

Quasi-Newton Methods

Zico Kolter

(notes by Ryan Tibshirani, Javier Peña, Zico Kolter)

Convex Optimization 10-725

Last time: primal-dual interior-point methods

Given the problem

$$\begin{aligned} \min_x \quad & f(x) \\ \text{subject to} \quad & h(x) \leq 0 \\ & Ax = b \end{aligned}$$

where f , $h = (h_1, \dots, h_m)$, all convex and twice differentiable, and strong duality holds. **Central path** equations:

$$r(x, u, v) = \begin{pmatrix} \nabla f(x) + Dh(x)^T u + A^T v \\ -\text{diag}(u)h(x) - 1/t \\ Ax - b \end{pmatrix} = 0$$

subject to $u > 0$, $h(x) < 0$

Primal dual interior point method: repeat updates

$$(x^+, u^+, v^+) = (x, u, v) + s(\Delta x, \Delta u, \Delta v)$$

where $(\Delta x, \Delta u, \Delta v)$ is defined by Newton step:

$$\begin{bmatrix} H_{\text{pd}}(x) & Dh(x)^T & A^T \\ -\text{diag}(u)Dh(x) & -\text{diag}(h(x)) & 0 \\ A & 0 & 0 \end{bmatrix} \begin{pmatrix} \Delta x \\ \Delta u \\ \Delta v \end{pmatrix} = -r(x, u, v)$$

and $H_{\text{pd}}(x) = \nabla^2 f(x) + \sum_{i=1}^m u_i \nabla^2 h_i(x)$

- Step size $s > 0$ is chosen by backtracking, while maintaining $u > 0$, $h(x) < 0$
- Primal-dual iterates are not necessarily feasible
- But often converges faster than barrier method

Outline

Today:

- Quasi-Newton motivation
- SR1, BFGS, DFP, Broyden class
- Convergence analysis
- Limited memory BFGS
- Stochastic quasi-Newton

Gradient descent and Newton revisited

Back to unconstrained, smooth convex optimization

$$\min_x f(x)$$

where f is convex, twice differentiable, and $\text{dom}(f) = \mathbb{R}^n$. Recall **gradient descent** update:

$$x^+ = x - t\nabla f(x)$$

and **Newton's method** update:

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

- Newton's method has (local) quadratic convergence, versus linear convergence of gradient descent
- But Newton iterations are much more expensive ...

Quasi-Newton methods

Two main steps in Newton iteration:

- Compute Hessian $\nabla^2 f(x)$
- Solve the system $\nabla^2 f(x)\Delta x = -\nabla f(x)$

Each of these two steps could be expensive

Quasi-Newton methods repeat updates of the form

$$x^+ = x + t\Delta x$$

where direction Δx is defined by linear system

$$B\Delta x = -\nabla f(x)$$

for some approximation B of $\nabla^2 f(x)$. We want B to be easy to compute, and $B\Delta x = g$ to be easy to solve

Some history

- In the mid 1950s, W. Davidon was a mathematician/physicist at Argonne National Lab
- He was using coordinate descent on an optimization problem and his computer kept crashing before finishing
- He figured out a way to accelerate the computation, leading to the first quasi-Newton method (soon Fletcher and Powell followed up on his work)
- Although Davidon's contribution was a major breakthrough in optimization, his original paper was rejected
- In 1991, after more than 30 years, his paper was published in the first issue of the SIAM Journal on Optimization
- In addition to his remarkable work in optimization, Davidon was a peace activist (see the book "The Burglary")

Quasi-Newton template

Let $x^{(0)} \in \mathbb{R}^n$, $B^{(0)} \succ 0$. For $k = 1, 2, 3, \dots$, repeat:

1. Solve $B^{(k-1)} \Delta x^{(k-1)} = -\nabla f(x^{(k-1)})$
2. Update $x^{(k)} = x^{(k-1)} + t_k \Delta x^{(k-1)}$
3. Compute $B^{(k)}$ from $B^{(k-1)}$

Different quasi-Newton methods implement Step 3 differently. As we will see, commonly we can compute $(B^{(k)})^{-1}$ from $(B^{(k-1)})^{-1}$

