

Variational Problems in Quasi-Newton Methods

Osman Güler, Filiz Gürtuna and Olena Shevchenko*

July 18, 2006

Abstract

It has been known since the early 1970s that the Hessian matrices in quasi-Newton methods can be updated by variational means, in several different ways. The usual formulation of these variational problems uses a coordinate system, and the symmetry of the Hessian matrices are enforced as explicit constraints. As a result, the variational problems seem complicated. In this paper, we give a very simple, coordinate-free, and unified treatment of these variational problems by working directly in the vector space of symmetric matrices. Most of our variational problems become simple, least squares problems. All variational problems are convex programming problems, and we consider their duals. A novel feature of our work is that we construct several new variational problems whose optimal solutions coincide with quasi-Newton update matrices. These new variational problems may be useful in suggesting new quasi-Newton methods in the future.

Key words. Quasi-Newton, DFP, BFGS, variational problems, duality, sparse problems.

Abbreviated title: Variational problems in Quasi-Newton methods

AMS(MOS) subject classifications: primary: 90C53, 90C30, 90C46, 65K10, 49N15; secondary: 90C06, 65K05, 49M37, 52A41.

*Department of Mathematics and Statistics, University of Maryland Baltimore County, Baltimore, Maryland 21250, USA. E-mail: {guler,gurtuna1,olenshe1}@math.umbc.edu. Research partially supported by the National Science Foundation under grant DMS-0411955

1 Introduction

The main point of this paper is to provide clear, geometric formulations and solutions to the several variational problems in quasi-Newton methods. By working directly in the space of symmetric matrices, we handle the symmetry constraint implicitly. This approach simplifies both the formulations and the proofs given in the literature, such as the ones given in [15], [14], [3] and [11].

In quasi-Newton methods for unconstrained minimization, an approximation of the true Hessian is used in the Newton step. Then, at each iteration, the approximation is updated to reflect the information coming from the new point or the gradient at the new point. There are two popular and successful updates, DFP (Davidon [6], Fletcher and Powell [13]) and BFGS (Broyden [2], Fletcher [10], Goldfarb [14], Shanno [19]). Both of these updates perform very well in practice but BFGS is known to have better performance both practically and theoretically. The main condition imposed on the updated matrix is that it satisfies the quasi-Newton equation. Additionally, the matrix is required to be symmetric since the Hessian is symmetric. For a detailed explanation of these, we refer the reader to [18] and the survey paper of Dennis and Moré [7].

The basic idea in obtaining the update formula is that in addition to the requirement that the new matrix satisfies the quasi-Newton equation, the matrix should also be "close" to the current matrix [15], [14]. The reason for this "closeness" requirement is so as not to lose any information that the current approximation may have about the true Hessian. Different measures lead to different update formulas. The originally proposed measure is the weighted Frobenius norm. It is known that BFGS and DFP updates have this property, i.e. being the minimum norm change to the current matrix [7]. Later, Byrd and Nocedal [3] introduced another measure involving trace and determinant in their convergence analysis of BFGS. Then, Fletcher [11] showed that BFGS and DFP updates can be obtained from the minimization of this measure as well.

In what follows, B represents approximations of the Hessian and H represents approximations of the inverse Hessian. Superscripts DFP and BFGS will be used to refer to the corresponding updates. For easy reference, these formulas are presented below

$$B_{k+1}^{\text{DFP}} := (I - \gamma_k y_k s_k^T) B_k (I - \gamma_k s_k y_k^T) + \gamma_k y_k y_k^T, \quad (1.1)$$

$$B_{k+1}^{\text{BFGS}} := B_k - \frac{B_k s_k s_k^T B_k}{s_k^T B_k s_k} + \gamma_k y_k y_k^T, \quad (1.2)$$

$$H_{k+1}^{\text{DFP}} := H_k - \frac{H_k y_k y_k^T H_k}{y_k^T H_k y_k} + \gamma_k s_k s_k^T, \quad (1.3)$$

$$H_{k+1}^{\text{BFGS}} := (I - \gamma_k s_k y_k^T) H_k (I - \gamma_k y_k s_k^T) + \gamma_k s_k s_k^T, \quad (1.4)$$

where $\gamma_k = 1/y_k^T s_k$.

