

Chapter 11 Quasi-Newton Methods

An Introduction to Optimization

Spring, 2014

Wei-Ta Chu

Introduction

- ▶ In Newton's method, for a general nonlinear objective function, convergence to a solution cannot be guaranteed from an arbitrary initial point $\mathbf{x}^{(0)}$.
- ▶ The idea behind Newton's method is to locally approximate the function f being minimized, at every iteration, by a quadratic function. The minimizer for the quadratic approximation is used as the starting point for the next iteration.

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

- ▶ Guarantee that the algorithm has the descent property by modifying as follows

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

where α_k is chosen to ensure that

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

Introduction

- ▶ For example, we may choose $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha_k \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)})$. We can then determine an appropriate value of α_k by performing a line search in the direction $-\mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}$. Note that although the line search is simply the minimization of the real variable function $\phi_k(\alpha) = f(\mathbf{x}^{(k)} - \alpha \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)})$, it is not a trivial problem to solve.
- ▶ A computational drawback of Newton's method is the need to evaluate $\mathbf{F}(\mathbf{x}^{(k)})$ and solve the equation $\mathbf{F}(\mathbf{x}^{(k)}) \mathbf{d}^{(k)} = -\mathbf{g}^{(k)}$. To avoid the computation of $\mathbf{F}(\mathbf{x}^{(k)})^{-1}$, the quasi-Newton methods use an approximation to $\mathbf{F}(\mathbf{x}^{(k)})^{-1}$ in place of the true inverse.

Introduction

- ▶ Consider the formula

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \mathbf{H}_k \mathbf{g}^{(k)}$$

where \mathbf{H}_k is an $n \times n$ matrix and α is a positive search parameter. Expanding f about $\mathbf{x}^{(k)}$ yields

$$\begin{aligned} f(\mathbf{x}^{(k+1)}) &= f(\mathbf{x}^{(k)}) + \mathbf{g}^{(k)T}(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) + o(\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|) \\ &= f(\mathbf{x}^{(k)}) - \alpha \mathbf{g}^{(k)T} \mathbf{H}_k \mathbf{g}^{(k)} + o(\|\mathbf{H}_k \mathbf{g}^{(k)}\| \alpha) \end{aligned}$$

As α tends to zero, the second term on the right-hand side dominates the third. Thus, to guarantee a decrease in f for small α , we have to have

$$\mathbf{g}^{(k)T} \mathbf{H}_k \mathbf{g}^{(k)} > 0$$

A simple way to ensure this is to require that \mathbf{H}_k be positive definite.

Introduction

- ▶ **Proposition 11.1:** Let $f \in \mathcal{C}^1$, $\mathbf{x}^{(k)} \in \mathbb{R}^n$, $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) \neq \mathbf{0}$, and \mathbf{H}_k an $n \times n$ real symmetric positive definite matrix. If we set $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha \mathbf{H}_k \mathbf{g}^{(k)}$, where $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \mathbf{H}_k \mathbf{g}^{(k)})$, then $\alpha_k > 0$ and $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$

$$\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)}) \quad \nabla f(\mathbf{x}) = \mathbf{Q}\mathbf{x} - \mathbf{b}$$

Approximating the Inverse Hessian

- ▶ Let $\mathbf{H}_0, \mathbf{H}_1, \mathbf{H}_2, \dots$ be successive approximations of the inverse $\mathbf{F}(\mathbf{x}^{(k)})^{-1}$ of the Hessian.
- ▶ Suppose first that the Hessian matrix $\mathbf{F}(\mathbf{x})$ of the objective function f is constant and independent of \mathbf{x} . In other words, the objective function is quadratic, with Hessian $\mathbf{F}(\mathbf{x}) = \mathbf{Q}$ for all \mathbf{x} , where $\mathbf{Q} = \mathbf{Q}^T$. Then,

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

Let

$$\begin{aligned}\Delta \mathbf{g}^{(k)} &\triangleq \mathbf{g}^{(k+1)} - \mathbf{g}^{(k)} \\ \Delta \mathbf{x}^{(k)} &\triangleq \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\end{aligned}$$

Then, we may write

$$\Delta \mathbf{g}^{(k)} = \mathbf{Q}\Delta \mathbf{x}^{(k)}$$

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

Approximating the Inverse Hessian

- ▶ We start with a real symmetric positive definite matrix \mathbf{H}_0 .
Note that given k , the matrix \mathbf{Q}^{-1} satisfies

$$\mathbf{Q}^{-1}\Delta\mathbf{g}^{(i)} = \Delta\mathbf{x}^{(i)} \quad 0 \leq i \leq k$$

- ▶ Therefore, we also impose the requirement that the approximation \mathbf{H}_{k+1} of the Hessian satisfy

$$\mathbf{H}_{k+1}\Delta\mathbf{g}^{(i)} = \Delta\mathbf{x}^{(i)} \quad 0 \leq i \leq k$$