Basic idea: as $B^{(k-1)}$ already contains info about the Hessian, use suitable matrix update to form $B^{(k)}$

Reasonable requirement for $B^{(k)}$:

$$\nabla f(x^{(k)}) = \nabla f(x^{(k-1)}) + B^{(k)}(x^{(k)} - x^{(k-1)})$$

Secant equation

We can equivalently write latter condition as

$$\nabla f(x^+) = \nabla f(x) + B^+(x^+ - x)$$

Letting $y = \nabla f(x^+) - \nabla f(x)$, and $s = x^+ - x$ this becomes

$$B^+s = y$$

This is called the **secant equation**

In addition to the secant equation, we want:

- B^+ to be symmetric
- B^+ to be “close” to B
- $B \succ 0 \Rightarrow B^+ \succ 0$

Symmetric rank-one update

Let's try an update of the form

$$B^+ = B + a u u^T$$

The secant equation $B^+ s = y$ yields

$$(a u^T s) u = y - B s$$

This only holds if u is a multiple of $y - B s$. Putting $u = y - B s$, we solve the above, $a = 1 / (y - B s)^T s$, which leads to

$$B^+ = B + \frac{(y - B s)(y - B s)^T}{(y - B s)^T s}$$

called the **symmetric rank-one** (SR1) update

How can we solve $B^+ \Delta x^+ = -\nabla f(x^+)$, in order to take next step? In addition to propagating B to B^+ , let's **propagate inverses**, i.e., $C = B^{-1}$ to $C^+ = (B^+)^{-1}$

Sherman-Morrison formula:

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

Thus for the SR1 update the inverse is also easily updated:

$$C^+ = C + \frac{(s - Cy)(s - Cy)^T}{(s - Cy)^T y}$$

In general, SR1 is simple and cheap, but has key shortcoming: it does not preserve positive definiteness

Broyden-Fletcher-Goldfarb-Shanno update

Instead of a rank-one update to B , let's try a rank-two update

$$B^+ = B + auu^T + bvv^T$$

Using secant equation $B^+s = y$ gives

$$y - Bs = (au^T s)u + (bv^T s)v$$

Setting $u = y$, $v = Bs$ and solving for a, b we get

$$B^+ = B - \frac{Bss^T B}{s^T Bs} + \frac{yy^T}{y^T s}$$

called the **Broyden-Fletcher-Goldfarb-Shanno** (BFGS) update

Woodbury formula (generalization of Sherman-Morrison):

$$(A + UDV)^{-1} = A^{-1} - A^{-1}U(D^{-1} + VA^{-1}U)^{-1}VA^{-1}$$

Applied to our case, with

$$U = V^T = \begin{bmatrix} Bs & y \end{bmatrix}, \quad D = \begin{bmatrix} -1/(s^T Bs) & 0 \\ 0 & 1/(y^T s) \end{bmatrix}$$

then after some algebra we get a rank-two update on C :

$$C^+ = \left(I - \frac{sy^T}{y^T s} \right) C \left(I - \frac{ys^T}{y^T s} \right) + \frac{ss^T}{y^T s}$$

The BFGS update is thus still quite cheap, $O(n^2)$ per update

Positive definiteness of BFGS update

Importantly, unlike SR1, the BFGS update **preserves positive definiteness** under appropriate conditions

Assume $y^T s = (\nabla f(x^+) - \nabla f(x))^T (x^+ - x) > 0$ (recall that e.g. strict convexity will imply this condition) and $C \succ 0$

Then consider the term

$$x^T C^+ x = \left(x - \frac{s^T x}{y^T s} y \right)^T C \left(x - \frac{s^T x}{y^T s} y \right) + \frac{(s^T x)^2}{y^T s}$$

Both terms are nonnegative; second term is only zero when $s^T x = 0$, and in that case first term is only zero when $x = 0$

Davidon-Fletcher-Powell update

Alternatively, compute a rank-two update directly on inverse C

$$C^+ = C + auu^T + bvv^T.$$

Using secant equation $s = C^+y$, setting $u = s$, $v = Cy$, and solving for a, b gives

$$C^+ = C - \frac{Cyy^TC}{y^TCy} + \frac{ss^T}{y^Ts}$$

Called the Davidon-Fletcher-Powell (DFP) update

Pre-dates BFGS, with same beneficial properties (preserves positive definiteness of Hessian, $O(n^2)$ computation), but not often used anymore

