

Quasi-Newton methods

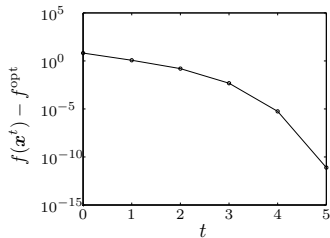
Yuanming Shi and Ye Shi

ShanghaiTech University

Newton's method

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

$$\mathbf{x}^{t+1} = \mathbf{x}^t - (\nabla^2 f(\mathbf{x}^t))^{-1} \nabla f(\mathbf{x}^t)$$



- quadratic convergence: attains ε accuracy within $O(\log \log \frac{1}{\varepsilon})$ iterations
- typically requires storing and inverting Hessian $\nabla^2 f(\mathbf{x}) \in \mathbb{R}^{n \times n}$
- a single iteration may last forever; prohibitive storage requirement

Quasi-Newton methods

key idea: approximate the Hessian matrix using only gradient information

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \underbrace{\mathbf{H}_t}_{\text{surrogate of } (\nabla^2 f(\mathbf{x}^t))^{-1}} \nabla f(\mathbf{x}^t)$$

challenges: how to find a good approximation $\mathbf{H}_t \succ \mathbf{0}$ of $(\nabla^2 f(\mathbf{x}^t))^{-1}$

- using only gradient information
- using limited memory
- achieving super-linear convergence

Criterion for choosing H_t

Consider the following *approximate quadratic model* of $f(\cdot)$:

$$f_t(\mathbf{x}) := f(\mathbf{x}^{t+1}) + \langle \nabla f(\mathbf{x}^{t+1}), \mathbf{x} - \mathbf{x}^{t+1} \rangle + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{t+1})^\top \mathbf{H}_{t+1}^{-1} (\mathbf{x} - \mathbf{x}^{t+1})$$

which satisfies

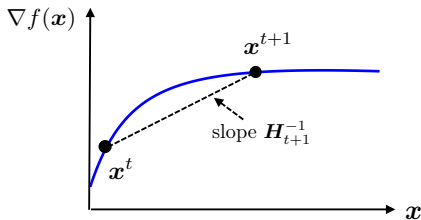
$$\nabla f_t(\mathbf{x}) = \nabla f(\mathbf{x}^{t+1}) + \mathbf{H}_{t+1}^{-1} (\mathbf{x} - \mathbf{x}^{t+1})$$

One reasonable criterion: **gradient matching** for the latest two iterates:

$$\nabla f_t(\mathbf{x}^t) = \nabla f(\mathbf{x}^t) \tag{13.1a}$$

$$\nabla f_t(\mathbf{x}^{t+1}) = \nabla f(\mathbf{x}^{t+1}) \tag{13.1b}$$

Secant equation



(13.1b) holds automatically. To satisfy (13.1a), one requires

$$\begin{aligned} \nabla f(\mathbf{x}^{t+1}) + \mathbf{H}_{t+1}^{-1}(\mathbf{x}^t - \mathbf{x}^{t+1}) &= \nabla f(\mathbf{x}^t) \\ \Leftrightarrow \underbrace{\mathbf{H}_{t+1}^{-1}(\mathbf{x}^{t+1} - \mathbf{x}^t)}_{\text{secant equation}} &= \nabla f(\mathbf{x}^{t+1}) - \nabla f(\mathbf{x}^t) \end{aligned}$$

- the secant equation requires that \mathbf{H}_{t+1}^{-1} maps the displacement $\mathbf{x}^{t+1} - \mathbf{x}^t$ into the change of gradients $\nabla f(\mathbf{x}^{t+1}) - \nabla f(\mathbf{x}^t)$

Secant equation

$$\mathbf{H}_{t+1} \underbrace{(\nabla f(\mathbf{x}^{t+1}) - \nabla f(\mathbf{x}^t))}_{=: \mathbf{y}_t} = \underbrace{\mathbf{x}^{t+1} - \mathbf{x}^t}_{=: \mathbf{s}_t} \quad (13.2)$$

- only possible when $\mathbf{s}_t^\top \mathbf{y}_t > 0$, since

$$\mathbf{s}_t^\top \mathbf{y}_t = \mathbf{y}_t^\top \mathbf{H}_{t+1} \mathbf{y}_t > 0$$

- admit an infinite number of solutions, since the degrees of freedom $O(n^2)$ in choosing \mathbf{H}_{t+1}^{-1} far exceeds the number of constraints n in (13.2)
- which \mathbf{H}_{t+1}^{-1} shall we choose?