- ▶ If n steps are involved, then moving in n directions $\Delta\mathbf{x}^{(0)}, \Delta\mathbf{x}^{(1)}, \dots, \Delta\mathbf{x}^{(n-1)}$ yields

$$\mathbf{H}_n\Delta\mathbf{g}^{(0)} = \Delta\mathbf{x}^{(0)}$$

$$\mathbf{H}_n\Delta\mathbf{g}^{(1)} = \Delta\mathbf{x}^{(1)}$$

$$\vdots$$

$$\mathbf{H}_n\Delta\mathbf{g}^{(n-1)} = \Delta\mathbf{x}^{(n-1)}$$

Approximating the Inverse Hessian

- ▶ This set of equations can be represented as

$$\mathbf{H}_n[\Delta \mathbf{g}^{(0)}, \Delta \mathbf{g}^{(1)}, \dots, \Delta \mathbf{g}^{(n-1)}] = [\Delta \mathbf{x}^{(0)}, \Delta \mathbf{x}^{(1)}, \dots, \Delta \mathbf{x}^{(n-1)}]$$

Note that \mathbf{Q} satisfies

$$\mathbf{Q}[\Delta \mathbf{x}^{(0)}, \Delta \mathbf{x}^{(1)}, \dots, \Delta \mathbf{x}^{(n-1)}] = [\Delta \mathbf{g}^{(0)}, \Delta \mathbf{g}^{(1)}, \dots, \Delta \mathbf{g}^{(n-1)}]$$

and

$$\mathbf{Q}^{-1}[\Delta \mathbf{g}^{(0)}, \Delta \mathbf{g}^{(1)}, \dots, \Delta \mathbf{g}^{(n-1)}] = [\Delta \mathbf{x}^{(0)}, \Delta \mathbf{x}^{(1)}, \dots, \Delta \mathbf{x}^{(n-1)}]$$

Therefore, if $[\Delta \mathbf{g}^{(0)}, \Delta \mathbf{g}^{(1)}, \dots, \Delta \mathbf{g}^{(n-1)}]$ is nonsingular, then \mathbf{Q}^{-1} is determined uniquely after n steps, via

$$\mathbf{Q}^{-1} = \mathbf{H}_n = [\Delta \mathbf{x}^{(0)}, \Delta \mathbf{x}^{(1)}, \dots, \Delta \mathbf{x}^{(n-1)}][\Delta \mathbf{g}^{(0)}, \Delta \mathbf{g}^{(1)}, \dots, \Delta \mathbf{g}^{(n-1)}]^{-1}$$

Approximating the Inverse Hessian

- ▶ We conclude that if \mathbf{H}_n satisfies the equations

$$\mathbf{H}_n \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}, 0 \leq i \leq n - 1$$

then the algorithm $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{H}_k \mathbf{g}^{(k)}$,

$\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} - \alpha \mathbf{H}_k \mathbf{g}^{(k)})$, is guaranteed to solve problems with quadratic objective functions in $n + 1$ steps, because the update $\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} - \alpha_n \mathbf{H}_n \mathbf{g}^{(n)}$ is equivalent to Newton's algorithm.

Approximating the Inverse Hessian

- ▶ The quasi-Newton algorithms have the form

$$\begin{aligned} \mathbf{d}^{(k)} &= -\mathbf{H}_k \mathbf{g}^{(k)} \\ \alpha_k &= \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) \\ \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)} \end{aligned}$$

where the matrices $\mathbf{H}_0, \mathbf{H}_1, \dots$ are symmetric. In the quadratic case these matrices are required to satisfy

$$\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}, 0 \leq i \leq k$$

where $\Delta \mathbf{x}^{(i)} = \mathbf{x}^{(i+1)} - \mathbf{x}^{(i)} = \alpha_i \mathbf{d}^{(i)}$ and $\Delta \mathbf{g}^{(i)} = \mathbf{g}^{(i+1)} - \mathbf{g}^{(i)} = \mathbf{Q} \Delta \mathbf{x}^{(i)}$

It turns out that quasi-Newton methods are also conjugate direction methods.

Approximating the Inverse Hessian

- ▶ Theorem 11.1: Consider a quasi-Newton algorithm applied to a quadratic function with Hessian $\mathbf{Q} = \mathbf{Q}^T$ such that for $0 \leq k < n - 1$

$$\mathbf{H}_{k+1} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}, 0 \leq i \leq k$$

where $\mathbf{H}_{k+1} = \mathbf{H}_{k+1}^T$. If $\alpha_i \neq 0$, $0 \leq i \leq k$, then $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k+1)}$ are \mathbf{Q} -conjugate.

- ▶ Proof: We proceed by induction. We begin with the $k = 0$ case: that $\mathbf{d}^{(0)}$ and $\mathbf{d}^{(1)}$ are \mathbf{Q} -conjugate. Because $\alpha_0 \neq 0$, we can write

