

# Quasi-Newton method

**Quasi-Newton methods** are methods used to either find zeroes or local maxima and minima of functions, as an alternative to Newton's method. They can be used if the Jacobian or Hessian is unavailable or is too expensive to compute at every iteration. The "full" Newton's method requires the Jacobian in order to search for zeros, or the Hessian for finding extrema.

## Contents

- Search for zeros: root finding
- Search for extrema: optimization
- Relationship to matrix inversion
- Notable implementations
- See also
- References
- Further reading

In vector calculus, the **Jacobian matrix** of a vector-valued function in several variables is the matrix of all its first-order partial derivatives. When this matrix is square, that is, when the function takes the same number of variables as input as the number of



## Search for zeros: root finding

Newton's method to find zeroes of a function *g* of multiple variables is given by  $x_{n+1} = x_n - [J_g(x_n)]^{-1}g(x_n)$ , where  $[J_g(x_n)]^{-1}$  is the left inverse of the Jacobian matrix *J<sub>g</sub>(x<sub>n</sub>)* of *g* evaluated for *x<sub>n</sub>*.

Strictly speaking, any method that replaces the exact Jacobian *J<sub>g</sub>(x<sub>n</sub>)* with an approximation is a quasi-Newton method.<sup>[1]</sup> For instance, the chord method (where *J<sub>g</sub>(x<sub>n</sub>)* is replaced by *J<sub>g</sub>(x<sub>0</sub>)* for all iterations) is a simple example. The methods given below for optimization refer to an important subclass of quasi-Newton methods, secant methods.<sup>[2]</sup>

Using methods developed to find extrema in order to find zeroes is not always a good idea, as the majority of the methods used to find extrema require that the matrix that is used is symmetrical. While this holds in the context of the search for extrema, it rarely holds when searching for zeroes. Broyden's "good" and "bad" methods are two methods commonly used to find extrema that can also be applied to find zeroes. Other methods that can be used are the column-updating method, the inverse column-updating method, the quasi-Newton least squares method and the quasi-Newton inverse least squares method.

More recently quasi-Newton methods have been applied to find the solution of multiple coupled systems of equations (e.g. fluid–structure interaction problems or interaction problems in physics). They allow the solution to be found by solving each constituent system separately (which is simpler than the global system) in a cyclic, iterative fashion until the solution of the global system is found.<sup>[2][3]</sup>

## Search for extrema: optimization

Noting that the search for a minimum or maximum of a scalar-valued function is nothing else than the search for the zeroes of the gradient of that function, quasi-Newton methods can be readily applied to find extrema of a function. In other words, if *g* is the gradient of *f*, then searching for the zeroes of the vector-valued function *g* corresponds to the search for the extrema of the scalar-valued function *f*; the Jacobian

of  $g$  now becomes the Hessian of  $f$ . The main difference is that the Hessian matrix is a symmetric matrix, unlike the Jacobian when searching for zeroes. Most quasi-Newton methods used in optimization exploit this property.

In optimization, **quasi-Newton methods** (a special case of **variable-metric methods**) are algorithms for finding local maxima and minima of functions. Quasi-Newton methods are based on Newton's method to find the stationary point of a function, where the gradient is 0. Newton's method assumes that the function can be locally approximated as a quadratic in the region around the optimum, and uses the first and second derivatives to find the stationary point. In higher dimensions, Newton's method uses the gradient and the Hessian matrix of second derivatives of the function to be minimized.

In quasi-Newton methods the Hessian matrix does not need to be computed. The Hessian is updated by analyzing successive gradient vectors instead. Quasi-Newton methods are a generalization of the secant method to find the root of the first derivative for multidimensional problems. In multiple dimensions the secant equation is under-determined, and quasi-Newton methods differ in how they constrain the solution, typically by adding a simple low-rank update to the current estimate of the Hessian.

The first quasi-Newton algorithm was proposed by William C. Davidon, a physicist working at Argonne National Laboratory. He developed the first quasi-Newton algorithm in 1959: the DFP updating formula, which was later popularized by Fletcher and Powell in 1963, but is rarely used today. The most common quasi-Newton algorithms are currently the SR1 formula (for "symmetric rank-one"), the BHHH method, the widespread BFGS method (suggested independently by Broyden, Fletcher, Goldfarb, and Shanno, in 1970), and its low-memory extension L-BFGS. The Broyden's class is a linear combination of the DFP and BFGS methods.

The SR1 formula does not guarantee the update matrix to maintain positive-definiteness and can be used for indefinite problems. The Broyden's method does not require the update matrix to be symmetric and is used to find the root of a general system of equations (rather than the gradient) by updating the Jacobian (rather than the Hessian).

One of the chief advantages of quasi-Newton methods over Newton's method is that the Hessian matrix (or, in the case of quasi-Newton methods, its approximation)  $B$  does not need to be inverted. Newton's method, and its derivatives such as interior point methods, require the Hessian to be inverted, which is typically implemented by solving a system of linear equations and is often quite costly. In contrast, quasi-Newton methods usually generate an estimate of  $B^{-1}$  directly.