It is worth noting the strong relations between B^{DFP} and H^{BFGS} and between H^{DFP} and B^{BFGS} . In fact, as will be mentioned later, these two pairs come from minimizing two different measures.

The paper is organized as follows. In § 2, we formulate the various norm minimization problems in quasi-Newton methods as geometric least squares problems in \mathbb{S}^n , the vector space $n \times n$ of symmetric matrices. In § 3, we formulate the duals of these least squares problems in two different ways and show that the optimal solutions to the dual problems also produce the quasi-Newton update formulas. In § 4, we treat in the same manner the primal and dual variational problems involving the minimization of the measure function $\langle H_k, B \rangle - \ln \det B$ of Byrd and Nocedal [3] and Fletcher [11]. In § 5, we treat some variational problems arising in quasi-Newton methods for sparse problems. We also discuss some algorithmic issues. In the Appendix (§ 6), we gather several results used in the main body of the paper.

2 Least squares problems in quasi-Newton methods

It is well-known that the approximations to the Hessian matrices in various quasi-Newton methods are updated using the solutions of some optimization problems. In the optimization problems used to obtain the DFP and BFGS updates, for example, the constraints have the form

$$\{X \in \mathbb{R}^{n \times n} : Xs = y, X^T = X\}, \quad (2.1)$$

where s and y are given vectors in \mathbb{R}^n . The first affine equation $Xs = y$ is called the *secant equation* (or *quasi-Newton equation*), and the constraint $X^T = X$ is included since a Hessian matrix (or its inverse) is always symmetric, and so should be any approximations to it.

In the traditional approach to solving these optimization problems, one first writes down the Lagrangian function by associating a multiplier $\lambda \in \mathbb{R}^n$ for the secant equation and by treating the symmetry equation as $n(n-1)/2$ scalar equations $X_{ij} = X_{ji}$, $1 \leq i < j \leq n$. One then proceeds to write the Karush-Kuhn-Tucker conditions and to solve the resulting equations. Although the required calculations are not difficult, some parts are unnecessary, and the resulting formula for the optimal solution seems somewhat obscure.

In this paper, we solve such problems in a simpler, geometric fashion, which also makes the solutions look very natural. We first note that the symmetry constraints can be *eliminated* entirely simply by working in the vector space \mathbb{S}^n of $n \times n$ symmetric matrices instead of the vector space $\mathbb{R}^{n \times n}$. We endow the vector space $\mathbb{R}^{n \times n}$ with the trace inner product

$$\langle X, Y \rangle := \text{tr}(X^T Y) = \sum_{i,j=1}^n X_{ij} Y_{ij},$$

which induces the trace inner product

$$\langle X, Y \rangle = \text{tr}(XY)$$

in \mathbb{S}^n . Both vector spaces become Euclidean spaces with these inner products.

We first need a preliminary result which is interesting in its own right.

Lemma 2.1. *Let $s, y \in \mathbb{R}^n$, $s \neq 0$. Consider the affine subspace $\mathcal{A} := \{X \in \mathbb{S}^n : Xs = y\}$ in the vector space \mathbb{S}^n . The linear subspace corresponding to \mathcal{A} is $\mathcal{L} = \{X \in \mathbb{S}^n : Xs = 0\}$. Let $\{u_i\}_1^n$ be a basis of \mathbb{R}^n , and define the matrices $S_i = su_i^T + u_i s^T$, $i = 1, \dots, n$. The matrices $\{S_i\}_1^n$ are linearly independent and \mathcal{L} is the intersection of n hyperplanes in \mathbb{S}^n ,*

$$\mathcal{L} = \{X \in \mathbb{S}^n : \langle X, S_i \rangle = 0, \quad i = 1, \dots, n\}.$$

Moreover,

$$\mathcal{L}^\perp = \text{span}\{S_1, \dots, S_n\} = \{s\lambda^T + \lambda s^T : \lambda \in \mathbb{R}^n\}. \quad (2.2)$$

