

Unconstrained Optimization V: Quasi-Newton methods (continued)

Note Title

4/23/2022

The BFGS and related Quasi-Newton Methods



William C. Davidon

Broyden, Fletcher, Goldfarb, Shanno



Read the historical remarks at the beginning of CH 6 of NBW, it's interesting.

In UO_1, I motivated quasi-Newton method as an analog of the secant method in 1-D for solving equations, giving rise to the secant equation

$$B_p(x_p - x_{p-1}) = \nabla f_p - \nabla f_{p-1}. \quad - (\text{Sec})$$

Let's recall the rationale for (sec) : If we want to solve the system of nonlinear eqs $F(x) = 0$, $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$, and don't want to compute derivatives. Think :

$$F(x) \approx F(x_k) + \cancel{DF(x_k)}(x-x_k) = L_k(x), \text{ and}$$

require B_k to satisfy $L_k(x_{k+1}) = F(x_{k+1})$, which is equivalent to

$$F(x_k) + B_k(x_{k+1} - x_k) = F(x_{k+1}), \text{ or } B_k(x_k - x_{k+1}) = F(x_k) - F(x_{k+1}).$$

↑
this is (sec) when $F = \nabla f$.

In the context of solving equation, B_k shouldn't be constrained to be positive definite or even symmetric. But in the optimization context, a line search of the form:

$$x_{k+1} = x_k - \alpha_k B_k^{-1} \nabla f(x_k)$$

can guarantee a decrease in the objective if

(i) $B_k > 0$, (ii) α_k is chosen to satisfy Armijo sufficient decrease cond.

[Note : $-B_k^{-1} \nabla f_k$ is also the unique minimizer of $m_k(p) := f_k + \nabla f_k^T p + \frac{1}{2} p^T B_k^{-1} p$.
 Indeed, when $\nabla^2 f(x_k) \succ 0$, it is a little strange to think of modelling the objective near x_k by a convex quadratic.]

$$\text{Change } k \text{ to } k+1 \text{ in (sec): } B_{k+1} \underbrace{(x_{k+1} - x_k)}_{=: s_k} = \underbrace{\nabla f(x_{k+1}) - \nabla f(x_k)}_{=: y_k} \quad - (\text{sec}_{k+1})$$

When f non-convex, it is by no means obvious that $\exists B_{k+1} > 0$ st. (sec_{k+1}) is satisfied.

(If f is such that $\nabla^2 f(x) \succ 0 \ \forall x$, then
 $\nabla f(x) - \nabla f(y) = \nabla^2 f(\xi)(x-y)$ for some $\xi \in \overline{xy}$.
so $(x-y)^T(\nabla f(x) - \nabla f(y)) > 0 \ \forall x \neq y$.

Then

$$\exists B > 0 \text{ st. } B(x-y) = \nabla f(x) - \nabla f(y) \quad (\text{your HW #1}).$$

NOTE: $\exists B_{k+1} > 0$ s.t. $B_{k+1}s_k = y_k \Leftrightarrow s_k^T y_k > 0$

\Rightarrow multiply s_k^T to both sides
 \Leftarrow your HW #1

N&W call $s_k^T y_k > 0$ the **curvature condition**, which is confusing, considered the eponymous second Wolfe condition $\nabla f_{k+1}^T p_k \geq c_2 \nabla f_k^T p_k$.

Here's the surprise: If the step size α_k in $x_{k+1} = x_k - \alpha_k B_k^{-1} \nabla f(x_k)$ is chosen st. it satisfies the second Wolfe condition, then $s_k^T y_k > 0$ is satisfied.

$$\begin{aligned}
 \text{The proof is easy though: } & \nabla f_{k+1}^T p_k \geq c_2 \nabla f_k^T p_k \\
 \Rightarrow & \nabla f_{k+1}^T (\alpha_k p_k) \geq c_2 \nabla f_k^T (\alpha_k p_k) \\
 \Rightarrow & \underbrace{(\nabla f_{k+1} - \nabla f_k)^T}_{y_k} \underbrace{(\alpha_k p_k)}_{s_k} \geq (c_2 - 1) \nabla f_k^T (\alpha_k p_k) > 0.
 \end{aligned}$$

So if x_{k+1} is formed by $x_{k+1} = x_k - \alpha_k B_k^{-1} \nabla f_k$ with α_k satisfying the second Wolfe condition, then we can always find $B_{k+1} > 0$ that satisfies the secant equation

$$B_{k+1} s_k = y_k \quad \leftarrow n \text{ linear constraints}$$

\uparrow
 $n(n+1)/2$ degrees of freedom

So there is a $\frac{n(n+1)}{2} - n = \frac{n(n-1)}{2}$ dimensional space of solutions. (The positive-definiteness constraint does not absorb these d.o.f.)

To determine B_{k+1} uniquely and **cheaply**, find B_{k+1} to be a solution close to B_k in some sense, two approaches have been proposed and studied

(i) Find B_{k+1} that solves : $\min_B \|B - B_k\|$ st. $B = B^T$, $B s_k = y_k$, $B > 0$.

for some matrix norm $\|\cdot\|$.

(ii) Set $B_{R+1} = B_R + (\text{low rank update})$ st. $B_{R+1} = B_{R+1}^T$, $B_{R+1} s_R = y_R$,
 $B_{R+1} > 0$

To begin, note that any rank 1 update of B_R that satisfies the secant condition must be st.

$$B_{R+1} = B_R + uC^T, \quad (B_R + uC^T)s_R = y_R$$
$$\Downarrow$$
$$uC^T s_R = y_R - B_R s_R \Rightarrow u = \frac{y_R - B_R s_R}{C^T s_R}$$

$$\text{so } B_{R+1} = B_R + (y_R - B_R s_R)C^T / C^T s_R.$$

C has to chosen st. $C^T s_R \neq 0$

Also, C should be chosen so that B_{R+1} is symmetric when B_R is symmetric.