Broyden class

SR1, BFGS, and DFP are some of numerous possible quasi-Newton updates. The **Broyden class** of updates is defined by:

$$B^+ = (1 - \phi)B_{\text{BFGS}}^+ + \phi B_{\text{DFP}}^+, \quad \phi \in \mathbb{R}$$

By putting $v = y/(y^T s) - Bs/(s^T Bs)$, we can rewrite the above as

$$B^+ = B - \frac{Bss^T B}{s^T Bs} + \frac{yy^T}{y^T s} + \phi(s^T Bs)vv^T$$

Note:

- BFGS corresponds to $\phi = 0$
- DFP corresponds to $\phi = 1$
- SR1 corresponds to $\phi = y^T s / (y^T s - s^T Bs)$

Convergence analysis

Assume that f convex, twice differentiable, having $\text{dom}(f) = \mathbb{R}^n$, and additionally

- ∇f is Lipschitz with parameter L
- f is strongly convex with parameter m
- $\nabla^2 f$ is Lipschitz with parameter M

(same conditions as in the analysis of Newton's method)

Theorem: Both BFGS and DFP, with backtracking line search, converge globally. Furthermore, for all $k \geq k_0$,

$$\|x^{(k)} - x^*\|_2 \leq c_k \|x^{(k-1)} - x^*\|_2$$

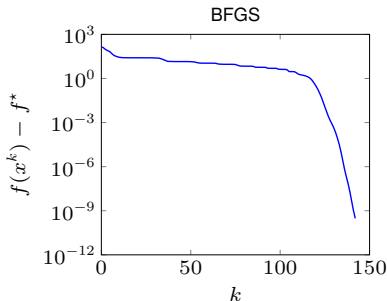
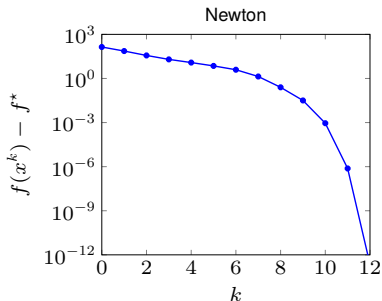
where $c_k \rightarrow 0$ as $k \rightarrow \infty$. Here k_0, c_k depend on L, m, M

This is called **local superlinear convergence**

Example: Newton versus BFGS

Example from Vandenberghe's lecture notes: Newton versus BFGS on LP barrier problem, for $n = 100$, $m = 500$

$$\min_x c^T x - \sum_{i=1}^m \log(b_i - a_i^T x)$$



Recall Newton update is $O(n^3)$, quasi-Newton update is $O(n^2)$.
But quasi-Newton converges in less than 100 times the iterations

Implicit-form quasi-Newton

For large problems, quasi-Newton updates can become too costly

Basic idea: instead of explicitly computing and storing C , compute an implicit version of C by maintaining all pairs (y, s)

Recall BFGS updates C via

$$C^+ = \left(I - \frac{sy^T}{y^T s} \right) C \left(I - \frac{ys^T}{y^T s} \right) + \frac{ss^T}{y^T s}$$

Observe this leads to

$$C^+ g = p + (\alpha - \beta)s, \quad \text{where}$$
$$\alpha = \frac{s^T g}{y^T s}, \quad q = g - \alpha y, \quad p = Cq, \quad \beta = \frac{y^T p}{y^T s}$$

We see that C^+g can be computed via two loops of length k (if C^+ is the approximation to the inverse Hessian after k iterations):

1. Let $q = -\nabla f(x^{(k)})$
2. For $i = k - 1, \dots, 0$:
 - (a) Compute $\alpha_i = (s^{(i)})^T q / ((y^{(i)})^T s^{(i)})$
 - (b) Update $q = q - \alpha y^{(i)}$
3. Let $p = C^{(0)}q$
4. For $i = 0, \dots, k - 1$:
 - (a) Compute $\beta = (y^{(i)})^T p / ((y^{(i)})^T s^{(i)})$
 - (b) Update $p = p + (\alpha_i - \beta)s^{(i)}$
5. Return p