As in Newton's method, one uses a second-order approximation to find the minimum of a function  $f(x)$ . The Taylor series of  $f(x)$  around an iterate is

$$f(x_k + \Delta x) \approx f(x_k) + \nabla f(x_k)^T \Delta x + \frac{1}{2} \Delta x^T B \Delta x,$$

where  $(\nabla f)$  is the gradient, and  $B$  an approximation to the Hessian matrix. The gradient of this approximation (with respect to  $\Delta x$ ) is

$$\nabla f(x_k + \Delta x) \approx \nabla f(x_k) + B \Delta x,$$

and setting this gradient to zero (which is the goal of optimization) provides the Newton step:

$$\Delta x = -B^{-1} \nabla f(x_k).$$

The Hessian approximation  $B$  is chosen to satisfy

$$\nabla f(x_k + \Delta x) = \nabla f(x_k) + B \Delta x,$$

5/20/2020

Quasi-Newton method - Wikipedia

which is called the *secant equation* (the Taylor series of the gradient itself). In more than one dimension ***B*** is underdetermined. In one dimension, solving for ***B*** and applying the Newton's step with the updated value is equivalent to the secant method. The various quasi-Newton methods differ in their choice of the solution to the secant equation (in one dimension, all the variants are equivalent). Most methods (but with exceptions, such as Broyden's method) seek a symmetric solution ( $B^T = B$ ); furthermore, the variants listed below can be motivated by finding an update  $B_{k+1}$  that is as close as possible to  $B_k$  in some norm; that is,  $B_{k+1} = \operatorname{argmin}_B \|B - B_k\|_V$ , where  $V$  is some positive-definite matrix that defines the norm. An approximate initial value  $B_0 = \beta I$  is often sufficient to achieve rapid convergence, although there is no general strategy to choose  $\beta$  <sup>[4]</sup>. Note that  $B_0$  should be positive-definite. The unknown  $x_k$  is updated applying the Newton's step calculated using the current approximate Hessian matrix  $B_k$ :

- $\Delta x_k = -\alpha_k B_k^{-1} \nabla f(x_k)$ , with  $\alpha$  chosen to satisfy the Wolfe conditions;
- $x_{k+1} = x_k + \Delta x_k$ ;
- The gradient computed at the new point  $\nabla f(x_{k+1})$ , and

$$y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$$

is used to update the approximate Hessian  $B_{k+1}$ , or directly its inverse  $H_{k+1} = B_{k+1}^{-1}$  using the Sherman–Morrison formula.

- A key property of the BFGS and DFP updates is that if  $B_k$  is positive-definite, and  $\alpha_k$  is chosen to satisfy the Wolfe conditions, then  $B_{k+1}$  is also positive-definite.

The most popular update formulas are:

Method	$B_{k+1} =$	$H_{k+1} = B_{k+1}^{-1} =$
<u>BFGS</u>	$B_k + \frac{y_k y_k^T}{y_k^T \Delta x_k} - \frac{B_k \Delta x_k (B_k \Delta x_k)^T}{\Delta x_k^T B_k \Delta x_k}$	$\left(I - \frac{\Delta x_k y_k^T}{y_k^T \Delta x_k}\right) H_k \left(I - \frac{y_k \Delta x_k^T}{y_k^T \Delta x_k}\right) + \frac{\Delta x_k \Delta x_k^T}{y_k^T \Delta x_k}$
<u>Broyden</u>	$B_k + \frac{y_k - B_k \Delta x_k}{\Delta x_k^T \Delta x_k} \Delta x_k^T$	$H_k + \frac{(\Delta x_k - H_k y_k) \Delta x_k^T H_k}{\Delta x_k^T H_k y_k}$
Broyden family	$(1 - \varphi_k) B_{k+1}^{\text{BFGS}} + \varphi_k B_{k+1}^{\text{DFP}}, \quad \varphi \in [0, 1]$	
<u>DFP</u>	$\left(I - \frac{y_k \Delta x_k^T}{y_k^T \Delta x_k}\right) B_k \left(I - \frac{\Delta x_k y_k^T}{y_k^T \Delta x_k}\right) + \frac{y_k y_k^T}{y_k^T \Delta x_k}$	$H_k + \frac{\Delta x_k \Delta x_k^T}{\Delta x_k^T y_k} - \frac{H_k y_k y_k^T H_k}{y_k^T H_k y_k}$
<u>SR1</u>	$B_k + \frac{(y_k - B_k \Delta x_k)(y_k - B_k \Delta x_k)^T}{(y_k - B_k \Delta x_k)^T \Delta x_k}$	$H_k + \frac{(\Delta x_k - H_k y_k)(\Delta x_k - H_k y_k)^T}{(\Delta x_k - H_k y_k)^T y_k}$

Other methods are Pearson's method, McCormick's method, the Powell symmetric Broyden (PSB) method and Greenstadt's method.<sup>[2]</sup>

## Relationship to matrix inversion

When ***f*** is a convex quadratic function with positive-definite Hessian ***B***, one would expect the matrices ***H<sub>k</sub>*** generated by a quasi-Newton method to converge to the inverse Hessian ***H*** = ***B***<sup>−1</sup>. This is indeed the case for the class of quasi-Newton methods based on least-change updates.<sup>[5]</sup>

## Notable implementations