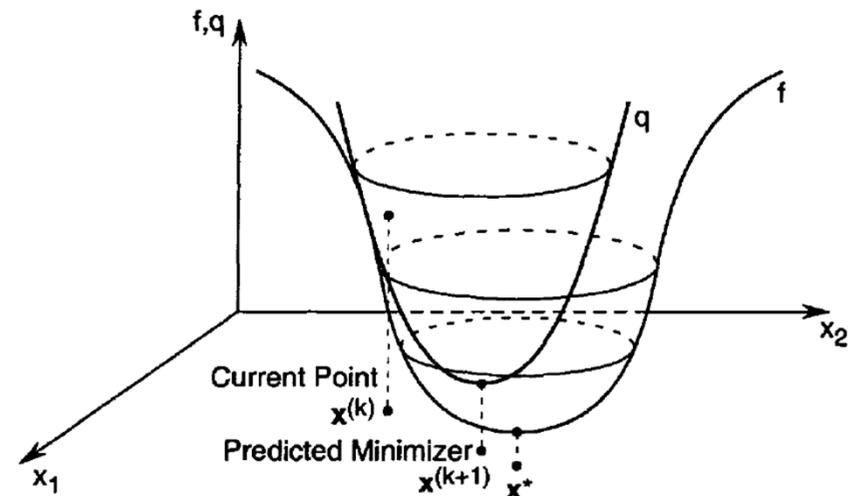
Where, for simplicity, we use the notation $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$

- ▶ Applying the FONC to q yields
- $$\mathbf{0} = \nabla q(\mathbf{x}) = \mathbf{g}^{(k)} + F(\mathbf{x}^{(k)})(\mathbf{x} - \mathbf{x}^{(k)})$$
- ▶ If $F(\mathbf{x}^{(k)}) > 0$, then q achieves a minimum at

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

$$0 = q'(x) = f'(x^{(k)}) + f''(x^{(k)})(x - x^{(k)})$$

$$x^{(k+1)} = x^{(k)} - \frac{f'(x^{(k)})}{f''(x^{(k)})}$$



Example

- ▶ 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 $\mathbf{x}^{(0)} = [3, -1, 0, 1]^T$. Perform three iterations.

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

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

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

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

Example

► Iteration 1.

$$\mathbf{g}^{(0)} = [306, -114, -2, -310]^T$$

$$\mathbf{F}(\mathbf{x}^{(0)}) = \begin{bmatrix} 482 & 20 & 0 & -480 \\ 20 & 212 & -24 & 0 \\ 0 & -24 & 58 & -10 \\ -480 & 0 & -10 & 490 \end{bmatrix} \quad \mathbf{F}(\mathbf{x}^{(0)})^{-1} = \begin{bmatrix} 0.1126 & -0.0089 & 0.0154 & 0.1106 \\ -0.0089 & 0.0057 & 0.0008 & -0.0087 \\ 0.0154 & 0.0008 & 0.0203 & 0.0155 \\ 0.1106 & -0.0087 & 0.0155 & 0.1107 \end{bmatrix}$$

$$\mathbf{F}(\mathbf{x}^{(0)})^{-1} \mathbf{g}^{(0)} = [1.4127, -0.1587, 0.2540, 0.2540]^T$$

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

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

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

Example

► Iteration 2.

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

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

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

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

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

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

Example

► Iteration 3.

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

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

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

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

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

Introduction

- ▶ 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 $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{d}^{(k)}$
- ▶ Step 1 requires the solution of an $n \times n$ system of linear equations. Thus, an efficient method for solving systems of linear equations is essential when using Newton's method.
- ▶ As in the one-variable case, Newton's method can be viewed as a technique for iteratively solving the equation

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

where $\mathbf{x} \in R^n$ and $\mathbf{g} : R^n \rightarrow R^n$. In this case $\mathbf{F}(\mathbf{x})$ is the Jacobian matrix of \mathbf{g} at \mathbf{x} ; that is, $\mathbf{F}(\mathbf{x})$ is the $n \times n$ matrix whose (i, j) entry is $(\partial g_i / \partial x_j)(\mathbf{x})$, $i, j = 1, 2, \dots, n$

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 $F(\mathbf{x}^{(k)})$ is not positive definite (recall Figure 7.7)
- ▶ Even if $F(\mathbf{x}^{(k)}) > 0$, Newton's method may not be a descent method; that is, it is possible that $f(\mathbf{x}^{(k+1)}) > f(\mathbf{x}^{(k)})$
 - ▶ This may occur if our starting point is far away from the solution
- ▶ Despite these drawbacks, Newton's method has superior convergence properties when the starting point is near the solution.

