

QUASI-NEWTON METHODS

David F. Gleich

February 29, 2012

The idea behind Quasi-Newton methods is to make an optimization algorithm *with only a function value and gradient* converge more quickly than steepest descent. That is, a Quasi-Newton method does not require a means to evaluate the Hessian matrix at the current iterate, as in a Newton method. Instead, the algorithm constructs a matrix that resembles the Hessian as it proceeds.

In fact, there are many ways of doing this, and so there is really a family of Quasi-Newton methods.

The material here is from Chapter 6 in Nocedal and Wright, and Section 12.3 in Griva, Sofer, and Nash.

1 QUASI-NEWTON IN ONE VARIABLE: THE SECANT METHOD

In a one dimensional problem, approximating the Hessian simplifies to approximating the second derivative: $f''(x) \approx \frac{f'(x+h) - f'(x)}{h}$. Thus, the fact that this is possible is not unreasonable. Using a related approximation in a one-dimensional optimization algorithm results in a procedure called the *Secant method*:

$$\underbrace{\text{“}x_{k+1} = x_k - \frac{1}{f''(x_k)} f'(x_k)\text{”}}_{\text{One dimensional Newton}} \quad \rightarrow \quad x_{k+1} = x_k - \underbrace{\frac{(x_k - x_{k-1})}{f'(x_k) - f'(x_{k-1})} f'(x_k)}_{\approx 1/f''(x_k)}$$

This new update is trying to approximate the Newton update by approximating the second derivative information.

The secant method converges superlinearly, under appropriate conditions; so this idea checks out in one-dimension.

2 QUASI-NEWTON IN GENERAL

Quasi-Newton methods are line-search methods that compute the search direction by trying to approximate the Newton direction:

$$\text{“}\mathbf{H}(\mathbf{x}_k)\mathbf{p} = -\mathbf{g}\text{”}$$

without using the matrix $\mathbf{H}(\mathbf{x}_k)$. They work by computing

$$\mathbf{B}_k \quad \text{“that behaves like”} \quad \mathbf{H}(\mathbf{x}_k).$$

Once we compute \mathbf{x}_{k+1} , then we update $\mathbf{B}_k \rightarrow \mathbf{B}_{k+1}$. Thus, a Quasi-Newton method has the general iteration:

```
initialize  $\mathbf{B}_0$ , and  $k = 0$ 
for  $k = 0, \dots$  and while  $\mathbf{x}_k$  does not satisfy the conditions we want ...
    solve for the search direction  $\mathbf{B}_k \mathbf{p}_k = -\mathbf{g}$ 
    compute a line search  $\alpha_k$ 
    update  $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \mathbf{p}_k$ 
    update  $\mathbf{B}_{k+1}$  based on  $\mathbf{x}_{k+1}$ 
```

We can derive different Quasi-Newton methods by changing how we update \mathbf{B}_{k+1} from \mathbf{B}_k .

3 THE SECANT CONDITION

While there are many ways of updating \mathbf{B}_{k+1} from \mathbf{B}_k , a random choice is unlikely to provide any benefit, and may making things considerably worse. Thus, we want to start from a principled approach.

Recall that the Newton direction $\mathbf{H}_k \mathbf{p}_k = -\mathbf{g}$ arises as the unconstrained minimizer of

$$m_k^N(\mathbf{p}) = f_k + \mathbf{g}_k^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{H}_k \mathbf{p}$$

when \mathbf{H}_k is positive definite.

The model for Quasi-Newton methods uses \mathbf{B}_k instead of \mathbf{H}_k :

$$m_k^Q(\mathbf{p}) = f_k + \mathbf{g}_k^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{B}_k \mathbf{p}$$

so one common requirement for \mathbf{B}_k is that it remains positive definite. This requirement is relaxed for some Quasi-Newton methods.

However, all Quasi-Newton methods require:

$$\nabla m_{k+1}^Q(0) = \mathbf{g}(\mathbf{x}_{k+1})$$

and

$$\nabla m_{k+1}^Q(-\alpha_k \mathbf{p}_k) = \mathbf{g}(\mathbf{x}_k).$$

In other words, a Quasi-Newton method has the property that the gradient of the model function $m_{k+1}^Q(\mathbf{p})$ has the same gradient as f at \mathbf{x}_k and \mathbf{x}_{k+1} .

This requirement imposes some conditions on \mathbf{B}_{k+1} :

$$\nabla m_{k+1}^Q(-\alpha_k \mathbf{p}_k) = \mathbf{g}(\mathbf{x}_{k+1}) - \alpha_k \mathbf{B}_{k+1} \mathbf{p}_k = \mathbf{g}(\mathbf{x}_k) \quad \longrightarrow \quad \mathbf{B}_{k+1} \alpha_k \mathbf{p}_k = \mathbf{g}(\mathbf{x}_{k+1}) - \mathbf{g}(\mathbf{x}_k).$$

Note that $\alpha_k \mathbf{p}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$. If we define:

$$\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k \quad \text{and} \quad \mathbf{y}_k = \mathbf{g}(\mathbf{x}_{k+1}) - \mathbf{g}(\mathbf{x}_k).$$

Then Quasi-Newton methods require:

$$\mathbf{B}_{k+1} \mathbf{s}_k = \mathbf{y}_k,$$

which is called the *secant condition*.

If we write this out for a one-dimensional problem:

$$b_{k+1}(x_{k+1} - x_k) = f'(x_{k+1}) - f'(x_k).$$

This equation is identical to the approximation of $f''(x_k)$ used in the secant method.