Easy to see that $C = y_R - B_R s_R$, i.e.

$$B_{R+1} = B_R + (y_R - B_R s_R)(y_R - B_R s_R)^T / (y_R - B_R s_R)^T s_R \quad -(SR1)$$

is the only solution provided $(Y_k - B_k S_k)^T S_k \neq 0$.

If $Y_k = B_k S_k$, then $B_{k+1} = B_k$ is the only solution.

If $Y_k \neq B_k S_k$ but $(Y_k - B_k S_k)^T S_k = 0$, then there is no solution.

} See N&W P144

This method, known as **symmetric-rank-1**, or SR1, has the shortcoming of breaking down in the last case (and related numerical instabilities).

Also, there is no way to guarantee $B_{k+1} > 0$ (even when $B_k > 0$.)

But, somewhat surprisingly, with a simple safeguard to adequately prevent the breakdown and numerical instabilities, the SR1 method is actually useful in problems with constraints. (More on this later.)

See Sec 6.2 for more details on SR1.

Since rank-1 update cannot fully solves the problem, researchers tried rank 2.

M. Powell had the following idea (1970) :

(For notational simplicity, drop the subscripts k . Write B_k as B , B_{k+1} as \bar{B} .
 Y_k as y , S_k as s .)

$$\bar{B}^{(1)} = B + (y - Bs)C^T / C^T s$$

$$\bar{B}^{(2)} = \frac{1}{2} (\bar{B}^{(1)} + \bar{B}^{(1)T})$$

(satisfies secant cond., not symmetric)

(symmetrized, but doesn't satisfy secant cond.)

Iterate: $\bar{B}^{(2j+1)} = \bar{B}^{(2j)} + (y - \bar{B}^{(2j)}s)C^T / C^T s$

$$\bar{B}^{(2j+2)} = \frac{1}{2} (\bar{B}^{(2j+1)} + \bar{B}^{(2j+1)T})$$

It turns out:

$\{\bar{B}^{(j)}\}$ has a limit given by

$$\bar{B} = B + \frac{(y - Bs)C^T + C(y - Bs)^T}{C^T s} - \frac{(y - Bs)^T s}{(C^T s)^2} CC^T \quad (P)$$

It looks like a rank-3 update but is actually a rank-2 update. (why?)
Again, C has to be chosen s.t. $C^T s \neq 0$.

I'll omit the proof of (P). It is easy to check that \bar{B} in (P) satisfies both
(i) the secant condition $\bar{B}s = y$ and (ii) B symmetric $\Rightarrow \bar{B}$ symmetric.

Ex: check (i). (ii) is trivial.)

Michael Powell
FRS FAA



The low-rank update approach is connected to the "closest" update approach by the following result:

Thm: Let $B \in \mathbb{R}^{n \times n}$, $B = B^T$, $C, S, Y \in \mathbb{R}^n$ with $C^T S \neq 0$.

\bar{B} defined by (P) is the unique solution to :

$$\min \|M(\hat{B} - B)M\|_F \text{ s.t. } \hat{B} = \hat{B}^T, \hat{B}S = Y,$$

where M is any non-singular symmetric matrix s.t. $Mc = M^{-1}S$.

$$\text{Proof: (P)} \Rightarrow \bar{B} - B = \frac{(Y - BS)C^T + C(Y - BS)^T}{C^T S} - \frac{(Y - BS)^T S}{(C^T S)^2} CC^T$$

Take any symmetric \hat{B} with $\hat{B}S = Y$, pre- and post-multiply the above by M , and write $Z := Mc = M^{-1}S$. Then

$$\begin{aligned} M(\bar{B} - B)M &= \frac{M(\hat{B} - B)S C^T M + M_c((\hat{B} - B)S)^T M}{Z^T Z} - \frac{((\hat{B} - B)S)^T S}{(Z^T Z)^2} M C C^T M \\ &= \frac{M(\hat{B} - B)M M^{-1} S C^T M + M_c S^T M^T M (\hat{B} - B)M}{Z^T Z} - \frac{S^T M^T M (\hat{B} - B) M M^{-1} S}{(Z^T Z)^2} M C C^T M \end{aligned}$$

$$\overline{E} = \underbrace{M(\bar{B} - B)M}_{\bar{z}^T \bar{z}} = \underbrace{M(\hat{B} - B)M}_{\bar{z}^T \bar{z}} \underbrace{z z^T + z z^T M(\hat{B} - B)M}_{\bar{z}^T \bar{z}} - \frac{\bar{z}^T M(\hat{B} - B)M \bar{z}}{(\bar{z}^T \bar{z})^2} z z^T$$

$$\bar{E} = \frac{E z z^T + z z^T E}{z^T z} - \frac{z^T E z}{(z^T z)^2} z z^T$$

We have

$$v \perp z \Rightarrow \bar{E} v = \frac{z z^T}{z^T z} E v \Rightarrow \|\bar{E} v\|_2 \leq \|E v\|_2$$

so $\|\bar{E}\|_F \leq \|E\|_F$, as desired.

(Ex: Explain this.)

The uniqueness follows from the strict convexity of $f: \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$, $f(A) = \|M(A-B)M\|_F$ restricted to the affine subspace (hence convex) $\{\hat{B}: \hat{B}^T = \hat{B}, \hat{B}S = y\}$.

Q.E.D.

Remark:

One may directly solve the minimization problem in the theorem to retrieve the update formula (P).