Newton's method works well if $f''(x) > 0$ everywhere. However, if $f''(x) < 0$ for some x , Newton's method may fail to converge to the minimizer.

Analysis of Newton's Method

- ▶ The convergence analysis of Newton's method when f is a quadratic function is straightforward. Newton's method reaches the point \mathbf{x}^* such that $\nabla f(\mathbf{x}^*) = \mathbf{0}$ in just one step starting from any initial point $\mathbf{x}^{(0)}$.

- ▶ Suppose that $\mathbf{Q} = \mathbf{Q}^T$ is invertible and $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{x}^T \mathbf{b}$
Then, $\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x}) = \mathbf{Q} \mathbf{x} - \mathbf{b}$ and $\mathbf{F}(\mathbf{x}) = \mathbf{Q}$

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

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

- ▶ Therefore, for the quadratic case the order of convergence of Newton's algorithm is ∞ for any initial point $\mathbf{x}^{(0)}$

Analysis of Newton's Method

- ▶ Theorem 9.1: Suppose that $f \in \mathcal{C}^3$ and $\mathbf{x}^* \in R^n$ 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 converge to \mathbf{x}^* with an order of convergence at least 2.

- ▶ Proof: The Taylor series expansion of ∇f about $\mathbf{x}^{(0)}$ yields

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

Because by assumption $f \in \mathcal{C}^3$ and $\mathbf{F}(\mathbf{x}^*)$ is invertible, there exist constants $\epsilon > 0$, $c_1 > 0$ and $c_2 > 0$ such that if $\mathbf{x}^{(0)}$, $\mathbf{x} \in \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\| \leq \epsilon\}$, we have

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

and by Lemma 5.3, $\mathbf{F}(\mathbf{x})^{-1}$ exists and satisfies

$$\|\mathbf{F}(\mathbf{x})^{-1}\| \leq c_2$$

Analysis of Newton's Method

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

$$\|\mathbf{F}(\mathbf{x})^{-1}\| \leq c_2$$

- ▶ The first inequality holds because the remainder term in the Taylor series expansion contains third derivatives of f that are continuous and hence bounded on $\{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\| \leq \epsilon\}$
- ▶ Suppose that $\mathbf{x}^{(0)} \in \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}^*\| \leq \epsilon\}$. Then, substituting $\mathbf{x} = \mathbf{x}^*$ in the inequality above and using the assumption that $\nabla f(\mathbf{x}^*) = \mathbf{0}$ we get

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

Analysis of Newton's Method

- ▶ Subtracting \mathbf{x}^* from both sides of Newton's algorithm and taking norms yields

$$\begin{aligned}\|\mathbf{x}^{(1)} - \mathbf{x}^*\| &= \|\mathbf{x}^{(0)} - \mathbf{x}^* - \mathbf{F}(\mathbf{x}^{(0)})^{-1} \nabla f(\mathbf{x}^{(0)})\| \\ &= \|\mathbf{F}(\mathbf{x}^{(0)})^{-1}(\mathbf{F}(\mathbf{x}^{(0)})(\mathbf{x}^{(0)} - \mathbf{x}^*) - \nabla f(\mathbf{x}^{(0)})\| \\ &\leq \|\mathbf{F}(\mathbf{x}^{(0)})^{-1}\| \|\mathbf{F}(\mathbf{x}^{(0)})(\mathbf{x}^{(0)} - \mathbf{x}^*) - \nabla f(\mathbf{x}^{(0)})\|\end{aligned}$$

- ▶ Applying the inequalities above involving the constants c_1 and c_2

$$\|\mathbf{x}^{(1)} - \mathbf{x}^*\| \leq c_1 c_2 \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2$$

- ▶ Suppose that $\mathbf{x}^{(0)}$ is such that

$$\|\mathbf{x}^{(0)} - \mathbf{x}^*\| \leq \frac{\alpha}{c_1 c_2} \quad \alpha \in (0, 1)$$

Then

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

Analysis of Newton's Method

- ▶ By induction, we obtain

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq c_1 c_2 \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2$$

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq \alpha \|\mathbf{x}^{(k)} - \mathbf{x}^*\|$$

Hence, $\lim_{k \rightarrow \infty} \|\mathbf{x}^{(k)} - \mathbf{x}^*\| = 0$ and therefore the sequence $\{\mathbf{x}_k\}$ converges to \mathbf{x}^* . The order of convergence is at least 2 because

$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| \leq c_1 c_2 \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2$. That is, $\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\| = O(\|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2)$

Analysis of Newton's Method

- ▶ Theorem 9.2: 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 $\alpha \in (0, \bar{\alpha})$

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

Analysis of Newton's Method

- **Proof:** Let $\phi(\alpha) = f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$, then using the chain rule, we obtain

$$\phi'(\alpha) = \nabla f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})^T \mathbf{d}^{(k)}$$