Quiz Is it always possible to find such a \mathbf{B}_{k+1} ? Suppose that \mathbf{B}_k is symmetric, positive definite. Show that we need $\mathbf{y}_k^T \mathbf{x}_k > 0$ in order for \mathbf{B}_{k+1} to be positive definite. If $\mathbf{B}_k = 1$ for a one dimensional problem, find a function where this isn't true.

4 FINDING THE UPDATE

We are getting closer to figuring out how to find such an update. There are many ways to derive the following updates, I'll just list them and state their properties.

4.1 DAVIDSON, FLETCHER, POWELL (DFP)

$$\text{Let } \rho = \frac{1}{\mathbf{y}_k^T \mathbf{s}_k}.$$

$$\mathbf{B}_{k+1} = (\mathbf{I} - \rho_k \mathbf{s} \mathbf{y}^T) \mathbf{B}_k (\mathbf{I} - \rho_k \mathbf{s} \mathbf{y}^T) + \rho_k \mathbf{y}_k \mathbf{y}_k^T.$$

Clearly this matrix is symmetric when \mathbf{B}_k is. Also, \mathbf{B}_{k+1} is positive definite.

Quiz Show that \mathbf{B}_{k+1} is positive definite.

This choice of \mathbf{B}_{k+1} has the following optimality property:

$$\begin{aligned} &\text{minimize} && \|\mathbf{B} - \mathbf{B}_k\|_W \\ &\text{subject to} && \mathbf{B}^T = \mathbf{B}, \mathbf{B} \mathbf{s}_k = \mathbf{y}_k \end{aligned}$$

where \mathbf{W} is a weight based on the average Hessian.

4.2 BROYDEN, FLETCHER, GOLDFARB, SHANNO (BFGS) – “STANDARD”

Because we compute the search direction by solving a system with the approximate Hessian matrix:

$$\mathbf{B}_k \mathbf{p}_k = -\mathbf{g}_k,$$

the BFGS update constructs an approximation of the *inverse Hessian* instead. Suppose that

$$\mathbf{T}_k \text{ “behaves like” } \mathbf{H}(\mathbf{x})^{-1}.$$

Then

$$\mathbf{T}_{k+1} \mathbf{y}_k = \mathbf{s}_k$$

is the secant condition for the inverse. This helps because now we can find search directions via

$$\mathbf{p}_k = -\mathbf{T}_k \mathbf{g}_k,$$

via a matrix-vector multiplication instead of a linear solve.

The BFGS method uses the update:

$$\mathbf{T}_{k+1} = (\mathbf{I} - \rho_k \mathbf{s} \mathbf{y}^T) \mathbf{T}_k (\mathbf{I} - \rho_k \mathbf{s} \mathbf{y}^T) + \rho_k \mathbf{s}_k \mathbf{s}_k^T.$$

By the same proof, this update is also positive definite.

This choice has the following optimality property:

$$\begin{aligned} & \text{minimize} && \|\mathbf{T} - \mathbf{T}_k\|_W \\ & \text{subject to} && \mathbf{T}^T = \mathbf{T}, \mathbf{T} \mathbf{y}_k = \mathbf{s}_k \end{aligned}$$

where \mathbf{W} is a weight based on the average Hessian.

4.3 SYMMETRIC RANK-1 (SR1) – FOR TRUST REGION METHODS

Both of the previous updates were rank-2 changes to \mathbf{B}_k (or \mathbf{T}_k). The SR1 method is a rank-1 update to \mathbf{B}_k . Unfortunately, this update will not preserve positive definiteness. Nonetheless, it’s frequently used in practice and is a reasonable choice for Trust Region methods that don’t require a positive definite approximate Hessian.

Any rank-1 symmetric matrix is:

$$\sigma \mathbf{v} \mathbf{v}^T$$

and so the update is:

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \sigma \mathbf{v} \mathbf{v}^T.$$

Applying the Secant equation constrains \mathbf{v} , and we have:

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \frac{(\mathbf{y}_k - \mathbf{B}_k \mathbf{s}_k)(\mathbf{y}_k - \mathbf{B}_k \mathbf{s}_k)^T}{(\mathbf{y}_k - \mathbf{B}_k \mathbf{s}_k)^T \mathbf{s}_k}$$

or

$$\mathbf{T}_{k+1} = \mathbf{T}_k + \frac{(\mathbf{s}_k - \mathbf{T}_k \mathbf{y}_k)(\mathbf{s}_k - \mathbf{T}_k \mathbf{y}_k)^T}{(\mathbf{s}_k - \mathbf{T}_k \mathbf{y}_k)^T \mathbf{y}_k}.$$

The SR1 method tends to generate better approximations to the true Hessian than the other methods. For instance, if the search directions \mathbf{p}_k are all linearly independent for $k = 1, \dots, n$, and $f(\mathbf{x})$ is a simple quadratic model, then \mathbf{T}_n is the inverse of the true Hessian.

4.4 BROYDEN CLASS

The Broyden class is a linear combination of the BFGS and the DFP method:

$$\mathbf{B}_{k+1} = (1 - \phi) \mathbf{B}_{k+1}^{\text{BFGS}} + \phi \mathbf{B}_{k+1}^{\text{DFP}}.$$

(This form requires the BFGS update for \mathbf{B} and not \mathbf{T} .)

There are all sorts of great properties of the Broyden class, e.g. for the right choice of parameters, it’ll reproduce the CG method.