

Optimization Algorithms: Theory and Software Implementation

Prof. Thirumulanathan D

Department of Mathematics

Institute of IIT Kanpur

Lecture: 30

This lecture concludes our discussion on the DFP method and introduces the BFGS method and the broader Broyden family of quasi-Newton algorithms. We begin by proving the final property of the DFP method concerning the preservation of positive definiteness.

Theorem: Consider minimizing a general function f using the DFP method where the step size α_k is chosen using an exact line search. If the matrix B_k is positive definite, then the updated matrix B_{k+1} is also positive definite.

Proof:

We aim to show that for any non-zero vector x , the quadratic form $x^T B_{k+1} x > 0$. Recall the DFP update formula:

$$B_{k+1} = B_k + (\delta_k \delta_k^T) / (\delta_k^T \gamma_k) - (B_k \gamma_k \gamma_k^T B_k) / (\gamma_k^T B_k \gamma_k)$$

The quadratic form is therefore:

$$x^T B_{k+1} x = x^T B_k x + (x^T \delta_k)^2 / (\delta_k^T \gamma_k) - (x^T B_k \gamma_k)^2 / (\gamma_k^T B_k \gamma_k)$$

(Refer Slide Time 3:22)

$B^{k+1} d^j = d^j \quad \forall j=0, 1, \dots, k$
 $\therefore B^n d^j = d^j \quad \forall j=0, 1, \dots, n-1$
 $B^n [d^0 | d^1 | \dots | d^{n-1}] = [d^0 | d^1 | \dots | d^{n-1}]$
 $B^n H [d^0 | d^1 | \dots | d^{n-1}] = [d^0 | d^1 | \dots | d^{n-1}]$
If d^0, \dots, d^{n-1} are linearly independent, then $B^n = H^{-1}$.
 d^0, \dots, d^{n-1} linearly independent $\Leftrightarrow (\sum_{i=0}^{n-1} \beta_i d^i = 0 \Rightarrow \beta_i = 0 \quad \forall i)$
If $\sum \beta_i d^i = 0$, then $d^{jT} H (\sum \beta_i d^i) = 0 \Rightarrow \beta_j (d^{jT} H d^j) = 0$
 $\therefore \beta_j = 0$. True for any $j=0, 1, \dots, n-1$.
 $\Rightarrow d^0, \dots, d^{n-1}$ are linearly independent.
When $B^n = H^{-1}$, then $d^n = -H^{-1} g^n$. Newton's method. So, $x^{n+1} = x^*$. \square

Theorem: Consider minimizing f using DFP method, where α^k is chosen using exact line search. Then $B^k > 0 \Rightarrow B^{k+1} > 0$.

Proof: Assume $B^k > 0$. Then we need to show that $x^T B^{k+1} x > 0 \quad \forall x \neq 0$.

$$x^T B^{k+1} x = x^T B^k x + \frac{x^T \delta^k \delta^{kT} x}{\delta^{kT} \gamma^k} - \frac{x^T B^k \gamma^k \gamma^{kT} B^k x}{\gamma^{kT} B^k \gamma^k}$$

To analyze this expression, we define two vectors using the positive definite square root of B_k (which exists since B_k is positive definite):

$$a = B_k^{1/2} x$$

$$b = B_k^{1/2} \gamma_k$$

We can now rewrite the expression as:

$$\begin{aligned} x^T B_{k+1} x &= (a^T a)(b^T b) / (b^T b) - (a^T b)^2 / (b^T b) + (x^T \delta_k)^2 / (\delta_k^T \gamma_k) \\ &= [(a^T a)(b^T b) - (a^T b)^2] / (b^T b) + (x^T \delta_k)^2 / (\delta_k^T \gamma_k) \end{aligned}$$

By the **Cauchy-Schwarz inequality**, $(a^T b)^2 \leq (a^T a)(b^T b)$, so the first term is non-negative. It equals zero if and only if a is parallel to b , i.e., $a = \lambda b$ for some scalar λ . This implies $B_k^{1/2} x = \lambda B_k^{1/2} \gamma_k$, and since $B_k^{1/2}$ is invertible, $x = \lambda \gamma_k$.