Broyden-Fletcher-Goldfarb-Shanno (BFGS) method

Broyden, Fletcher, Goldfarb, Shanno



Closeness to H_t

In addition to the secant equation, choose H_{t+1} sufficiently close to H_t :

$$\begin{aligned} \text{minimize}_{\mathbf{H}} \quad & \|\mathbf{H} - \mathbf{H}_t\| \\ \text{subject to} \quad & \mathbf{H} = \mathbf{H}^\top \\ & \mathbf{H}\mathbf{y}_t = \mathbf{s}_t \end{aligned}$$

for some norm $\|\cdot\|$

- exploit past information regarding H_t
- choosing different norms $\|\cdot\|$ results in different quasi-Newton methods

Choice of norm in BFGS

Choosing $\|M\| := \|W^{1/2}MW^{1/2}\|_F$ for *any* weight matrix W obeying $Ws_t = y_t$, we get

$$\begin{aligned} \text{minimize}_H \quad & \|W^{1/2}(H - H_t)W^{1/2}\|_F \\ \text{subject to} \quad & H = H^\top \\ & Hy_t = s_t \end{aligned}$$

This admits a closed-form expression

$$\underbrace{H_{t+1} = (I - \rho_t s_t y_t^\top) H_t (I - \rho_t y_t s_t^\top) + \rho_t s_t s_t^\top}_{\text{BFGS update rule; } H_{t+1} \succeq 0 \text{ if } H_t \succeq 0} \quad (13.3)$$

$$\text{with } \rho_t = \frac{1}{y_t^\top s_t}$$

An alternative interpretation

\mathbf{H}_{t+1} is also the solution to

$$\begin{aligned} \text{minimize}_{\mathbf{H}} \quad & \underbrace{\langle \mathbf{H}_t, \mathbf{H}^{-1} \rangle - \log \det (\mathbf{H}_t \mathbf{H}^{-1}) - n}_{\text{KL divergence between } \mathcal{N}(\mathbf{0}, \mathbf{H}^{-1}) \text{ and } \mathcal{N}(\mathbf{0}, \mathbf{H}_t^{-1})} \\ \text{subject to} \quad & \mathbf{H} \mathbf{y}_t = \mathbf{s}_t \end{aligned}$$

- minimizing some sort of KL divergence subject to the secant equation constraints

BFGS methods

Algorithm 13.1 BFGS

- 1: **for** $t = 0, 1, \dots$ **do**
 - 2: $\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \mathbf{H}_t \nabla f(\mathbf{x}^t)$ (line search to determine η_t)
 - 3: $\mathbf{H}_{t+1} = (\mathbf{I} - \rho_t \mathbf{s}_t \mathbf{y}_t^\top) \mathbf{H}_t (\mathbf{I} - \rho_t \mathbf{y}_t \mathbf{s}_t^\top) + \rho_t \mathbf{s}_t \mathbf{s}_t^\top$, where $\mathbf{s}_t = \mathbf{x}^{t+1} - \mathbf{x}^t$, $\mathbf{y}_t = \nabla f(\mathbf{x}^{t+1}) - \nabla f(\mathbf{x}^t)$, and $\rho_t = \frac{1}{\mathbf{y}_t^\top \mathbf{s}_t}$
-

- each iteration costs $O(n^2)$ (in addition to computing gradients)
- no need to solve linear systems or invert matrices
- *no magic formula for initialization*; possible choices: approximate inverse Hessian at \mathbf{x}^0 , or identity matrix

Rank-2 update on H_t^{-1}

From the Sherman-Morrison-Woodbury formula
 $(\mathbf{A} + \mathbf{U}\mathbf{V}^\top)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{I} + \mathbf{V}^\top\mathbf{A}^{-1}\mathbf{U})^{-1}\mathbf{V}^\top\mathbf{A}^{-1}$, we can
show that the BFGS rule is equivalent to

$$\underbrace{H_{t+1}^{-1} = H_t^{-1} - \frac{1}{s_t^\top H_t^{-1} s_t} H_t^{-1} s_t s_t^\top H_t^{-1} + \rho_t y_t y_t^\top}_{\text{rank-2 update}}$$

Local superlinear convergence

Theorem 13.1 (informal)

Suppose f is strongly convex and has Lipschitz-continuous Hessian. Under mild conditions, BFGS achieves

$$\lim_{t \rightarrow \infty} \frac{\|\mathbf{x}^{t+1} - \mathbf{x}^*\|_2}{\|\mathbf{x}^t - \mathbf{x}^*\|_2} = 0$$

- *iteration complexity*: larger than Newton methods but smaller than gradient methods
- *asymptotic result*: holds when $t \rightarrow \infty$