$\mathbf{d}^{(0)} = \Delta \mathbf{x}^{(0)} / \alpha_0$. Hence,

$$\mathbf{d}^{(1)T} \mathbf{Q} \mathbf{d}^{(0)} = -\mathbf{g}^{(1)T} \mathbf{H}_1 \mathbf{Q} \mathbf{d}^{(0)}$$

but $\mathbf{g}^{(1)T} \mathbf{d}^{(0)} = 0$ as a consequence

$$= -\mathbf{g}^{(1)T} \mathbf{H}_1 \frac{\mathbf{Q} \Delta \mathbf{x}^{(0)}}{\alpha_0}$$

of $\alpha_0 > 0$ being the minimizer of

$$= -\mathbf{g}^{(1)T} \frac{\mathbf{H}_1 \Delta \mathbf{g}^{(0)}}{\alpha_0}$$

$\phi(\alpha) = f(\mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)})$. Hence,

$$= -\mathbf{g}^{(1)T} \frac{\Delta \mathbf{x}^{(0)}}{\alpha_0}$$

$$\mathbf{d}^{(1)T} \mathbf{Q} \mathbf{d}^{(0)} = 0$$

$$= -\mathbf{g}^{(1)T} \mathbf{d}^{(0)}$$

Approximating the Inverse Hessian

- ▶ Assume that the result is true for $k - 1$. We now prove that the result for k , that is, that $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k+1)}$ are Q -conjugate. It suffices to show that $\mathbf{d}^{(k+1)T} Q \mathbf{d}^{(i)} = 0, 0 \leq i \leq k$. Given $0 \leq i \leq k$ using the same algebraic steps as in the $k = 0$ case, and using the assumption that $\alpha_i \neq 0$, we obtain

$$\begin{aligned} \mathbf{d}^{(k+1)T} Q \mathbf{d}^{(i)} &= -\mathbf{g}^{(k+1)T} \mathbf{H}_{k+1} Q \mathbf{d}^{(i)} \\ &\vdots \\ &= -\mathbf{g}^{(k+1)T} \mathbf{d}^{(i)} \end{aligned}$$

Because $\mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k)}$ are Q -conjugate by assumption, we conclude from Lemma 10.2 that $\mathbf{g}^{(k+1)T} \mathbf{d}^{(i)} = 0$. Hence, $\mathbf{d}^{(k+1)T} Q \mathbf{d}^{(i)} = 0$, which completes the proof.

The Rank One Correction Formula

- ▶ In the *rank one correction formula*, the correction term is symmetric and has the form $a_k \mathbf{z}^{(k)} \mathbf{z}^{(k)T}$, where $a_k \in R$ and $\mathbf{z}^{(k)} \in R^n$. The update equation is

$$\mathbf{H}_{k+1} = \mathbf{H}_k + a_k \mathbf{z}^{(k)} \mathbf{z}^{(k)T}$$

Note that

$$\text{rank}(\mathbf{z}^{(k)} \mathbf{z}^{(k)T}) = \text{rank} \left(\begin{bmatrix} z_1^{(k)} \\ \vdots \\ z_n^{(k)} \end{bmatrix} \begin{bmatrix} z_1^{(k)} & \dots & z_n^{(k)} \end{bmatrix} \right) = 1$$

and hence the name *rank one correction* [also called *single-rank symmetric* (SRF) algorithm].

The product $\mathbf{z}^{(k)} \mathbf{z}^{(k)T}$ is sometimes referred to as the *dyadic product* or *outer product*. Observe that if \mathbf{H}_k is symmetric, then so is \mathbf{H}_{k+1} .

The Rank One Correction Formula

- ▶ Our goal now is to determine a_k and $\mathbf{z}^{(k)}$, given \mathbf{H}_k , $\Delta\mathbf{g}^{(k)}$, $\Delta\mathbf{x}^{(k)}$ so that the required relationship discussed in Section 11.2 is satisfied; namely $\mathbf{H}_{k+1}\Delta\mathbf{g}^{(i)} = \Delta\mathbf{x}^{(i)}$, $i = 1, \dots, k$.
- ▶ To begin, consider the condition $\mathbf{H}_{k+1}\Delta\mathbf{g}^{(k)} = \Delta\mathbf{x}^{(k)}$. In other words, given \mathbf{H}_k , $\Delta\mathbf{g}^{(k)}$, $\Delta\mathbf{x}^{(k)}$, we wish to find a_k and $\mathbf{z}^{(k)}$ to ensure that

$$\mathbf{H}_{k+1}\Delta\mathbf{g}^{(k)} = (\mathbf{H}_k + a_k\mathbf{z}^{(k)}\mathbf{z}^{(k)T})\Delta\mathbf{g}^{(k)} = \Delta\mathbf{x}^{(k)}$$

- ▶ First note that $\mathbf{z}^{(k)T}\Delta\mathbf{g}^{(k)}$ is a scalar. Thus,

$$\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)} = (a_k\mathbf{z}^{(k)T}\Delta\mathbf{g}^{(k)})\mathbf{z}^{(k)}$$

and hence

$$\mathbf{z}^{(k)} = \frac{\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)}}{a_k(\mathbf{z}^{(k)T}\Delta\mathbf{g}^{(k)})}$$

$$\mathbf{H}_{k+1}\Delta\mathbf{g}^{(k)} = (\mathbf{H}_k + a_k\mathbf{z}^{(k)}\mathbf{z}^{(k)T})\Delta\mathbf{g}^{(k)} = \Delta\mathbf{x}^{(k)}$$

The Rank One Correction Formula

- ▶ We can now determine

$$a_k\mathbf{z}^{(k)}\mathbf{z}^{(k)T} = \frac{(\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)})(\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)})^T}{a_k(\mathbf{z}^{(k)T}\Delta\mathbf{g}^{(k)})^2}$$