Limited memory BFGS

Limited memory BFGS (LBFGS) simply limits each of these loops to be length m :

1. Let $q = -\nabla f(x^{(k)})$
2. For $i = k - 1, \dots, k - m$:
 - (a) Compute $\alpha_i = (s^{(i)})^T q / ((y^{(i)})^T s^{(i)})$
 - (b) Update $q = q - \alpha y^{(i)}$
3. Let $p = \bar{C}^{(k-m)} q$
4. For $i = k - m, \dots, k - 1$:
 - (a) Compute $\beta = (y^{(i)})^T p / ((y^{(i)})^T s^{(i)})$
 - (b) Update $p = p + (\alpha_i - \beta) s^{(i)}$
5. Return p

In Step 3, $\bar{C}^{(k-m)}$ is our guess at $C^{(k-m)}$ (which is not stored). A popular choice is $\bar{C}^{(k-m)} = I$, more sophisticated choices exist

Stochastic quasi-Newton methods

Consider now the problem

$$\min_x \mathbb{E}_\xi[f(x, \xi)]$$

for a noise variable ξ . Tempting to extend previous ideas and take **stochastic quasi-Newton** updates of the form:

$$x^{(k)} = x^{(k-1)} - t_k C^{(k-1)} \nabla f(x^{(k-1)}, \xi_k)$$

But there are challenges:

- Can have at best sublinear convergence (recall lower bound by Nemirovski et al.) So is additional overhead of quasi-Newton, worth it, over plain SGD?
- Updates to C depend on consecutive gradient estimates; noise in the gradient estimates could be a hindrance

The most straightforward adaptation of quasi-Newton methods is to use BFGS (or LBFGS) with

$$s^{(k-1)} = x^{(k)} - x^{(k-1)}, y^{(k-1)} = \nabla f(x^{(k)}, \xi_k) - \nabla f(x^{(k-1)}, \xi_k)$$

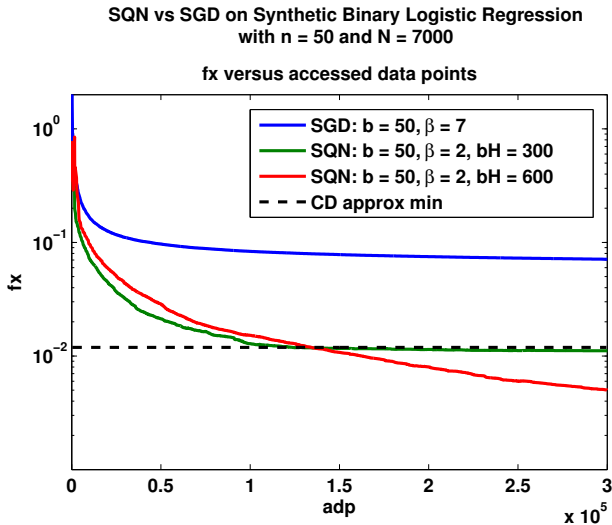
The key is to use the same noise variable ξ_k in the two stochastic gradients. This is due to Schraudolph et al. (2007)

More recently, Byrd et al. (2015) propose a stochastic version of LBFGS with three main changes:

- Perform an LBFGS update only every L iterations
- Compute s to be an average over L last search directions
- Compute y using Hessian approximation based on sampling

With proper tuning, either approach can give improvements over SGD

Example from Byrd et al. (2015):



References and further reading

- L. Bottou, F. Curtis, J. Nocedal (2016), “Optimization methods for large-scale machine learning”
- R. Byrd, S. Hansen, J. Nocedal, Y. Singer (2015), “A stochastic quasi-Newton method for large-scale optimization”
- J. Dennis and R. Schnabel (1996), “Numerical methods for unconstrained optimization and nonlinear equations”
- J. Nocedal and S. Wright (2006), “Numerical optimization”, Chapters 6 and 7
- N. Schraudolph, J. Yu, S. Gunter (2007), “A stochastic quasi-Newton method for online convex optimization”
- L. Vandenberghe, Lecture notes for EE 236C, UCLA, Spring 2011-2012