Proof. The formula for \mathcal{L} is obvious. Notice that the equation $Xs = 0$ in $\mathbb{R}^{n \times n}$ is equivalent to the component equations

$$\langle X, su_i^T \rangle = 0, \quad i = 1, \dots, n,$$

since

$$0 = \langle u_i, Xs \rangle = \text{tr}(u_i^T Xs) = \text{tr}(Xsu_i^T) = \langle X, su_i^T \rangle, \quad i = 1, \dots, n.$$

Since X is symmetric, we also have

$$\langle X, u_i s^T \rangle = \text{tr}(Xu_i s^T) = \text{tr}(su_i^T X) = \text{tr}(u_i^T Xs) = \text{tr}(Xsu_i^T) = \langle X, su_i^T \rangle,$$

so that the equation $Xs = 0$ is also equivalent to the component equations

$$\langle X, u_i s^T \rangle = 0, \quad i = 1, \dots, n.$$

Consequently, $\mathcal{L} \subseteq \mathbb{S}^n$ can be written as an intersection of n hyperplanes $\langle X, S_i \rangle = 0$, $i = 1, \dots, n$.

The formula $\mathcal{L}^\perp = \text{span}\{S_1, \dots, S_n\}$ follows immediately. Any linear combination $\sum_{i=1}^n \delta_i S_i \in \mathcal{L}^\perp$ can be written as

$$\sum_{i=1}^n \delta_i (su_i^T + u_i s^T) = s\lambda^T + \lambda s^T$$

where $\lambda = \sum_{i=1}^n \delta_i u_i$.

The matrices $\{S_i\}_1^n$ are linearly independent: the equation $0 = \sum_{i=1}^n \delta_i S_i = \lambda s^T + s\lambda^T$ gives $0 = (s^T \lambda)\lambda + \|\lambda\|^2 s$, and taking the inner product of both sides with s yields $(s^T \lambda)^2 + \|\lambda\|^2 \cdot \|s\|^2 = 0$. Thus, $\|\lambda\|^2 \cdot \|s\|^2 = 0$, and since $s \neq 0$, we have $\lambda = 0$. Since $\lambda = \sum_{i=1}^n \delta_i u_i = 0$ and $\{u_i\}_1^n$ is a basis of \mathbb{R}^n , we have $\delta_i = 0$, $i = 1, \dots, n$. \square

Most of the variational problems encountered in quasi-Newton methods are closely related to the following generic least squares problem.

Theorem 2.2. *The solution \bar{X} to the minimization problem in \mathbb{S}^n*

$$\begin{aligned} \min \quad & \frac{1}{2} \|X\|^2 \\ \text{s. t.} \quad & Xs = y \end{aligned} \tag{2.3}$$

is given by

$$\bar{X} = \frac{sy^T + ys^T}{\langle s, s \rangle} - \frac{\langle y, s \rangle}{\langle s, s \rangle^2} ss^T. \tag{2.4}$$

Proof. Define $f(X) = \|X\|^2 = \langle X, X \rangle / 2$. We have $\nabla f(X) = X$ and $\nabla^2 f(X) = I$, and the function f is convex. It follows from Lemma 6.1 that the solution \bar{X} is characterized by the condition $\nabla f(\bar{X}) = \bar{X} \in \mathcal{L}^\perp$. Consequently, Lemma 2.1 implies that \bar{X} is characterized by the equation

$$\bar{X} = \lambda s^T + s\lambda^T,$$

for some $\lambda \in \mathbb{R}^n$. We have

$$\langle y, s \rangle = \langle \bar{X}s, s \rangle = \langle [\lambda s^T + s\lambda^T]s, s \rangle = 2\langle \lambda, s \rangle \|s\|^2,$$

and $\langle \lambda, s \rangle = \langle y, s \rangle / (2\|s\|^2)$. Substituting this in the equation $y = \bar{X}s = \langle \lambda, s \rangle s + \|s\|^2 \lambda$ gives

$$\lambda = \frac{1}{\|s\|^2} y - \frac{\langle y, s \rangle}{2\|s\|^4} s.$$

Finally, substituting this in $\bar{X} = \lambda s^T + s\lambda^T$ gives equation (2.4). \square