Hence,

$$\mathbf{H}_{k+1} = \mathbf{H}_{(k)} + \frac{(\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)})(\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)})^T}{a_k(\mathbf{z}^{(k)T}\Delta\mathbf{g}^{(k)})^2}$$

- ▶ The next step is to express the denominator of the second term on the right-hand side as a function of the given quantities

\mathbf{H}_k , $\Delta\mathbf{g}^{(k)}$, $\Delta\mathbf{x}^{(k)}$. Premultiply $\Delta\mathbf{x}^{(k)} - \mathbf{H}_k\Delta\mathbf{g}^{(k)} = (a_k\mathbf{z}^{(k)T}\Delta\mathbf{g}^{(k)})\mathbf{z}^{(k)}$ by $\Delta\mathbf{g}^{(k)T}$ to obtain

$$\Delta\mathbf{g}^{(k)T}\Delta\mathbf{x}^{(k)} - \Delta\mathbf{g}^{(k)T}\mathbf{H}_k\Delta\mathbf{g}^{(k)} = \Delta\mathbf{g}^{(k)T}a_k\mathbf{z}^{(k)}\mathbf{z}^{(k)T}\Delta\mathbf{g}^{(k)}$$

$$\Delta \mathbf{g}^{(k)T} \Delta \mathbf{x}^{(k)} - \Delta \mathbf{g}^{(k)T} \mathbf{H}_k \Delta \mathbf{g}^{(k)} = \Delta \mathbf{g}^{(k)T} a_k \mathbf{z}^{(k)} \mathbf{z}^{(k)T} \Delta \mathbf{g}^{(k)}$$

The Rank One Correction Formula

- ▶ Observe that a_k is a scalar and so is $\Delta \mathbf{g}^{(k)T} \mathbf{z}^{(k)} = \mathbf{z}^{(k)T} \Delta \mathbf{g}^{(k)}$.

Thus,

$$\Delta \mathbf{g}^{(k)T} \Delta \mathbf{x}^{(k)} - \Delta \mathbf{g}^{(k)T} \mathbf{H}_k \Delta \mathbf{g}^{(k)} = a_k (\mathbf{z}^{(k)T} \Delta \mathbf{g}^{(k)})^2$$

Taking this relation into account yields

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})^T}{\Delta \mathbf{g}^{(k)T} (\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})}$$

$$\mathbf{H}_{k+1} = \mathbf{H}_{(k)} + \frac{(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})^T}{a_k (\mathbf{z}^{(k)T} \Delta \mathbf{g}^{(k)})^2}$$

Rank One Algorithm

- ▶ 1. Set $k := 0$; select $\mathbf{x}^{(0)}$ and a real symmetric positive definite \mathbf{H}_0
- ▶ 2. If $\mathbf{g}^{(k)} = \mathbf{0}$, stop; else, $\mathbf{d}^{(k)} = -\mathbf{H}_k \mathbf{g}^{(k)}$
- ▶ 3. Compute

$$\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$$

- ▶ 4. Compute

$$\Delta \mathbf{x}^{(k)} = \alpha_k \mathbf{d}^{(k)}$$

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

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})^T}{\Delta \mathbf{g}^{(k)T}(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})}$$

- ▶ 5. Set $k := k + 1$; go to step 2.

Rank One Algorithm

- ▶ However, what we want is $\mathbf{H}_{k+1}\Delta\mathbf{g}^{(i)} = \Delta\mathbf{x}^{(i)}, i = 1, \dots, k$
- ▶ Theorem 11.2: For the rank one algorithm applied to the quadratic with Hessian $\mathbf{Q} = \mathbf{Q}^T$, we have $\mathbf{H}_{k+1}\Delta\mathbf{g}^{(i)} = \Delta\mathbf{x}^{(i)}$
 $0 \leq i \leq k$
- ▶ Proof.

Example

- ▶ Let $f(x_1, x_2) = x_1^2 + \frac{1}{2}x_2^2 + 3$. Apply the rank one correction algorithm to minimize f . Use $\mathbf{x}^{(0)} = [1, 2]^T$ and $\mathbf{H}_0 = \mathbf{I}_2$
- ▶ We can represent f as

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \mathbf{x} + 3$$

Thus,

$$\mathbf{g}^{(k)} = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \mathbf{x}^{(k)}$$

Because $\mathbf{H}_0 = \mathbf{I}_2$, $\mathbf{d}^{(0)} = -\mathbf{g}^{(0)} = [-2, -2]^T$

Example

- ▶ The objective function is quadratic, and hence

$$\begin{aligned}\alpha_0 &= \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)}) \\ &= -\frac{\mathbf{g}^{(0)T} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)T} \mathbf{Q} \mathbf{d}^{(0)}} = \frac{[2, 2] \begin{bmatrix} 2 \\ 2 \end{bmatrix}}{[2, 2] \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix}} = \frac{2}{3}\end{aligned}$$

and thus $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} = [-\frac{1}{3}, \frac{2}{3}]^T$

We then compute

$$\Delta \mathbf{x}^{(0)} = \alpha_0 \mathbf{d}^{(0)} = [-\frac{4}{3}, -\frac{4}{3}]^T$$

$$\mathbf{g}^{(1)} = \mathbf{Q} \mathbf{x}^{(1)} = [-\frac{2}{3}, \frac{2}{3}]^T$$

$$\Delta \mathbf{g}^{(0)} = \mathbf{g}^{(1)} - \mathbf{g}^{(0)} = [-\frac{8}{3}, -\frac{4}{3}]^T$$