An obvious choice of C is $C=S$ (so $M=I$ above), this gives the so-called Powell Symmetric Broyden update

$$\bar{B}_{PSB} = B + \frac{(y - Bs)S^T + S(y - Bs)^T}{S^T S} - \frac{(y - Bs)^T S}{(S^T S)^2} S S^T.$$

Problem with this formula : (i) $B > 0 \not\Rightarrow \bar{B} > 0$, (ii) lacks affine/scale invariance.

Surprisingly, if we choose $C=y$, the resulted update formula — credited to Davidon, Fletcher and Powell — enjoys both properties.

A simple twist of DFP resulted in the BFGS formula, which enjoys the same two properties, on top of being more effective than DFP.

See below. The development is a little long, but it's good to go through it.

Analyzing the $B \succ 0 \Rightarrow \bar{B} \succ 0$ property

$$\begin{aligned} \text{Recall } \bar{B} &= B + \underbrace{\frac{(Y - BS)C^T + C(Y - BS)^T}{C^T S}}_{\substack{1st \\ 3rd}} - \underbrace{\frac{(Y - BS)^T S}{(C^T S)^2} C C^T}}_{\substack{2nd \\ 3rd}} \\ &= \underbrace{PC^T}_{\substack{1st \\ 3rd}} + \underbrace{CQ^T}_{\substack{2nd}}, \text{ or } \underbrace{\tilde{P}C^T}_{\substack{1st}} + \underbrace{\tilde{C}\tilde{Q}^T}_{\substack{2nd \\ 3rd}}, \text{ or } \underbrace{VC^T}_{\substack{1st \\ 3rd}} + \underbrace{CV^T}_{\substack{2nd \\ 3rd}} \\ &\quad \frac{Y - BS}{C^T S} - \frac{1}{2} \frac{(Y - BS)^T S}{(C^T S)^2} \end{aligned}$$

Thm : If $B \succ 0$, then $\bar{B} \succ 0 \Leftrightarrow \det \bar{B} > 0$.

Proof : Of course (\Rightarrow) is always true. The converse has to do with rank 2.

$$(\Leftarrow) : \bar{B} = B + VC^T + CV^T = B + \frac{1}{2} [(V+C)(V+C)^T - (V-C)(V-C)^T]$$

Recall the eigenvalue interlacing thm at the end of UO-3 :

Lemma : If A is symmetric with eigenvalues $\lambda_1 \leq \dots \leq \lambda_n$, and $A^* = A + \sum_{i=1}^n \lambda_i u_i u_i^T$, then A^* has eigenvalues λ_i^* st. $\lambda_1 \leq \lambda_1^* \leq \lambda_2 \leq \dots \leq \lambda_n \leq \lambda_n^*$.

See, e.g., Wilkinson "The Algebraic Eigenvalue Problem" (1965)

Applying this result twice to \bar{B} shows that \bar{B} can have at most one non-positive eigenvalue.

But $0 < \det \bar{B} = \text{product of eigenvalues of } \bar{B}$,
so all the eigenvalues of \bar{B} must be positive, so $\bar{B} > 0$. Q.E.D.

To see what choices of C yield $\det \bar{B} > 0$, we need an expression for $\det \bar{B}$.

Lemma: For $v, w \in \mathbb{R}^n$, $\det(I + vw^T) = 1 + v^T w$.

Lemma (Sherman-Morrison) Let $u, v \in \mathbb{R}^n$, A non-singular. Then
 $A + uv^T$ is non-singular $\iff \sigma := 1 + v^T A^{-1} u \neq 0$.
If $\sigma \neq 0$, then
 $(A + uv^T)^{-1} = A^{-1} - \frac{1}{\sigma} A^{-1} u v^T A^{-1}$.

(See Dennis-More' 77 SIAM Rev paper for proofs, they are not hard.)

Lemma: Let $u_i \in \mathbb{R}^n$, $i=1, 2, 3, 4$. Then

$$\det(I + u_1 u_2^T + u_3 u_4^T) = (1 + u_1^T u_2)(1 + u_3^T u_4) - u_1^T u_4 u_2^T u_3.$$

Proof: Assume $u_1^T u_2 \neq -1$, so $I + u_1 u_2^T$ is nonsingular (by lemma above), and

$$I + u_1 u_2^T + u_3 u_4^T = (I + u_1 u_2^T) (I + (I + u_1 u_2^T)^{-1} u_3 u_4^T)$$

Then

$$\begin{aligned}
 \det(I + u_1 u_2^T + u_3 u_4^T) &= \det(I + u_1 u_2^T) \det(I + \underbrace{(I + u_1 u_2^T)^{-1} u_3 u_4^T}_{I - \frac{u_1 u_2^T}{1 + u_1^T u_2}}) \\
 &= (1 + u_1^T u_2) \det(I + u_3 u_4^T - \underbrace{\frac{u_1 u_2^T u_3 u_4^T / (1 + u_1^T u_2)}{1 + u_1^T u_2}}_{\frac{u_2^T u_3}{1 + u_1^T u_2} u_1 u_4^T}) \\
 &\quad \underbrace{\qquad\qquad\qquad}_{I + \left(u_3 - \frac{u_2^T u_3}{1 + u_1^T u_2} u_1\right) u_4^T} \\
 &= (1 + u_1^T u_2) \left(1 + \left(u_3 - \frac{u_2^T u_3}{1 + u_1^T u_2} u_1\right)^T u_4\right) \\
 &= (1 + u_1^T u_2) \left(1 + u_3^T u_4 - \frac{u_2^T u_3 u_1^T u_4}{1 + u_1^T u_2}\right) \\
 &= (1 + u_1^T u_2) (1 + u_3^T u_4) - (u_2^T u_3)(u_1^T u_4).
 \end{aligned}$$

Since the result holds for $u_1^T u_2 \neq -1$, a continuity argument shows that it holds in general.