Corollary 2.3. *Let $X_0 \in \mathbb{S}^n$ be given and W be a symmetric positive definite weight matrix. The optimal solution \bar{X} to the minimization problem (in \mathbb{S}^n)*

$$\begin{aligned} \min \quad & \frac{1}{2} \|W^{1/2}(X - X_0)W^{1/2}\|^2 \\ \text{s. t.} \quad & Xs = y, \end{aligned} \tag{2.5}$$

is given by

$$\begin{aligned} \bar{X} = X_0 + & \frac{W^{-1}s(y - X_0s)^T + (y - X_0s)s^T W^{-1}}{\langle s, W^{-1}s \rangle} \\ & - \frac{\langle y - X_0s, s \rangle}{\langle s, W^{-1}s \rangle^2} W^{-1}ss^T W^{-1}. \end{aligned}$$

Proof. With the following change of variables

$$\hat{X} := W^{1/2}(X - X_0)W^{1/2}, \quad \hat{y} := W^{1/2}(y - X_0s), \quad \hat{s} := W^{-1/2}s,$$

the problem (2.5) reduces to problem (2.3). After substituting the expressions for \hat{X} , \hat{y} , and \hat{s} into (2.4), we multiply the resulting equality by $W^{-1/2}$ from both sides to get the desired expression for \bar{X} . \square

The formulas for the DFP and BFGS updates also follow as easy corollaries of Theorem 2.2.

Corollary 2.4. *The DFP update matrix B_{k+1} in (1.1) is the optimal solution to the minimization problem in \mathbb{S}^n*

$$\begin{aligned} \min \quad & \frac{1}{2} \|W^{1/2}(B - B_k)W^{1/2}\|^2 \\ \text{s. t.} \quad & Bs_k = y_k, \end{aligned} \tag{2.6}$$

where $B_k \in \mathbb{S}^n$ is an approximation to the Hessian of f at iteration k , $y_k := \nabla f(x_{k+1}) - \nabla f(x_k)$, $s_k := x_{k+1} - x_k$, and W is any symmetric positive definite matrix satisfying $Wy_k = s_k$.

Proof. Using Corollary 2.3, we have

$$\begin{aligned} B_{k+1} = B_k + & \frac{W^{-1}s_k(y_k - B_k s_k)^T + (y_k - B_k s_k)s_k^T W^{-1}}{\langle s_k, W^{-1}s_k \rangle} \\ & - \frac{\langle y_k - B_k s_k, s_k \rangle}{\langle s_k, W^{-1}s_k \rangle^2} W^{-1}s_k s_k^T W^{-1}. \end{aligned}$$

The requirement $Wy_k = s_k$ (or $y_k = W^{-1}s_k$) simplifies the above expression and makes B_{k+1} independent of W . It is a routine to verify that the resulting formula for B_{k+1} is the same as the one obtained by expanding (1.1). \square

It is evident from (1.1) that $B_{k+1} = G^T G + F$ where $F = \gamma_k y_k y_k^T$ and $G = B_k^{1/2}(I - \gamma_k s_k y_k^T)$. Since both $G^T G$ and F are positive semi definite, so is B_{k+1} . Moreover, $Gd = 0$ for $d \neq 0$ if and only if d is a multiple of s_k , but then $\langle Fd, d \rangle > 0$. Thus, we see that if B_k is positive definite and $\langle s_k, y_k \rangle > 0$, then the matrix B_{k+1} is also positive definite.

Corollary 2.5. *The BFGS update matrix H_{k+1} in (1.4) is the optimal solution to the minimization problem in \mathbb{S}^n*

$$\begin{aligned} \min \quad & \frac{1}{2} \|W^{1/2}(H - H_k)W^{1/2}\|^2 \\ \text{s. t.} \quad & Hy_k = s_k, \end{aligned} \tag{2.7}$$

where s_k and y_k are defined as in Corollary 2.4, $H_k \in \mathbb{S}^n$ is a matrix approximating the inverse Hessian of f at iteration k , and W is any symmetric positive definite matrix satisfying $Ws_k = y_k$.