Hence, $\phi'(0) = \nabla f(\mathbf{x}^{(k)})^T \mathbf{d}^{(k)} = -\mathbf{g}^{(k)T} \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)} < 0$

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 $\alpha \in (0, \bar{\alpha})$

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

Analysis of Newton's Method

- ▶ Theorem 9.2 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 $-\mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}$

- ▶ 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 $\mathbf{F}(\mathbf{x}^{(k)}) \mathbf{d}^{(k)} = -\mathbf{g}^{(k)}$. In Chapters 10 and 11 we discuss this issue.
- ▶ The Hessian matrix may not be positive definite. In the next we describe a simple modification to overcome this problem.

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.

- ▶ ***Levenberg-Marquardt modification:***

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - (\mathbf{F}(\mathbf{x}^{(k)}) + \mu_k \mathbf{I})^{-1} \mathbf{g}^{(k)} \quad \mu_k \geq 0$$

- ▶ 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 are real, but may not all be positive.
- ▶ Consider the matrix $\mathbf{G} = \mathbf{F} + \mu \mathbf{I}$, where $\mu \geq 0$. Note that the eigenvalues of \mathbf{G} are $\lambda_1 + \mu, \dots, \lambda_n + \mu$.

Levenberg-Marquardt Modification

▶ Indeed,

$$\begin{aligned} G\mathbf{v}_i &= (\mathbf{F} + \mu\mathbf{I})\mathbf{v}_i \\ &= \mathbf{F}\mathbf{v}_i + \mu\mathbf{I}\mathbf{v}_i \\ &= \lambda_i\mathbf{v}_i + \mu\mathbf{v}_i \\ &= (\lambda_i + \mu)\mathbf{v}_i \end{aligned}$$

which shows that for all $i = 1, \dots, n$, \mathbf{v}_i is also an eigenvector of G with eigenvalue $\lambda_i + \mu$.

- ▶ If μ is sufficiently large, then all the eigenvalues of G are positive and G is positive definite.
- ▶ Accordingly, if the parameter μ_k in the Levenberg-Marquardt modification of Newton's algorithm is sufficiently large, then the search direction $\mathbf{d}^{(k)} = -(\mathbf{F}(\mathbf{x}^{(k)}) + \mu_k\mathbf{I})^{-1}\mathbf{g}^{(k)}$ always points in a descent direction.

Levenberg-Marquardt Modification

- ▶ If we further introduce a step size α_k

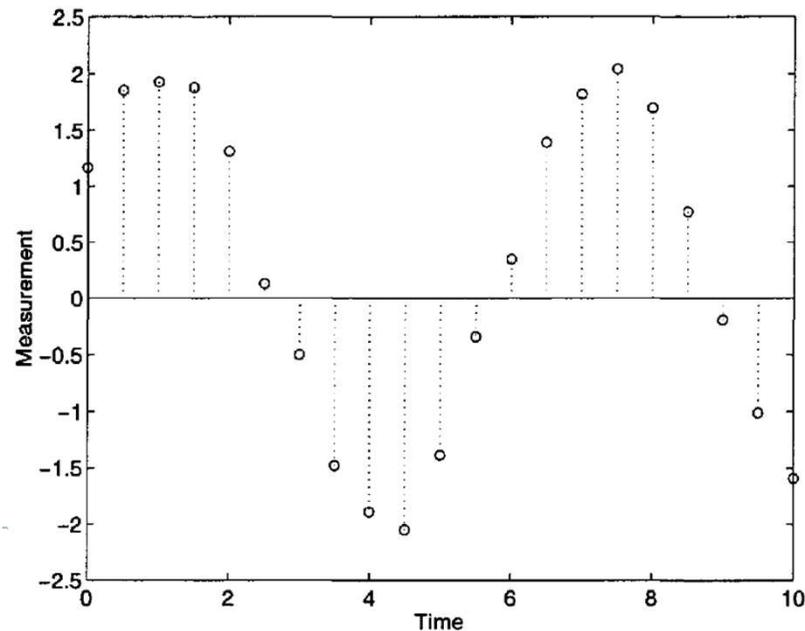
$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k (\mathbf{F}'(\mathbf{x}^{(k)}) + \mu_k \mathbf{I})^{-1} \mathbf{g}^{(k)}$$

then we are guaranteed that the descent property holds.