/

Q.E.D.

$$\begin{aligned} \text{Now the above formula to } \bar{B} &= B + \underbrace{\frac{(y - Bs)c^T + c(y - Bs)^T}{c^T s}}_{\text{rank 2}} - \underbrace{\frac{(y - Bs)^T s}{(c^T s)^2} c c^T}_{\text{rank 1}} \\ &= B [I + B^{-1} \begin{bmatrix} \text{rank 2} \end{bmatrix}] \end{aligned}$$

With some algebra, using the last lemma (tricky):

$$\det \bar{B} = \det B [(c^T B^{-1} y)^2 - c^T B^{-1} c y^T B^{-1} y + c^T B^{-1} c y^T s] / (c^T s)^2.$$

Since we assume $B > 0$, \sqrt{B} exists. Write $v := B^{-\frac{1}{2}} y$, $w = B^{-\frac{1}{2}} c$, then

$$\det \bar{B} = \det B [(v^T w)^2 - \|v\|^2 \|w\|^2 + \|w\|^2 y^T s] / \|c^T s\|^2,$$

$$\begin{aligned} \text{So } \bar{B} > 0 \Leftrightarrow \det \bar{B} > 0 \Leftrightarrow \|w\|^2 y^T s &> \underbrace{\|v\|^2 \|w\|^2 - (v^T w)^2}_{\geq 0 \text{ by Cauchy-Schwarz}} \\ &= 0 \text{ if } w \text{ is a multiple of } v \end{aligned}$$

Recall that the (two different, subtly related) curvature conditions guarantee $y^T s > 0$,
so

$\bar{B} > 0$ is guaranteed if w is a multiple of v
 $\Leftrightarrow c$ is a multiple of y

This leads to the DFP update

$$\begin{aligned}
 \bar{B}_{DFP} &= B + \frac{(y - Bs)y^T + y(y - Bs)^T}{y^T s} - \frac{(y - Bs)^T s}{(y^T s)^2} yy^T \\
 &= B + \frac{y^T s [2yy^T - Bs y^T - y s^T B] - y^T s yy^T + (s^T Bs) \cancel{yy^T}}{(y^T s)^2} \\
 &= B - \frac{Bs y^T}{y^T s} - \frac{y s^T B}{y^T s} + \frac{y s^T B s y^T}{(y^T s)^2} + \frac{yy^T}{y^T s} \\
 &= (I - \frac{y s^T}{y^T s}) B (I - \frac{s y^T}{y^T s}) + \frac{yy^T}{y^T s}
 \end{aligned}$$

Applying Sherman-Morrison twice to get:

$$\bar{B}_{DFP}^{-1} = B^{-1} + \underbrace{\frac{ss^T}{s^T y}}_{\text{also rank-2!}} - \frac{B^{-1} yy^T B^{-1}}{y^T B^{-1} y}$$

Ex : Prove this formula.

(Do not use the last expression of \bar{B}_{DFP} . $I - \frac{y s^T}{y^T s}$ and $I - \frac{s y^T}{y^T s}$ are singular!)

Again : $B > 0$ and $y^T s > 0 \Rightarrow \bar{B}_{DFP}, \bar{B}_{DFP}^{-1} > 0$.

Recap :

Assume $B_k \succ 0$ is an approximate Hessian at x_k .

With B_k , $x_{k+1} = x_k - \alpha_k B_k^{-1} \nabla f_k$ can be computed.

With α_k chosen to satisfy Wolfe's condition, $s_k^T y_k \geq 0$

$$x_{k+1} - x_k \approx \nabla f_{k+1} - \nabla f_k$$

Then the following problem :

Find $B_{k+1} = B_k + (\text{rank-2 update})$ st. $B_{k+1} \succ 0$ and

?

$$B_{k+1} s_k = y_k$$

has the following solution by DFP :

$$B_{k+1} = \left(I - \frac{y_k s_k^T}{y_k^T s_k} \right) B_k \left(I - \frac{s_k y_k^T}{y_k^T s_k} \right) + \frac{y_k y_k^T}{y_k^T s_k}$$

What if we aim to approximate/update the inverse Hessian instead of the Hessian? i.e. we ask the following question:

Assume $H_k \succ 0$ is an approximate inverse Hessian at x_k .

Define $x_{k+1} = x_k - \alpha_k H_k \nabla f_k$, α_k satisfies Wolfe's condns, so $s_k^T y_k \geq 0$.

Find $H_{k+1} = H_k + (\text{low rank update})$ st. $H_{k+1} \succ 0$, $H_{k+1} y_k = s_k$.
??

We already have an answer:

$$H_{k+1}^{\text{DFP}} = H_k + \frac{S_k S_k^T}{S_k^T Y_k} - \frac{H_k Y_k Y_k^T H_k}{Y_k^T H_k Y_k}$$

Nothing is new.

But wait... Note that the last question has the exact same algebraic structure as the original question but with y and s interchanged.

What does it mean?

$$B_{k+1}^{\text{DFP}} = \left(I - \frac{Y_k S_k^T}{Y_k^T S_k} \right) B_k \left(I - \frac{S_k Y_k^T}{Y_k^T S_k} \right) + \frac{Y_k Y_k^T}{Y_k^T S_k}$$

$\downarrow \quad y \leftrightarrow s, \quad B \leftrightarrow H$

$$H_{k+1}^{\text{BFGS}} = \left(I - \frac{S_k Y_k^T}{S_k^T Y_k} \right) H_k \left(I - \frac{Y_k S_k^T}{S_k^T Y_k} \right) + \frac{S_k S_k^T}{S_k^T Y_k}$$

NOT only is it new,
but it happens to
outperform DFP!