Proof. The proof is similar to the proof of Corollary 2.4. As in the DFP case above, if H_k is positive definite and $\langle y_k, s_k \rangle > 0$, then H_{k+1} is positive definite. \square

3 Dual least squares problems

Although the minimization problems 2.6 and 2.7 are well-known in the literature since the 1970s, it seems that the associated dual problems have not been studied so far, as far as we can determine. In this section, we give *two* dual problems for the least squares minimization problem. The first dual problem turns out to be another minimization problem whose solution *coincides* with the DFP and BFGS update formulas and the second dual problem gives a different geometric interpretation for the update formulas.

We first consider a primal–dual pair of geometric least squares problems, see Courant–Hilbert [5], pp. 252–257.

Theorem 3.1. *Let $x_0, y_0 \in E$ be given points in a Euclidean space E , and let $L \subseteq E$ be a linear subspace of E . The following least squares problems are duals of each other. Furthermore, they have the same optimal solution.*

$$(P) \quad \min_{x \in y_0 + L} \frac{1}{2} \|x - x_0\|^2 \qquad (D) \quad \min_{y \in x_0 + L^\perp} \frac{1}{2} \|y - y_0\|^2$$

Proof. Note that the problem (P) can be written as a minimax problem

$$\min_{x \in E} \max_{\lambda \in L^\perp} L(x, \lambda) := \frac{1}{2} \|x - x_0\|^2 + \langle y_0 - x, \lambda \rangle,$$

since $\max_{\lambda \in L^\perp} \langle x - y_0, \lambda \rangle = 0$ if $x \in y_0 + L$, and $+\infty$ otherwise. The dual problem with respect to the Lagrangian function $L(x, \lambda)$ is the maximin problem

$$\max_{\lambda \in L^\perp} \min_{x \in E} L(x, \lambda) := \frac{1}{2} \|x - x_0\|^2 + \langle y_0 - x, \lambda \rangle.$$

The inner minimum is achieved at the point $x^* = x_0 + \lambda$. Substituting this in L and rearranging the terms, we get $L(x^*, \lambda) = -\|\lambda + x_0 - y_0\|^2/2 + \|x_0 - y_0\|^2/2$. Thus the dual problem becomes, up to an additive constant $\|x_0 - y_0\|^2/2$,

$$\max_{\lambda \in L^\perp} -\frac{1}{2} \|\lambda + x_0 - y_0\|^2 = -\min_{\lambda \in L^\perp} \frac{1}{2} \|\lambda + x_0 - y_0\|^2.$$

With the change of variables $y = \lambda + x_0$, the right hand side problem is equivalent to (D).

Let x^*, y^* be the optimal solutions of (P) and (D), respectively. We have

$$x^* - x_0 \in L^\perp, \quad x^* - y_0 \in L, \quad y^* - y_0 \in L, \quad y^* - x_0 \in L^\perp,$$

where the first and third inclusions follow from Lemma 6.1. These imply $x^* - y^* = (x^* - x_0) - (y^* - x_0) \in L^\perp$ and $x^* - y^* = (x^* - y_0) - (y^* - y_0) \in L$. Thus, $x^* - y^* \in L \cap L^\perp = \{0\}$, that is, $x^* = y^*$. \square

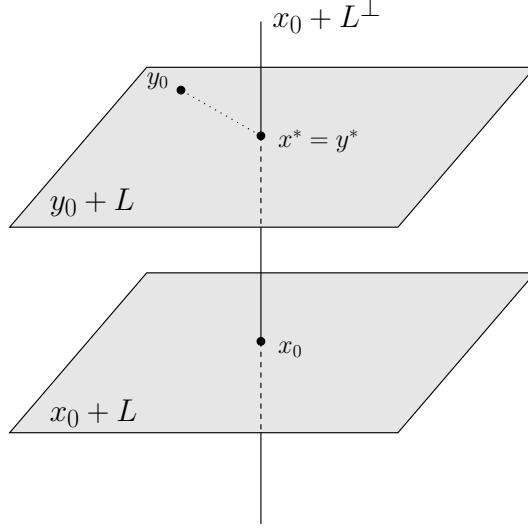


Figure 1: Illustration of Theorem 3.1

Remark 3.2. We emphasize the fact that the above pair of least squares problems (P) and (D) have the remarkable property that they have the same optimal solution. This is illustrated in Figure 1. Note that the primal problem (P) is the (orthogonal) projection of the point x_0 on the lower affine subspace $x_0 + L$ onto the upper affine subspace $y_0 + L$, whereas the dual problem (D) is the projection of the point y_0 on the upper affine subspace $x_0 + L$ onto the complementary affine subspace $x_0 + L^\perp$. As the proof above shows, the equality $x^* = y^*$ of the optimal solutions together with the Strong Duality Theorem (which holds true in this case since no Slater conditions are needed for affine constraints) implies the equality