- ▶ By letting $\mu_k \rightarrow 0$, the Levenberg-Marquardt modification approaches the behavior of the pure Newton's method.
- ▶ By letting $\mu_k \rightarrow \infty$, this algorithm approaches a pure gradient method with small step size.
- ▶ In practice, we may start with a small value of μ_k and increase it slowly until we find that the iteration is descent: $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$

Newton's Method for Nonlinear Least Squares

- ▶ Consider minimize $\sum_{i=1}^m (r_i(\mathbf{x}))^2$, where $r_i : R^n \rightarrow R$, $i = 1, \dots, m$ are given functions. This particular problem is called a ***nonlinear least-squares problem***.
- ▶ Suppose that we are given m measurements of a process at m points in time. Let t_1, \dots, t_m denote the measurement times and y_1, \dots, y_m the measurements values. Note that $t_1 = 0$ and $t_{21} = 10$. We wish to fit a sinusoid to the measurement data.



Newton's Method for Nonlinear Least Squares

- ▶ The equation of the sinusoid is

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

with appropriate choices of the parameters A, ω, ϕ .

- ▶ 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 obtain the least-squares problem with

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

Newton's Method for Nonlinear Least Squares

- ▶ Defining $\mathbf{r} = [r_1, \dots, r_m]^T$, we write the objective function as $f(\mathbf{x}) = \mathbf{r}(\mathbf{x})^T \mathbf{r}(\mathbf{x})$. 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

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

- ▶ Denote the Jacobian matrix of \mathbf{r} by

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

- ▶ Thus, the gradient of f can be represented as

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

Newton's Method for Nonlinear Least Squares

- ▶ 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

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

Newton's Method for Nonlinear Least Squares

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

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

- ▶ In some applications, the matrix $\mathbf{S}(\mathbf{x})$ involving the second derivatives of the function \mathbf{r} can be ignored because its components are negligibly small.
- ▶ In this case Newton's algorithm reduces to what is commonly called the ***Gauss-Newton method***:

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

Note that the Gauss-Newton method does not require calculation of the second derivatives of \mathbf{r}

Example

$$r_i(\mathbf{x}) = y_i - A \sin(\omega t + \phi) \quad i = 1, \dots, 21$$

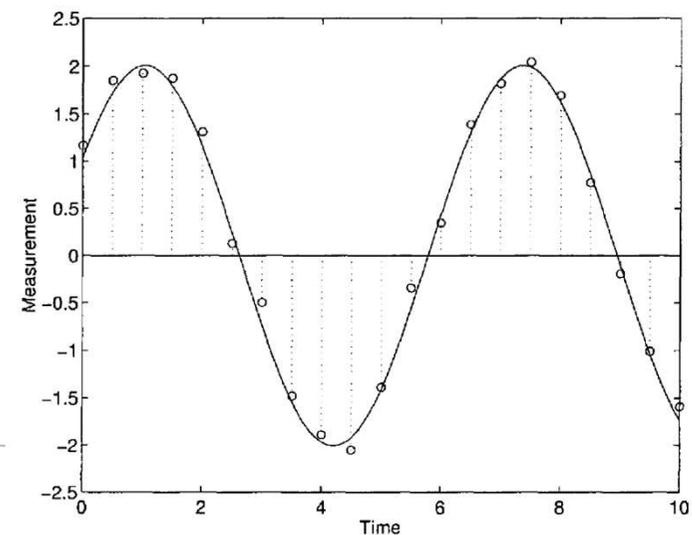
- ▶ The Jacobian matrix $J(\mathbf{x})$ in this problem is a 21×3 matrix with elements given by

$$(J(\mathbf{x}))_{(i,1)} = -\sin(\omega t_i + \phi)$$

$$(J(\mathbf{x}))_{(i,2)} = -t_i A \cos(\omega t_i + \phi) \quad i = 1, \dots, 21$$

$$(J(\mathbf{x}))_{(i,3)} = -A \cos(\omega t_i + \phi)$$

- ▶ We apply the Gauss-Newton algorithm to find the sinusoid of best fit.
- ▶ The parameters of this sinusoid are $A = 2.01, \omega = 0.992, \phi = 0.541$



Newton's Method for Nonlinear Least Squares

- ▶ A potential problem with the Gauss-Newton method is that the matrix $\mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x})$ may not be positive definite.
- ▶ This problem can be overcome using a Levenberg-Marquardt modification:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - (\mathbf{J}(\mathbf{x})^T \mathbf{J}(\mathbf{x}) + \mu_k \mathbf{I})^{-1} \mathbf{J}(\mathbf{x})^T \mathbf{r}(\mathbf{x})$$

- ▶ This is referred to in the literature as the ***Levenberg-Marquardt algorithm*** because the original modification was developed specifically for the nonlinear least-squares problem.
- ▶ An alternative interpretation of the Levenberg-Marquardt algorithm is to view the term $\mu_k \mathbf{I}$ as an approximation to $\mathbf{S}(\mathbf{x})$ in the Newton's algorithm.