Yet another surprise, DFP and BFGS methods, similar to Newton's method, have an affine invariant property.

(As mentioned, PSB does not share this property.)

AI guarantees that the ROC is not affected by ill-conditioning.

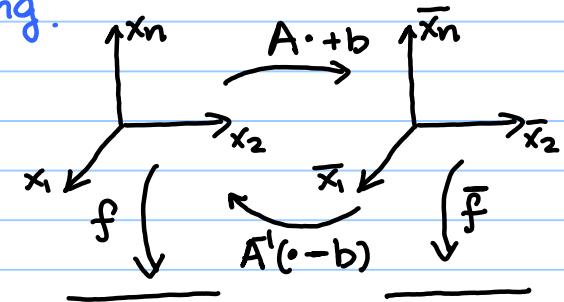
In the original coordinate system:

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

In the transformed coordinate system:

$$\bar{x}_{k+1} = \bar{x}_k - \bar{S}_k \bar{H}_k \nabla \bar{f}(\bar{x}_k)$$

$$\downarrow A^T(-b)$$



Chain rule:

$$\bar{f}(\bar{x}) = f(A^T(\bar{x}-b))$$

$$d\bar{f}(\bar{x}) = df(A^T(\bar{x}-b)) \cdot A^{-1}$$

$$\nabla \bar{f}(\bar{x}) = A^T \nabla f(A^T(\bar{x}-b))$$

$$d\nabla \bar{f}(\bar{x}) = A^T d\nabla f(A^T(\bar{x}-b)) A^{-1}$$

$$\underbrace{\nabla^2 \bar{f}(\bar{x})}_{\nabla^2 f(A^T(\bar{x}-b))} = A^T \underbrace{\nabla^2 f(A^T(\bar{x}-b))^{-1}}_{A^{-T}} A$$

$$[\nabla^2 \bar{f}(\bar{x})]^{-1} = A [\nabla^2 f(A^T(\bar{x}-b))]^{-1} A^T$$

? //

$$\therefore x_{k+1} = x_k - \bar{S}_k A^{-1} \bar{H}_k A^{-T} \nabla \bar{f}(\bar{x}_k)$$

The answer is affirmative if we can show:

(i) $H_k = A^{-1} \bar{H}_k A^{-T}$ (i.e. H_k satisfies the same transformation law as the true inverse Hessian.)

and

$$(ii) \bar{S}_k = S_k \quad \forall k.$$

Ex: (i) Assume for instance we have an affine invariant way to pick the initial approx. inverse Hessian, so that $\bar{H}_0 = A H_0 A^T$. (E.g. by choosing H_0 to be the honest Hessian.)

Prove by induction that $\bar{H}_k = A H_k A^T \quad \forall k \geq 0$,
for both DFP and BFGS updates.

(ii) seems very difficult as s_k is determined by some complicated algorithm for satisfying the Wolfe conditions.

But the line search function, $\bar{\Phi}(\alpha)$, in the transformed coordinates is :

$$\begin{aligned}\bar{\Phi}(\alpha) &= \bar{f}(\bar{x}_k - \alpha \bar{H}_k \nabla \bar{f}(\bar{x}_k)) \\ &= \bar{f}(\bar{x}_k - \alpha A H_k A^T A^{-T} \nabla f(x_k)) \quad \text{by (i) and the transformation law of } \nabla f. \\ &= f(A^{-1}(\bar{x}_k - b - \alpha A H_k \nabla f(x_k))) \\ &= f(x_k - \alpha H_k \nabla f(x_k)) = \Phi(\alpha).\end{aligned}$$

Since the line search function does not change, any algorithm applied to Φ and

$\tilde{\Phi}$ would yield the same step size, i.e. $s_k = \bar{s}_k$.

Technical Remarks :

The derivation of DFP and BFGS above follows the 1977 SIAM Review paper by Dennis-More. There is no mention of affine invariance in the paper.

NW has no mention of Powell's iteration leading to the general formula (P). It only mentions that DFP and BFGS can be obtained by solving the minimization problem $\min_B \|M(B - B^*)M\|_F$ st. $\hat{B} = \hat{B}^T$, $\hat{B}y = y$. — (V)

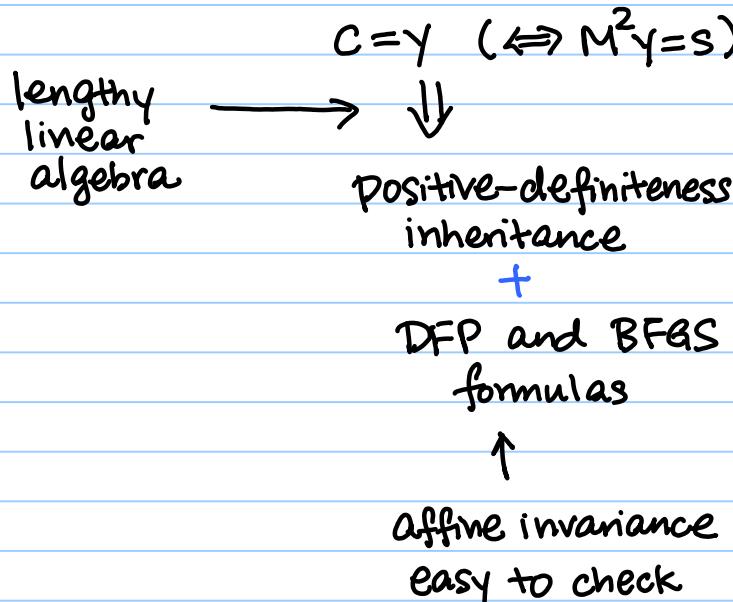
NW repeated 3 times (!) in sec 6.1, but with no clear justification, that a choice of a (symmetric nonsingular) M st. $M^2y = S$ is necessary for the resulted formula to satisfy affine invariance.