$$\frac{1}{2}\|x^* - x_0\|^2 = -\frac{1}{2}\|x^* - y_0\|^2 + \frac{1}{2}\|x_0 - y_0\|^2,$$

where the left hand side is the value of the minimax problem and the right hand side is the value of the maximin problem. This amounts to the equation

$$\|x_0 - y_0\|^2 = \|x_0 - x^*\|^2 + \|x^* - y_0\|^2, \quad x_0 - x^* \in L^\perp, x^* - y_0 \in L,$$

which is precisely the Pythagorean theorem applied to the triangle with vertices $\{x_0, y_0, x^*\}$.

We now use Theorem 3.1 to obtain the duals of the least squares problems (2.6) and (1.4). By Theorem 3.1, these dual problems are new variational problems whose solutions are the DFP and BFGS updates, respectively.

Note that in problem (2.6), if we make the change of variables

$$\hat{B} := W^{1/2}(B - B_k)W^{1/2}, \quad \hat{y}_k := W^{1/2}y_k, \quad \hat{s}_k := W^{-1/2}s_k,$$

we arrive at the least squares problem in \mathbb{S}^n

$$\begin{aligned} \min \quad & \frac{1}{2} \|\hat{B} - \hat{B}_k\|^2 \\ \text{s. t.} \quad & \hat{B} \hat{s}_k = \hat{y}_k. \end{aligned}$$

Theorem 3.1 and Lemma 2.1 can then be used to obtain a dual least squares problem which in turn can be transformed into another least squares problem in terms of the original variables. Here we obtain this final dual problem in a different fashion, by changing the inner product instead of the variables.

Let $W \in \mathbb{S}^n$ be a positive definite matrix. Consider W -norm on \mathbb{S}^n given by

$$\|X\|_W^2 := \text{tr}(W^{1/2} X W^{1/2})^2 = \text{tr}(W X W X),$$

and the corresponding inner-product

$$\langle X, Y \rangle_W := \text{tr}(W X W Y) = \text{tr}((W \otimes W) X Y) = \langle (W \otimes W) X, Y \rangle,$$

where $(W \otimes W) X := W X W$. In the Euclidean space $(\mathbb{S}^n, \|\cdot\|_W)$, the problem (2.6) becomes

$$\begin{aligned} \min \quad & \frac{1}{2} \|B - B_k\|_W^2 \\ \text{s. t.} \quad & B s_k = y_k, \end{aligned} \tag{3.1}$$

to which Theorem 3.1 applies. Let \bar{B} be any matrix in the affine constraint set $\mathcal{A} := \{B \in \mathbb{S}^n : B s_k = y_k\}$. Then $\mathcal{A} = \bar{B} + \mathcal{L}$ where $\mathcal{L} = \{B : B s_k = 0\}$. In order to determine the dual problem in this setting, we need to compute the orthogonal complement of \mathcal{L} . This is done in the lemma below, which is an analogue of Lemma 2.1.

Lemma 3.3. *Let $W \in \mathbb{S}^n$ be a positive definite matrix and $0 \neq s \in \mathbb{R}^n$. The orthogonal complement of the linear subspace $\mathcal{L} = \{B : B s = 0\}$ in the Euclidean space $(\mathbb{S}^n, \|\cdot\|_W)$ is*

$$\mathcal{L}^\perp = \{\lambda(W^{-1}s)^T + (W^{-1}s)\lambda^T : \lambda \in \mathbb{R}^n\}.$$