Key observation

The BFGS update rule achieves

$$\lim_{t \rightarrow \infty} \frac{\|(\mathbf{H}_t^{-1} - \nabla^2 f(\mathbf{x}^*))(\mathbf{x}^{t+1} - \mathbf{x}^t)\|_2}{\|\mathbf{x}^{t+1} - \mathbf{x}^t\|_2} = 0$$

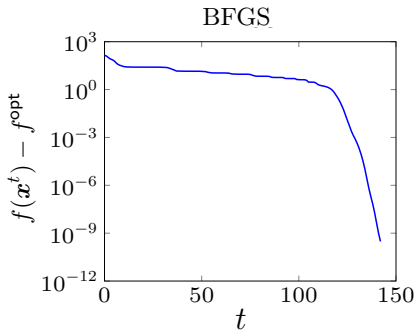
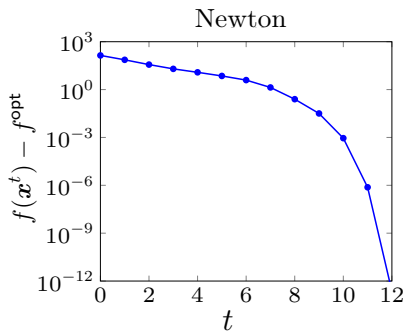
Implications

- even though \mathbf{H}_t^{-1} may not converge to $\nabla^2 f(\mathbf{x}^*)$, it becomes an increasingly more accurate approximation of $\nabla^2 f(\mathbf{x}^*)$ along the search direction $\mathbf{x}^{t+1} - \mathbf{x}^t$
- asymptotically, $\mathbf{x}^{t+1} - \mathbf{x}^t \approx \underbrace{-(\nabla^2 f(\mathbf{x}^t))^{-1} \nabla f(\mathbf{x}^t)}_{\text{Newton search direction}}$

Numerical example

— EE236C lecture notes

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^n} \quad \mathbf{c}^\top \mathbf{x} - \sum_{i=1}^N \log(b_i - \mathbf{a}_i^\top \mathbf{x})$$



$$n = 100, N = 500$$

Limited-memory quasi-Newton methods

Hessian matrices are usually dense. For large-scale problems, even storing the (inverse) Hessian matrices is prohibitive

Instead of storing full Hessian approximations, one may want to maintain more parsimonious approximation of the Hessians, using only a few vectors

Limited-memory BFGS (L-BFGS)

$$\underbrace{H_{t+1} = V_t^\top H_t V_t + \rho_t \mathbf{s}_t \mathbf{s}_t^\top}_{\text{BFGS update rule}} \quad \text{with } V_t = I - \rho_t \mathbf{y}_t \mathbf{s}_t^\top$$

key idea: maintain a modified version of H_t *implicitly* by storing m (e.g. 20) most recent vector pairs $(\mathbf{s}_t, \mathbf{y}_t)$

Limited-memory BFGS (L-BFGS)

L-BFGS maintains

$$\begin{aligned} \mathbf{H}_t^L &= \mathbf{V}_{t-1}^\top \cdots \mathbf{V}_{t-m}^\top \mathbf{H}_{t,0}^L \mathbf{V}_{t-m} \cdots \mathbf{V}_{t-1} \\ &\quad + \rho_{t-m} \mathbf{V}_{t-1}^\top \cdots \mathbf{V}_{t-m+1}^\top \mathbf{s}_{t-m} \mathbf{s}_{t-m}^\top \mathbf{V}_{t-m+1} \cdots \mathbf{V}_{t-1} \\ &\quad + \rho_{t-m+1} \mathbf{V}_{t-1}^\top \cdots \mathbf{V}_{t-m+2}^\top \mathbf{s}_{t-m+1} \mathbf{s}_{t-m+1}^\top \mathbf{V}_{t-m+1} \cdots \mathbf{V}_{t-1} \\ &\quad + \cdots + \rho_{t-1} \mathbf{s}_{t-1} \mathbf{s}_{t-1}^\top \end{aligned}$$

- can be computed recursively
- initialization $\mathbf{H}_{t,0}^L$ may vary from iteration to iteration
- only needs to store $\{(\mathbf{s}_i, \mathbf{y}_i)\}_{t-m \leq i < t}$

Reference

- [1] "*Numerical optimization*, J. Nocedal, S. Wright, 2000.
- [2] "*Optimization methods for large-scale systems, EE236C lecture notes*," L. Vandenberghe, UCLA.
- [3] "*Optimization methods for large-scale machine learning*," L. Bottou et al., arXiv, 2016.
- [4] "*Convex optimization, EE364B lecture notes*," S. Boyd, Stanford.