And any such M leads to the DFP formula.

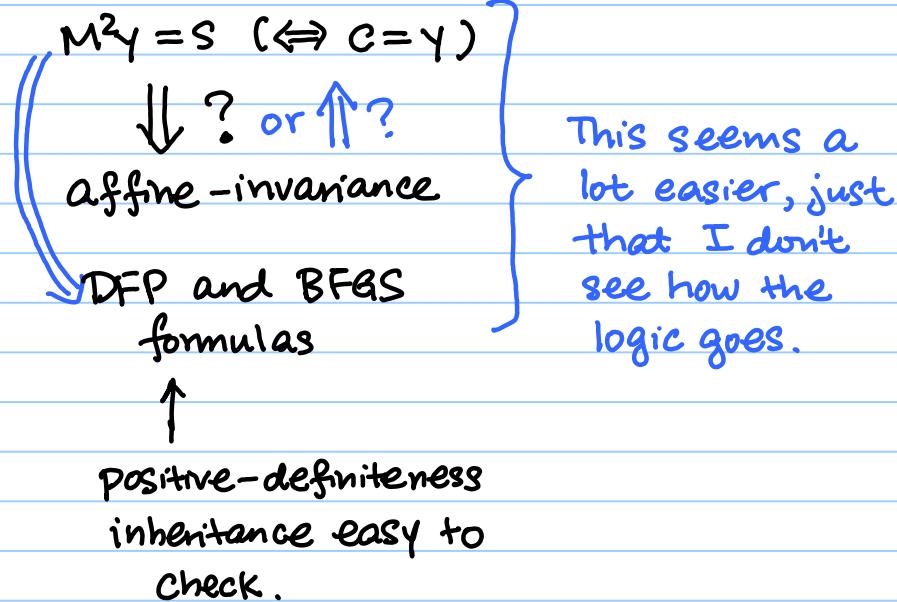
And it is explained in sec 6.1 that the positive-definiteness inheritance property is easy to check (pg 141) once we have the DFP or BFGS formulas.

The derivation of Dennis-Moré is very different. We went through a lengthy linear algebra development to see that the choice of $C = Y$ in (P) — which corresponds to a M with $M^2Y = S$ in the variational formulation (V), according to the theorem on Page 7 above — would imply positive-definiteness inheritance.

Dennis-Moré



N&W



Computational Experiments

Let's try different versions of BFGS, GD, Newton on the generalized Rosenbrock function

$$f(x) = \sum_{i=1}^{n-1} a (x_{i+1} - x_i^2)^2 + (1-x_i)^2$$

It is non-convex degree 4 polynomial, with a global minimizer at $[1, \dots, 1]^T$, and at least a local minimizer near $[-1, 1, \dots, 1]^T$.

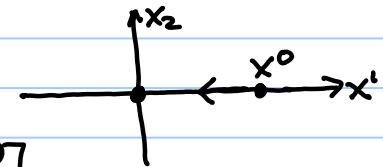
See class demo :

- a naive implementation of BFGS with Wolfe line search and $H_0 = I$.
vs
Matlab's `fminunc()`
- empirical global convergence of BFGS (but there is no theorem for it.)

Note: If $H_0 = \beta I$, BFGS and DFP are invariant under rigid transformations (like GD). So the same examples we use to show that GD can get stuck at a saddle point applies to BFGS or DFP.

$$f(x) = x_1^2 - x_2^2$$

$$H_0 = \beta I, x^0 = \begin{bmatrix} * \\ 0 \end{bmatrix} \Rightarrow \text{the BFGS/DFP iterates } x^k \text{ converge to the saddle pt. } \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$



BFGS and SR-1 methods are known to be remarkable robust in practice, but unfortunately there is no known theoretical result that the iterates of these quasi-Newton methods approach a stationary point from any starting point and any (suitable) initial Hessian approximation.

But a global convergence result can be proved in the convex case.

Assumption 6.1.

- (i) The objective function f is twice continuously differentiable.
- (ii) The level set $\mathcal{L} = \{x \in \mathbb{R}^n \mid f(x) \leq f(x_0)\}$ is convex, and there exist positive constants m and M such that

$$m\|z\|^2 \leq z^T \underbrace{G(x)}_{\nabla^2 f(x)} z \leq M\|z\|^2 \quad (6.39)$$

for all $z \in \mathbb{R}^n$ and $x \in \mathcal{L}$.

Theorem 6.5.

Let B_0 be any symmetric positive definite initial matrix, and let x_0 be a starting point for which Assumption 6.1 is satisfied. Then the sequence $\{x_k\}$ generated by Algorithm 6.1 (with $\epsilon = 0$) converges to the minimizer x^* of f .

The proof is a beautiful application of Zoutendijk's result.
See N&W, sec 6.4.

↑
BFGS with Wolfe line search,
Any H \ddot{o} g \ddot{o} O

The setting of the above theorem can also guarantee that the iterates converge fast enough so that

$$\sum_k \|x_k - x^*\| < \infty.$$

(Proof found in Dennis-Moré)

This isn't quite superlinear convergence, but this condition can somehow be "bootstrapped" to show superlinear convergence :

Assumption 6.2.

The Hessian matrix G is Lipschitz continuous at x^* , that is,

$$\|G(x) - G(x^*)\| \leq L\|x - x^*\|,$$

for all x near x^* , where L is a positive constant.

Theorem 6.6.