Example

► Because

$$\Delta \mathbf{g}^{(0)T} (\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)}) = \left[-\frac{8}{3}, -\frac{4}{3}\right] \begin{bmatrix} \frac{4}{3} \\ 0 \end{bmatrix} = -\frac{32}{9}$$

We obtain

$$\mathbf{H}_1 = \mathbf{H}_0 + \frac{(\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)})(\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)})^T}{\Delta \mathbf{g}^{(0)T} (\Delta \mathbf{x}^{(0)} - \mathbf{H}_0 \Delta \mathbf{g}^{(0)})} = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{bmatrix}$$

Therefore,

$$\mathbf{d}^{(1)} = -\mathbf{H}_1 \mathbf{g}^{(1)} = \left[\frac{1}{3}, -\frac{2}{3}\right]^T$$

$$\alpha_1 = -\frac{\mathbf{g}^{(1)T} \mathbf{d}^{(1)}}{\mathbf{d}^{(1)T} \mathbf{Q} \mathbf{d}^{(1)}} = 1$$

We now compute $\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{d}^{(1)} = [0, 0]^T$

Note that $\mathbf{g}^{(2)} = \mathbf{0}$, and therefore $\mathbf{x}^{(2)} = \mathbf{x}^*$. As expected, the algorithm solves the problem in two steps.

Note that the directions $\mathbf{d}^{(0)}, \mathbf{d}^{(1)}$ are \mathbf{Q} -conjugate, in accordance with

The Rank One Correction Formula

- ▶ Unfortunately, the rank one correction algorithm is not very satisfactory for several reasons.
 - ▶ The matrix \mathbf{H}_{k+1} that the rank one algorithm generates may not be positive definite and thus $\mathbf{d}^{(k+1)}$ may not be a descent direction. This happens even in the quadratic case.
 - ▶ If $\Delta \mathbf{g}^{(k)T}(\Delta \mathbf{x}^{(k)} - \mathbf{H}_k \Delta \mathbf{g}^{(k)})$ is close to zero, then there may be numerical problems in evaluating \mathbf{H}_{k+1} .
- ▶ Fortunately, alternative algorithms have been developed for updating \mathbf{H}_k . In particular, if we use a “rank two” update, then \mathbf{H}_k is guaranteed to be positive definite for all k , provided that the line search is exact.

The DFP Algorithm

- ▶ This algorithm was developed by Davidon (1959), Fletcher, and Powell (1963).
- ▶ The DFP algorithm is also known as the *variable metric algorithm*.
- ▶ DFP Algorithm
 - ▶ 1. Set $k := 0$; select $\mathbf{x}^{(0)}$ and a real symmetric positive definite \mathbf{H}_0
 - ▶ 2. If $\mathbf{g}^{(k)} = \mathbf{0}$, stop; else, $\mathbf{d}^{(k)} = -\mathbf{H}_k \mathbf{g}^{(k)}$
 - ▶ 3. Compute $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)})$
 $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$

- ▶ 4. Compute

$$\begin{aligned}\Delta \mathbf{x}^{(k)} &= \alpha_k \mathbf{d}^{(k)} \\ \Delta \mathbf{g}^{(k)} &= \mathbf{g}^{(k+1)} - \mathbf{g}^{(k)} \\ \mathbf{H}_{k+1} &= \mathbf{H}_k + \frac{\Delta \mathbf{x}^{(k)} \Delta \mathbf{x}^{(k)T}}{\Delta \mathbf{x}^{(k)T} \Delta \mathbf{g}^{(k)}} - \frac{[\mathbf{H}_k \Delta \mathbf{g}^{(k)}][\mathbf{H}_k \Delta \mathbf{g}^{(k)}]^T}{\Delta \mathbf{g}^{(k)T} \mathbf{H}_k \Delta \mathbf{g}^{(k)}}\end{aligned}$$

- ▶ ²³ 5. Set $k := k + 1$; go to step 2.

The DFP Algorithm

- ▶ Theorem 11.3: In the DFP algorithm applied to the quadratic with Hessian $Q = Q^T$, we have $H_{k+1}\Delta g^{(i)} = \Delta x^{(i)}$, $0 \leq i \leq k$
- ▶ Theorem 11.4: Suppose that $g^{(k)} \neq 0$. In the DFP algorithm, if H_k is positive definite, then so is H_{k+1} .