We now examine the denominator $\delta_k^T \gamma_k$. Given that an exact line search is used, we have the property that $\delta_k^T g_{k+1} = 0$. Therefore:

$$\delta_k^T \gamma_k = \delta_k^T (g_{k+1} - g_k) = -\delta_k^T g_k$$

Since $\delta_k = \alpha_k d_k$ and the search direction is $d_k = -B_k g_k$, this becomes:

$$\delta_k^T \gamma_k = -\alpha_k (-B_k g_k)^T g_k = \alpha_k g_k^T B_k g_k$$

Because B_k is positive definite and $\alpha_k > 0$ for a minimization step, $g_k^T B_k g_k > 0$ for any non-zero g_k . Thus, $\delta_k^T \gamma_k > 0$.

Consequently, the second term, $(x^T \delta_k)^2 / (\delta_k^T \gamma_k)$, is always non-negative and is zero only if

$$x^T \delta_k = 0.$$

Now, if $x^T B_{k+1} x = 0$, both terms in its expression must be zero. The first term is zero only if $x = \lambda \gamma_k$. Substituting this into the second term yields $(\lambda \gamma_k^T \delta_k)^2 / (\delta_k^T \gamma_k) = \lambda^2 (\delta_k^T \gamma_k)$, which is greater than zero for any non-zero λ (and hence any non-zero x). This is a contradiction. Therefore, $x^T B_{k+1} x$ cannot be zero for any non-zero x , and it must be strictly positive. This proves that B_{k+1} is positive definite.

This property highlights the importance of using an exact line search with the DFP method to ensure the descent property is maintained throughout the iterations.

(Refer Slide Time 15:40)

$$x^T B^{k+1} x = x^T B^k x - \frac{(x^T B^k \eta^k)^2}{\eta^{kT} B^k \eta^k} + \frac{(\eta^{kT} x)^2}{\eta^{kT} \eta^k}$$

Define $a = (B^k)^{\frac{1}{2}} x$, $b = (B^k)^{\frac{1}{2}} \eta^k$.

$$\therefore x^T B^{k+1} x = \frac{(a^T a)(b^T b) - (a^T b)^2}{b^T b} + \frac{(\eta^{kT} x)^2}{\eta^{kT} \eta^k}$$

$|<a, b>| \leq \|a\| \|b\| \iff (a^T b)^2 \leq (a^T a)(b^T b)$. Cauchy-Schwarz inequality

$\therefore (a^T a)(b^T b) - (a^T b)^2 \geq 0$.

$\delta^{kT} \eta^k = \delta^{kT} (g^{k+1} - g^k)$. But $\delta^{kT} g^{k+1} = 0$ when we use exact line search.

$$\delta^{kT} \eta^k = -\delta^{kT} g^k = -x^k \delta^{kT} g^k = x^k g^k \delta^{kT} B^k g^k > 0 \quad \forall g^k \neq 0 \text{ when } B^k > 0$$

This proves that $x^T B^{k+1} x \geq 0 \quad \forall x \neq 0$.

$(a^T a)(b^T b) = (a^T b)^2$ iff $a = \lambda b$ for some scalar λ .

$\Rightarrow x = \lambda \eta^k, \lambda \neq 0$.

$\delta^{kT} x = \lambda \delta^{kT} \eta^k > 0$

$\therefore x^T B^{k+1} x > 0 \quad \forall x \neq 0 \Rightarrow B^{k+1} > 0$

We now transition to the BFGS method, named after Broyden, Fletcher, Goldfarb, and Shanno. The BFGS update formula is given by:

$$B_{k+1} = B_k + [1 + (\gamma_k^T B_k \gamma_k) / (\delta_k^T \gamma_k)] * (\delta_k \delta_k^T) / (\delta_k^T \gamma_k) - (\delta_k \gamma_k^T B_k + B_k \gamma_k \delta_k^T) / (\delta_k^T \gamma_k)$$

The overall algorithm structure remains identical to the DFP and rank-one correction methods; only the update rule for B_k changes. The BFGS method also belongs to the **Broyden family** of updates, which can be expressed as a convex combination of the DFP and BFGS updates:

$$B_{k+1} = \varphi * B\{DFP\}_{k+1} + (1 - \varphi) * B\{BFGS\}_{k+1}$$

for some φ in the interval $[0, 1]$.