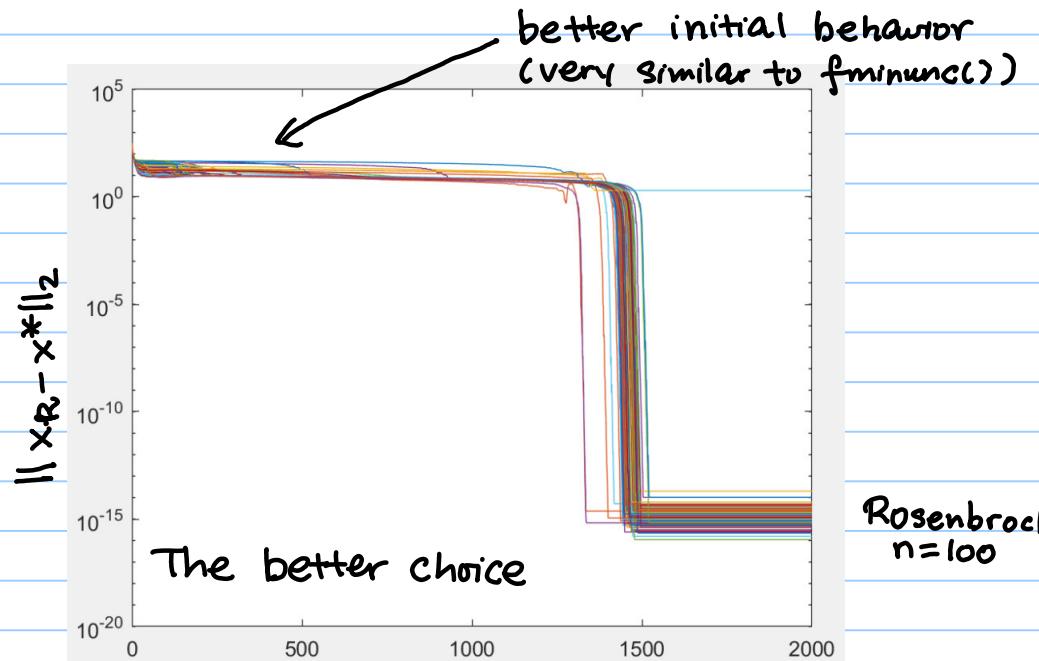
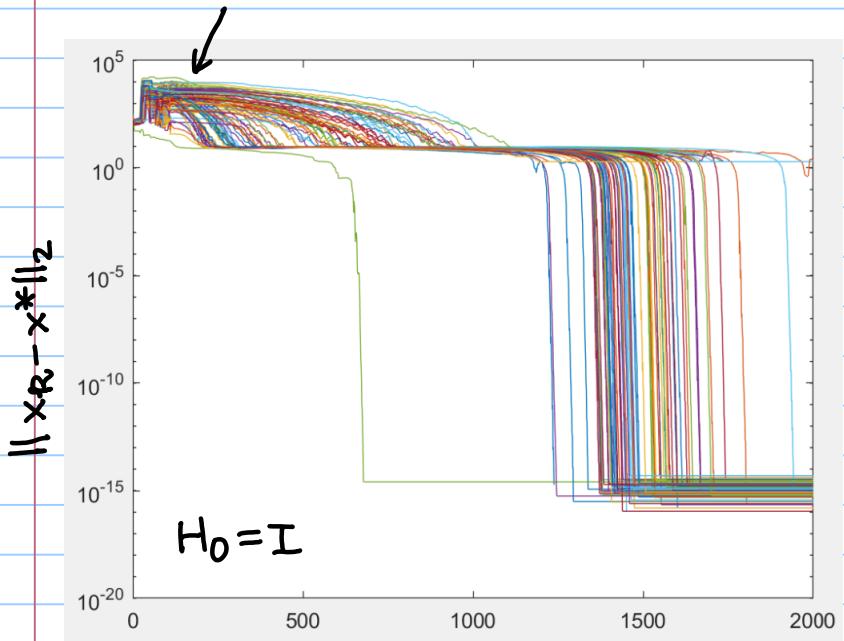
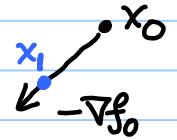
Suppose that f is twice continuously differentiable and that the iterates generated by the BFGS algorithm converge to a minimizer x^* at which Assumption 6.2 holds. Suppose also that (6.52) holds. Then x_k converges to x^* at a superlinear rate.

As you may expect, its proof is based on the fundamental result of Superlinear Convergence by Dennis-Moré presented in UO_4 (last week).

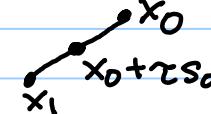
Choice of H_0

If we set $H_0 = I$, then the first step is gradient descent.

- A better choice :
1. Apply one step of GD to get x_1
 2. Compute $s_0 = x_1 - x_0$, $y_0 = \nabla f(x_1) - \nabla f(x_0)$
 3. set $H_0 \leftarrow (y_0^T s_0 / y_0^T y_0) I$.



Here's the rationale for the choice

$$\text{Consider } \int_0^1 \underbrace{\nabla^2 f(x_0 + \tau s_0) s_0}_{= \frac{d}{d\tau} \nabla f(x_0 + \tau s_0)} d\tau = \underbrace{\nabla f(x_1) - \nabla f(x_0)}_{= y_0}$$


$$x_1 - x_0 = s_0$$

Write $\bar{G}_0 := \int_0^1 \nabla^2 f(x_0 + \tau s_0) d\tau$ = average Hessian of f on $\overline{x_0 x_1}$.