Example

- ▶ Locate the minimizer of $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \mathbf{x} - \mathbf{x}^T \begin{bmatrix} -1 \\ 1 \end{bmatrix}$, $\mathbf{x} \in R^2$

Use the initial point $\mathbf{x}^{(0)} = [0, 0]^T$ and $\mathbf{H}_0 = \mathbf{I}_2$

- ▶ Note that in this case

$$\mathbf{g}^{(k)} = \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \mathbf{x}^{(k)} - \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

Hence, $\mathbf{g}^{(0)} = [1, -1]^T$

$$\mathbf{d}^{(0)} = -\mathbf{H}_0\mathbf{g}^{(0)} = -\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

Because f is a quadratic function,

$$\alpha_0 = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(0)} + \alpha\mathbf{d}^{(0)}) = -\frac{\mathbf{g}^{(0)T}\mathbf{d}^{(0)}}{\mathbf{d}^{(0)T}\mathbf{Q}\mathbf{d}^{(0)}} = 1$$

Example

▶ Therefore, $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} = [-1, 1]^T$

▶ We then compute $\Delta \mathbf{x}^{(0)} = \mathbf{x}^{(1)} - \mathbf{x}^{(0)} = [-1, 1]^T$

$$\mathbf{g}^{(1)} = \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} - \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

$$\Delta \mathbf{g}^{(0)} = \mathbf{g}^{(1)} - \mathbf{g}^{(0)} = [-2, 0]^T$$

▶ Observe that

$$\Delta \mathbf{x}^{(0)} \Delta \mathbf{x}^{(0)T} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$$

$$\Delta \mathbf{x}^{(0)T} \Delta \mathbf{g}^{(0)} = 2$$

$$\mathbf{H}_0 \Delta \mathbf{g}^{(0)} = \begin{bmatrix} -2 \\ 0 \end{bmatrix}$$

Thus,

$$(\mathbf{H}_0 \Delta \mathbf{g}^{(0)}) (\mathbf{H}_0 \Delta \mathbf{g}^{(0)})^T = \begin{bmatrix} 4 & 0 \\ 0 & 0 \end{bmatrix}$$

$$\Delta \mathbf{g}^{(0)T} \mathbf{H}_0 \Delta \mathbf{g}^{(0)} = 4$$

Example

$$\begin{aligned} \mathbf{H}_1 &= \mathbf{H}_0 + \frac{\Delta \mathbf{x}^{(0)} \Delta \mathbf{x}^{(0)T}}{\Delta \mathbf{x}^{(0)T} \Delta \mathbf{g}^{(0)}} - \frac{[\mathbf{H}_0 \Delta \mathbf{g}^{(0)}][\mathbf{H}_0 \Delta \mathbf{g}^{(0)}]^T}{\Delta \mathbf{g}^{(0)T} \mathbf{H}_0 \Delta \mathbf{g}^{(0)}} \\ &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} - \frac{1}{4} \begin{bmatrix} 4 & 0 \\ 0 & 0 \end{bmatrix} \\ &= \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{3}{2} \end{bmatrix} \end{aligned}$$

- ▶ We now compute $\mathbf{d}^{(1)} = -\mathbf{H}_1 \mathbf{g}^{(1)} = [0, 1]^T$ and

$$\alpha_1 = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(1)} + \alpha \mathbf{d}^{(1)}) = -\frac{\mathbf{g}^{(1)T} \mathbf{d}^{(1)}}{\mathbf{d}^{(1)T} \mathbf{Q} \mathbf{d}^{(1)}} = \frac{1}{2}$$

Hence, $\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_0 \mathbf{d}^{(1)} = [-1, 3/2]^T = \mathbf{x}^*$, because f is a quadratic function of two variables.

- ▶ Note that we have $\mathbf{d}^{(0)T} \mathbf{Q} \mathbf{d}^{(1)} = \mathbf{d}^{(1)T} \mathbf{Q} \mathbf{d}^{(0)} = 0$; that is, $\mathbf{d}^{(0)}$ and $\mathbf{d}^{(1)}$ are \mathbf{Q} -conjugate directions.

The DFP Algorithm

- ▶ The DFP algorithm is superior to the rank one algorithm in that it preserves the positive definiteness of H_k .
- ▶ However, it turns out that in the case of larger nonquadratic problems the algorithm has the tendency of sometimes getting stuck. This phenomenon is attributed to H_k becoming nearly singular.

The BFGS Algorithm

- ▶ Suggested by Broyden, Fletcher, Goldfarb, and Shanno.
- ▶ Recall that the updating formulas for the approximation of the inverse of the Hessian matrix were based on satisfying the equations

$$\mathbf{H}_{k+1}\Delta\mathbf{g}^{(i)} = \Delta\mathbf{x}^{(i)} \quad 0 \leq i \leq k$$

which were derived from $\Delta\mathbf{g}^{(i)} = \mathbf{Q}\Delta\mathbf{x}^{(i)}$, $0 \leq i \leq k$. We then formulated update formulas for the approximations to the inverse of the Hessian matrix \mathbf{Q}^{-1} .

- ▶ An alternative to approximating \mathbf{Q}^{-1} is to approximate \mathbf{Q} itself.

The BFGS Algorithm

- ▶ Let B_k be our estimate of Q at the k th step. We require B_{k+1} to satisfy $\Delta g^{(i)} = B_{k+1} \Delta x^{(i)}$, $0 \leq i \leq k$.
- ▶ Notice that this set of equations is similar to the previous set of equations for H_{k+1} , the only difference being that the roles of $\Delta x^{(i)}$ and $\Delta g^{(i)}$ are interchanged.
- ▶ Given any update formula for H_k , a corresponding update formula for B_k can be found by interchanging the roles of B_k and H_k and of $\Delta g^{(k)}$ and $\Delta x^{(k)}$. In particular, the BFGS update for B_k corresponds to the DFP update for H_k . Formulas related in this way are said to be *dual* or *complementary*.