Setting $\varphi = 1$ gives the **DFP**

update, $\varphi = 0$ gives the **BFGS** update, and values in between yield a hybrid method.

Numerical experiments on quadratic functions show that all these methods—rank-one correction, DFP, BFGS, and the Broyden family—exhibit identical performance for these problems. They all converge to the exact solution in at most n steps and exactly compute the inverse Hessian H^{-1} upon completion.

For more general non-quadratic functions, such as

$$f(x_1, x_2) = x_1^2 e^{x_2} + x_2^2 e^{x_1},$$

the performance of these methods is also very similar. They successfully converge to a minimizer from some initial points (e.g., $(1, 1)$ or $(-0.5, -0.5)$) but may converge to a saddle point from others (e.g., $(-\sqrt{2}, -\sqrt{2})$).

The rank-one correction method can sometimes be numerically unstable if the denominator $(\delta_k - B_k \gamma_k)^T \gamma_k$ becomes very small, which motivated the development of the more robust rank-two updates like DFP and BFGS.

This concludes our discussion on unconstrained optimization algorithms.

To summarize, we have covered:

1. **Gradient Descent:** $d_k = -\nabla f(x_k)$
2. **Conjugate Gradient (Fletcher-Reeves):** $d_0 = -g_0, d_k = -g_k + \beta_k d_{k-1}$
 where $\beta_k = (g_k^T g_k) / (g_{k-1}^T g_{k-1})$
3. **Newton's Method:** $d_k = -H_k^{-1} g_k$
4. **Damped Newton's Method:** Uses backtracking line search to choose α_k for the Newton direction.
5. **Modified Newton's Method:** $d_k = -(H_k + \lambda_k I)^{-1} g_k$, where λ_k is chosen to ensure the matrix is positive definite.
6. **Quasi-Newton Methods:** $d_k = -B_k g_k$, where B_k is updated to satisfy the quasi-Newton condition $B_{k+1} \gamma_k = \delta_k$ using various formulas (Rank-One, DFP, BFGS, Broyden Family).

(Refer Slide Time 30:50)

BFGS method (Broyden - Fletcher - Goldfarb - Shanno)

$$B^{k+1} = B^k + \left(1 + \frac{g_k^T B^k g_k}{\rho_k^T g_k}\right) \frac{\delta_k \delta_k^T}{\delta_k^T g_k} - \frac{(\delta_k^T B^k g_k + B^k g_k \delta_k^T)}{\delta_k^T g_k}$$

Broyden family:

$$B^{k+1} = \phi B_{DFP}^{k+1} + (1-\phi) B_{BFGS}^{k+1} \quad \text{for some } \phi \in [0, 1].$$

$\phi = 1 \Rightarrow$ DFP method
 $\phi = 0 \Rightarrow$ BFGS method
 $\phi \in (0, 1) \Rightarrow$ convex combination

Rank-one correction, DFP method, BFGS method, Broyden family.

Gradient descent: $d^k = -g^k$
 Conjugate Gradient: $d^0 = -g^0, d^k = -g^k + \left(\frac{g^{k+1} g^{k+1}}{g^k g^k}\right) d^{k-1}$
 Newton's method: $d^k = -(H^k)^{-1} g^k, \alpha^k = 1$
 Damped Newton: $-d_0 - \alpha^k$ by backtracking
 Modified Newton: $d^k = -(H^k + \lambda^k)^{-1} g^k, \alpha^k = 1$
 Quasi-Newton: $d^k = -B^k g^k, B^{k+1} = \dots$

These foundational concepts are essential as we proceed to constrained optimization problems in the subsequent weeks. Thank you