

Quasi-Newton Methods

Friday, March 26, 2021

9:56 AM

following Nocedal + Wright textbook

gold-standard for large-scale smooth, unconstrained optimization

Quadratic approximation at the current iterate x_k (let $p = x - x_k$)

$$m_k(p) = \underbrace{f_k}_{f(x_k)} + \underbrace{\langle \nabla f_k, p \rangle}_{\nabla f(x_k)} + \frac{1}{2} \langle p | B_k | p \rangle$$

$$B_k \succ 0$$

$\Rightarrow m$ is strongly convex

$$\tilde{x}_{k+1} = x_k + p_k, \quad p_k = \arg\min_p m_k(p)$$

Ex: Gradient Descent

$$B_k = L \cdot I, \quad L \text{ is Lipsch. constant of } \nabla f$$

(i.e., $\nabla^2 f(x_k) \leq L \cdot I$)

Ex: Newton's Method

$$B_k = \nabla^2 f(x_k)$$

$$p = -B_k^{-1} \cdot \nabla f_k$$

Sometimes do linesearch, $x_{k+1} = x_k + \alpha \cdot p^*$

$$\alpha \leq 1$$

Quasi-Newton method

means use the above framework w/ B_k such that

① B_k is more accurate than $L \cdot I$

② B_k is cheaper than Newton

both forming B_k and inverting B_k

Math trick is to construct B_{k+1} by updating B_k

To find x_{k+2} , minimize $m_{k+1}(p) := f_{k+1} + \langle \nabla f_{k+1}, p \rangle + \frac{1}{2} \langle p | B_{k+1} | p \rangle$

$$\nabla m_{k+1}(p) = \nabla f_{k+1} + B_{k+1} \cdot p$$

$$\text{so } \nabla m_{k+1}(0) = \nabla f_{k+1}$$

we've automatically enforced:

$$\text{① } m_{k+1}(0) = f_{k+1}$$

$$\text{② } \nabla m_{k+1}(0) = \nabla f_{k+1}$$

} x_{k+1}

Freedom in B_{k+1}

let's impose

$$\text{③ } \nabla m_{k+1}(\text{previous iterate}) = \nabla f_k$$

} x_k

$$x_{k+1} = x_k + p_k$$

Notation

$$S_k = x_{k+1} - x_k$$

$$y_k = \nabla f_{k+1} - \nabla f_k$$

$$p = x - x_{k+1}$$

\uparrow x_k previous iterate, $p = x_k - x_{k+1} = -S_k$

$$(3) \nabla m_{k+1}(-S_k) = \nabla f_{k+1} - B_{k+1} \cdot S_k \stackrel{\text{desired}}{=} \nabla f_k$$

$$B_{k+1} \cdot S_k = y_k \quad \text{"secant equation"}$$

Most/all quasi-Newton methods choose B_{k+1} to satisfy the secant equation.

Observe: want $B_{k+1} \succ 0$

$$\underbrace{\langle S_k | B_{k+1} | S_k \rangle}_{>0} = \langle S_k | y_k \rangle$$

$$\langle S_k | y_k \rangle > 0 \quad \text{curvature condition}$$

$$\text{i.e., } \underbrace{\langle x_{k+1} - x_k, \nabla f_{k+1} - \nabla f_k \rangle}_{>0} \stackrel{?}{>} 0$$

Strictly monotone happens automatically if f is strictly convex.

So if f isn't strictly convex, that complicates quasi-Newton method.
eg. add a linesearch

B_{k+1} chosen to solve secant equations. $x \in \mathbb{R}^n$ so $B_{k+1} \in \mathcal{S}^{n \times n}$

deg. of freedom of B_{k+1} is $\frac{n \cdot (n+1)}{2}$ \nwarrow n constraints

Ex: $n=1$, 1 d.o.f., 1 constraint. B_{k+1} is completely determined

Secant method

Ex: $n=2$, 3 d.o.f., 2 constraints.

$\Rightarrow n > 1$, there's a whole family of reasonable quasi-Newton methods

How to choose B_{k+1} ?

Standard ways: either solve $B_{k+1} = \underset{B \succ 0}{\operatorname{argmin}} \|B - B_k\|_w^2$
 $B \cdot S_k = y_k$

Notation

B approximates $\nabla^2 f$

H approximates $(\nabla^2 f)^{-1}$

$$\text{or } B_{k+1}^{-1} = \underset{B^{-1} \succ 0}{\operatorname{argmin}} \|H_{k+1} - H_k\|_w^2$$

$$B_{k+1}^{-1} y_k = S_k$$

$\|B\|_w$ is some norm, chosen so problem \uparrow has a closed-form solution.

$$B_k^{-1} = H_k$$

Broyden class

Issue of inverting B_k : not an issue,

1) use H_k not B_k

$$B_{k+1} = B_k + u \cdot u^T$$

and/or 2) B_{k+1} is a low-rank update of B_k

so use Sherman-Morrison-Woodbury formula

BFGS

Works w/ H update above, special choice of weighted norm

$$H_{k+1} = (I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) + \rho_k s_k s_k^T, \quad \rho_k = \frac{1}{\langle y_k, s_k \rangle}$$

$$H_0 = \frac{\langle y_1, s_1 \rangle}{\langle y_1, y_1 \rangle} \cdot I \quad (\text{aka Barzilai-Borwein or "spectral" stepsize})$$

Benefits:

- more accurate than gradient descent
- cheaper than Newton ($O(n^2)$ to update H_{k+1} , rather than $O(n^3)$ for Newton to invert $\nabla^2 f_k$)

Disadvantages

- still needs $O(n^2)$ memory

❗

Theory: iffy.

If f is nonconvex, does BFGS always converge to a stationary pt.?
unknown!

Thm: Nocedal & Wright, Thm 6.5

If $0 < \mu I \leq \nabla^2 f(x) \leq L \cdot I \quad \forall x$, then BFGS sequence converges to the global minimizer.

Thm: local convergence is often super-linear
(see book for details)

Limited-memory BFGS "L-BFGS"

Recall BFGS update $\rho_k = \frac{1}{\langle y_k, s_k \rangle}$, let $V_k = I - \rho_k y_k s_k^T$

$$H_{k+1} = V_k^T H_k V_k + \rho_k s_k s_k^T$$

$$s_k = x_{k+1} - x_k$$

$$y_k = \nabla f_{k+1} - \nabla f_k$$

so to compute matrix-vector multiply (eg., w/ $x = \nabla f(x_{k+1})$)

$$H_{k+1} \cdot x = V_k^T H_k \underbrace{V_k x}_z + \rho_k \underbrace{s_k s_k^T x}_{\text{easy}}$$

$z = V_k x$ easy

What about $H_k \cdot z$?

compute recursively! $H_{k+1} \cdot x$ computed using $\{s_k, y_k, H_k\}$

For BFGS, recurse until $k=0$, $H_0 = \text{const} \cdot I$.

For **L-BFGS**, only recurse back $m \ll n$ steps,

and set "base case" to $H_0^{(k)} = \gamma_k \cdot I$, $\gamma_k = \frac{\langle s_{k-1}, y_{k-1} \rangle}{\|y_{k-1}\|^2}$

Only need to store $O(mn)$ data

$\{s_i, y_i\}_{i=k-m+1}^{i=k}$ and flop count for $H_{k+1} \cdot x$
is $O(mn)$

Typically $m \in \{3, 20\}$

Cool fact: $m=0$ "memoryless" BFGS, w/ exact linesearch,
is the HS variant of nonlinear CG.