The BFGS Algorithm

- ▶ Recall that the DFP update for the approximation \mathbf{H}_k of the inverse Hessian is

$$\mathbf{H}_{k+1}^{DFP} = \mathbf{H}_k + \frac{\Delta \mathbf{x}^{(k)} \Delta \mathbf{x}^{(k)T}}{\Delta \mathbf{x}^{(k)T} \Delta \mathbf{g}^{(k)}} - \frac{[\mathbf{H}_k \Delta \mathbf{g}^{(k)}][\mathbf{H}_k \Delta \mathbf{g}^{(k)}]^T}{\Delta \mathbf{g}^{(k)T} \mathbf{H}_k \Delta \mathbf{g}^{(k)}}$$

- ▶ Using the complementarity concept, we can easily obtain an update equation for the approximation \mathbf{B}_k of the Hessian

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \frac{\Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)T}}{\Delta \mathbf{g}^{(k)T} \Delta \mathbf{x}^{(k)}} - \frac{[\mathbf{B}_k \Delta \mathbf{x}^{(k)}][\mathbf{B}_k \Delta \mathbf{x}^{(k)}]^T}{\Delta \mathbf{x}^{(k)T} \mathbf{B}_k \Delta \mathbf{x}^{(k)}}$$

- ▶ To obtain the BFGS update for the approximation of the inverse Hessian, we take the inverse of \mathbf{B}_{k+1} to obtain

$$\begin{aligned} \mathbf{H}_{k+1}^{BFGS} &= (\mathbf{B}_{k+1})^{-1} \\ &= \left(\mathbf{B}_k + \frac{\Delta \mathbf{g}^{(k)} \Delta \mathbf{g}^{(k)T}}{\Delta \mathbf{g}^{(k)T} \Delta \mathbf{x}^{(k)}} - \frac{[\mathbf{B}_k \Delta \mathbf{x}^{(k)}][\mathbf{B}_k \Delta \mathbf{x}^{(k)}]^T}{\Delta \mathbf{x}^{(k)T} \mathbf{B}_k \Delta \mathbf{x}^{(k)}} \right)^{-1} \end{aligned}$$

The BFGS Algorithm

- ▶ Lemma 11.1 *Sherman-Morrison formula*: Let A be a nonsingular matrix. Let u and v be column vectors such that $1 + v^T A u \neq 0$. Then, $A + uv^T$ is nonsingular, and its inverse can be written in terms of A^{-1} using the following formula:

$$(A + uv^T)^{-1} = A^{-1} - \frac{(A^{-1}u)(v^T A^{-1})}{1 + v^T A^{-1}u}$$

- ▶ From Lemma 11.1 it follows that if A^{-1} is known, then the inverse of the matrix A augmented by a rank one matrix can be obtained by a modification of the matrix A^{-1} .

The BFGS Algorithm

- ▶ Applying Lemma 11.1 twice to B_{k+1} yields

$$\mathbf{H}_{k+1}^{BFGS} = \mathbf{H}_k + \left(1 + \frac{\Delta \mathbf{g}^{(k)T} \mathbf{H}_k \Delta \mathbf{g}^{(k)}}{\Delta \mathbf{g}^{(k)T} \Delta \mathbf{x}^{(k)}} \right) \frac{\Delta \mathbf{x}^{(k)} \Delta \mathbf{x}^{(k)T}}{\Delta \mathbf{x}^{(k)T} \Delta \mathbf{g}^{(k)}} - \frac{\mathbf{H}_k \Delta \mathbf{g}^{(k)} \Delta \mathbf{x}^{(k)T} + (\mathbf{H}_k \Delta \mathbf{g}^{(k)} \Delta \mathbf{x}^{(k)T})^T}{\Delta \mathbf{g}^{(k)T} \Delta \mathbf{x}^{(k)}}$$

- ▶ Recall that for the quadratic case the DFP algorithm satisfies $\mathbf{H}_{k+1}^{DFP} \Delta \mathbf{g}^{(i)} = \mathbf{x}^{(i)}, 0 \leq i \leq k$. Therefore, the BFGS update for B_k satisfies $B_{k+1} \Delta \mathbf{x}^{(i)} = \mathbf{g}^{(i)}, 0 \leq i \leq k$. By construction of the BFGS formula for \mathbf{H}_{k+1}^{BFGS} , we conclude that $\mathbf{H}_{k+1}^{BFGS} \Delta \mathbf{g}^{(i)} = \Delta \mathbf{x}^{(i)}, 0 \leq i \leq k$. Hence, the BFGS algorithm enjoys all the properties of quasi-Newton methods, including the conjugate directions property. Moreover, the BFGS algorithm also inherits the positive definiteness property of the DFP algorithm; that is, if $\mathbf{g}^{(k)} \neq \mathbf{0}$ and $\mathbf{H}_k > 0$, then $\mathbf{H}_{k+1}^{BFGS} > 0$.