Since $\bar{G}_0 s_0 = y_0$, if $\bar{G}_0 > 0$ then $\bar{G}_0^{1/2}$ exists, and

$$\frac{y_0^T s_0}{y_0^T y_0} = \frac{s_0^T \bar{G}_0 s_0}{s_0^T \bar{G}_0 \bar{G}_0 s_0} = \frac{s_0^T \bar{G}_0^{1/2} \bar{G}_0^{1/2} s_0}{s_0^T \bar{G}_0^{1/2} \bar{G}_0 \bar{G}_0^{1/2} s_0} \stackrel{\text{set } z := G_0^{1/2} s_0}{=} \frac{z^T z}{z^T \bar{G}_0 z}.$$

\approx One of the eigenvalues of \bar{G}_0^{-1}

So setting $H_0 = \frac{y_0^T s_0}{y_0^T y_0}$ attempts to make the size of H_0 similar to that of $\nabla^2 f(x_0)^{-1}$.

$$\nabla^2 f(x_0)^{-1}$$

This choice of H_0 is reported to be the most successful (among other choices) in practice. It is also used (in a more dynamic) way in the limited BFGS method, to be presented next.

Limited Memory BFGS

$$\text{BFGS : } x_{R+1} = x_R - \alpha_R H_R \nabla f_R$$

$$H_{R+1} = V_R^T H_R V_R + P_R S_R S_R^T \quad (*)$$

$$P_R = 1/\gamma_R^T S_R, \quad V_R = I - P_R \gamma_R \gamma_R^T, \quad S_R = x_{R+1} - x_R, \quad \gamma_R = \nabla f_{R+1} - \nabla f_R.$$

It requires $O(n^2)$ memory to store H_R . ($O(n^3)$ time if you use (*) verbatim, but only $O(n^2)$ time if you use it wisely.)

Note that if we store $S_0, \gamma_0, S_1, \gamma_1, \dots, S_{R-1}, \gamma_{R-1}$, we can use them to compute $H_R \nabla f_R$ recursively as :

$$H_R q = (I - P_R \gamma_R \gamma_R^T)^T H_{R-1} \underbrace{(I - P_R \gamma_R \gamma_R^T) q}_{q - P_R \gamma_R (S_R^T q)} + P_R S_R S_R^T q = S_R (S_R^T q)$$

↑ ↑ ↑ ↑
 applying it to recur! if $H_0 = \beta I$, $O(n)$ time
 a vector takes if $H_0 = \beta I$, applying it to takes $O(n)$ time
 $O(n)$ time a vector also

But it is not helpful when $k \approx n$, as storing $s_i, y_i, i=0, \dots, k-1$ is as costly as storing H_k .

The idea of L-BFGS is to store only $s_i, y_i, i=k-m, \dots, k-1$, and use them to implicitly represent a modified version of H_k :

$$H_k = (V_{k-1}^T \cdots (V_{k-m+1}^T H_k^0 V_{k-m} + p_{k-m} s_{k-m} s_{k-m}^T) V_{k-m+1} + p_{k-m+1} s_{k-m+1} s_{k-m+1}^T) \cdots V_{k-1} + p_{k-1} s_{k-1} s_{k-1}^T$$

\uparrow
Some choice of initial inverse Hessian approx. at x_k

When $k \leq m$, this is almost the same H_k as before.

When $k > m$, this is a truncated version of the original H_k .

Also, H_k^0 is reassigned as some (crude) inverse Hessian approx. at x_k ,

As in BFGS,

$$\text{set } H_k^0 = \left(\frac{s_{k-1}^T y_{k-1}}{y_{k-1}^T y_{k-1}} \right) I. \quad (7.20)$$

$x_k \quad \cdot x_{k-1}$

(Unlike BFGS, this H_k^0 varies with k .)

We may rewrite H_k as

$$\begin{aligned} H_k &= (V_{k-1}^T \cdots V_{k-m}^T) H_k^0 (V_{k-m} \cdots V_{k-1}) \\ &\quad + \rho_{k-m} (V_{k-1}^T \cdots V_{k-m+1}^T) S_{k-m} S_{k-m}^T (V_{k-m+1} \cdots V_{k-1}) \\ &\quad + \rho_{k-m+1} (V_{k-1}^T \cdots V_{k-m+2}^T) S_{k-m+1} S_{k-m+1}^T (V_{k-m+2} \cdots V_{k-1}) \\ &\quad + \cdots \\ &\quad + \rho_{k-1} S_{k-1} S_{k-1}^T. \end{aligned} \quad (\rho_i = 1/y_i^T s_i, V_i = I - \rho_i y_i s_i^T)$$

From this expression, the following efficient procedure for computing $H_k \nabla f_k$ can be derived.

Algorithm 7.4 (L-BFGS two-loop recursion).

```
q ← ∇fk;
for i = k - 1, k - 2, ..., k - m
    αi ← ρisiTq;
    q ← q - αiyi;
end (for)
r ← Hk0q;
for i = k - m, k - m + 1, ..., k - 1
    β ← ρiyiTr;
    r ← r + si(αi - β)
end (for)
stop with result Hk∇fk = r.
```

It's a bit tricky to see that it does what is supposed to do.

But it is easy to see that it takes O(mn) time, and O(mn) storage.

Algorithm 7.5 (L-BFGS).

Choose starting point x_0 , integer $m > 0$;

$k \leftarrow 0$;

repeat

 Choose H_k^0 (for example, by using (7.20));

 Compute $p_k \leftarrow -H_k \nabla f_k$ from Algorithm 7.4;

 Compute $x_{k+1} \leftarrow x_k + \alpha_k p_k$, where α_k is chosen to
 satisfy the Wolfe conditions;

if $k > m$

 Discard the vector pair $\{s_{k-m}, y_{k-m}\}$ from storage;

 Compute and save $s_k \leftarrow x_{k+1} - x_k$, $y_k = \nabla f_{k+1} - \nabla f_k$;

$k \leftarrow k + 1$;

until convergence.

\curvearrowleft see above