Example

- ▶ The BFGS formula is often far more efficient than the DFP formula.
- ▶ Use the BFGS method to minimize $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Q} \mathbf{x} - \mathbf{x}^T \mathbf{b} + \log(\pi)$

$$\mathbf{Q} = \begin{bmatrix} 5 & -3 \\ -3 & 2 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

- ▶ Take $\mathbf{H}_0 = \mathbf{I}_2$ and $\mathbf{x}_0 = [0, 0]^T$. Verify that $\mathbf{H}_2 = \mathbf{Q}^{-1}$.
- ▶ We have $\mathbf{d}^{(0)} = -\mathbf{g}^{(0)} = -(\mathbf{Q}\mathbf{x}^{(0)} - \mathbf{b}) = \mathbf{b} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$

The objective function is a quadratic, and hence we can use the following formula to compute α_0

$$\alpha_0 = -\frac{\mathbf{g}^{(0)T} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)T} \mathbf{Q} \mathbf{d}^{(0)}} = \frac{1}{2}$$

Example

► Therefore, $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} = \begin{bmatrix} 0 \\ 1/2 \end{bmatrix}$

To compute $\mathbf{H}_1 = \mathbf{H}_1^{BFGS}$, we need the following quantities:

$$\begin{aligned}\Delta \mathbf{x}^{(0)} &= \mathbf{x}^{(1)} - \mathbf{x}^{(0)} = \begin{bmatrix} 0 \\ 1/2 \end{bmatrix} \\ \mathbf{g}^{(1)} &= \mathbf{Q}\mathbf{x}^{(1)} - \mathbf{b} = \begin{bmatrix} -3/2 \\ 0 \end{bmatrix} \\ \Delta \mathbf{g}^{(0)} &= \mathbf{g}^{(1)} - \mathbf{g}^{(0)} = \begin{bmatrix} -3/2 \\ 1 \end{bmatrix}\end{aligned}$$

Therefore,

$$\begin{aligned}\mathbf{H}_1 &= \mathbf{H}_0 + \left(1 + \frac{\Delta \mathbf{g}^{(0)T} \mathbf{H}_0 \Delta \mathbf{g}^{(0)}}{\Delta \mathbf{g}^{(0)T} \Delta \mathbf{x}^{(0)}} \right) \frac{\Delta \mathbf{x}^{(0)} \Delta \mathbf{x}^{(0)T}}{\Delta \mathbf{x}^{(0)T} \Delta \mathbf{g}^{(0)}} \\ &\quad - \frac{\Delta \mathbf{x}^{(0)} \Delta \mathbf{g}^{(0)T} \mathbf{H}_0 + \mathbf{H}_0 \Delta \mathbf{g}^{(0)} \Delta \mathbf{x}^{(0)T}}{\Delta \mathbf{g}^{(0)T} \Delta \mathbf{x}^{(0)}} = \begin{bmatrix} 1 & 3/2 \\ 3/2 & 11/4 \end{bmatrix}\end{aligned}$$

Example

▶ Hence, we have $\mathbf{d}^{(1)} = -\mathbf{H}_1 \mathbf{g}^{(1)} = \begin{bmatrix} 3/2 \\ 9/4 \end{bmatrix}$

$$\alpha_1 = -\frac{\mathbf{g}^{(1)T} \mathbf{d}^{(1)}}{\mathbf{d}^{(1)T} \mathbf{Q} \mathbf{d}^{(1)}} = 2$$

Therefore,

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{d}^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

- ▶ Because our objective function is a quadratic on R^2 , $\mathbf{x}^{(2)}$ is the minimizer. Notice that the gradient at $\mathbf{x}^{(2)}$ is $\mathbf{0}$; that is, $\mathbf{g}^{(2)} = \mathbf{0}$

Example

- ▶ To verify that $\mathbf{H}_2 = \mathbf{Q}^{-1}$, we compute

$$\Delta \mathbf{x}^{(1)} = \mathbf{x}^{(2)} - \mathbf{x}^{(1)} = \begin{bmatrix} 3 \\ 9/2 \end{bmatrix}$$

$$\Delta \mathbf{g}^{(1)} = \mathbf{g}^{(2)} - \mathbf{g}^{(1)} = \begin{bmatrix} 3/2 \\ 0 \end{bmatrix}$$

$$\begin{aligned} \mathbf{H}_2 &= \mathbf{H}_1 + \left(1 + \frac{\Delta \mathbf{g}^{(1)T} \mathbf{H}_1 \Delta \mathbf{g}^{(1)}}{\Delta \mathbf{g}^{(1)T} \Delta \mathbf{x}^{(1)}} \right) \frac{\Delta \mathbf{x}^{(1)} \Delta \mathbf{x}^{(1)T}}{\Delta \mathbf{x}^{(1)T} \Delta \mathbf{g}^{(1)}} \\ &\quad - \frac{\Delta \mathbf{x}^{(1)} \Delta \mathbf{g}^{(1)T} \mathbf{H}_1 + \mathbf{H}_1 \Delta \mathbf{g}^{(1)} \Delta \mathbf{x}^{(1)T}}{\Delta \mathbf{g}^{(1)T} \Delta \mathbf{x}^{(1)}} = \begin{bmatrix} 2 & 3 \\ 3 & 5 \end{bmatrix} \end{aligned}$$

$$\Rightarrow \mathbf{H}_2 \mathbf{Q} = \mathbf{Q} \mathbf{H}_2 = \mathbf{I}_2 \quad \Rightarrow \quad \mathbf{H}_2 = \mathbf{Q}^{-1}$$