Proof. Let $\{u_i\}$ be a basis of \mathbb{R}^n . \mathcal{L} is characterized by the component equations $u_i^T B s = s^T B u_i = 0$, $i = 1, \dots, n$, or equally well by the equations

$$\begin{aligned} 0 &= \langle B, u_i s^T + s u_i^T \rangle = \langle B, W^{-1}(u_i s^T + s u_i^T) W^{-1} \rangle_W \\ &= \langle B, (W^{-1} u_i)(W^{-1} s)^T + (W^{-1} s)(W^{-1} u_i)^T \rangle_W, \quad i = 1, \dots, n. \end{aligned}$$

It follows that

$$\begin{aligned} \mathcal{L}^\perp &= \text{span}\{(W^{-1} u_i)(W^{-1} s)^T + (W^{-1} s)(W^{-1} u_i)^T : i = 1, \dots, n\} \\ &= \{\lambda(W^{-1} s)^T + (W^{-1} s)\lambda^T : \lambda \in \mathbb{R}^n\}. \end{aligned}$$

□

Corollary 3.4. *The matrix B_{k+1} in the DFP update formula (1.1) is the optimal solution to the least squares problem*

$$\min \left\{ \frac{1}{2} \|B_k - \bar{B} + \lambda y_k^T + y_k \lambda^T\|_W^2 : \lambda \in \mathbb{R}^n \right\}, \quad (3.2)$$

where y_k , s_k , B_k , and W satisfy the conditions in Corollary 2.4, and $\bar{B} \in \mathbb{S}^n$ is any matrix satisfying the secant equation $\bar{B}s_k = y_k$. In particular, we may choose

$$\bar{B} = \frac{y_k y_k^T}{\langle s_k, y_k \rangle}.$$

Proof. The affine constraint set in (3.1) is $\bar{B} + \mathcal{L}$ where $\mathcal{L} = \{B : Bs_k = 0\}$, and we have $W^{-1}s_k = y_k$. The proof follows immediately from Theorem 3.1 and Lemma 3.3. \square

Similarly, we also have

Corollary 3.5. *The matrix H_{k+1} in the BFGS update formula (1.4) is the optimal solution to the least squares problem*

$$\min \left\{ \frac{1}{2} \|H_k - \bar{H} + \lambda s_k^T + s_k \lambda^T\|_W^2 : \lambda \in \mathbb{R}^n \right\}, \quad (3.3)$$

where y_k , s_k , B_k , and W satisfy the conditions in Corollary 2.5, and $\bar{H} \in \mathbb{S}^n$ is any matrix satisfying the secant equation $\bar{H}y_k = s_k$. In particular, we may choose

$$\bar{H} = \frac{s_k s_k^T}{\langle s_k, y_k \rangle}.$$

3.1 Another dualization of the least squares problems

We now give a pair of different dual problems for the minimization problems (2.6) and (2.7). They provide a different interpretation for the DFP and BFGS updates.

Theorem 3.6. *Let y_k , s_k , B_k , and W satisfy the conditions in Corollary 2.4. Let $\bar{B} \in \mathbb{S}^n$ be any matrix satisfying $Bs_k = y_k$ (say $\bar{B} = (y_k y_k^T) / \langle s_k, y_k \rangle$). The problem*

$$\begin{aligned} \min \quad & \langle B_k - \bar{B}, Y \rangle_W \\ \text{s. t.} \quad & \|Y\|_W \leq 1, \\ & Y = \lambda y_k^T + y_k \lambda^T, \lambda \in \mathbb{R}^n, \end{aligned} \quad (3.4)$$

is dual to the minimization problem (2.6) for the DFP update. The DFP update formula is given by

$$B_{k+1} = B_k + \alpha \bar{Y}$$

where \bar{Y} is the solution to (3.4) and α is chosen so that the secant equation $B_{k+1}s_k = y_k$ is satisfied.

Proof. We write problem (2.6) in the form $\min\{\|B - B_k\|_W : B \in \bar{B} + \mathcal{L}\}$ where $\mathcal{L} = \{B : Bs_k = 0\}$. We have

$$\begin{aligned} \min_{B \in \bar{B} + \mathcal{L}} \|B - B_k\|_W &= \min_{B \in \bar{B} + \mathcal{L}} \max_{\|Y\|_W \leq 1} \langle B - B_k, Y \rangle_W \\ &= \max_{\|Y\|_W \leq 1} \min_{B \in \bar{B} + \mathcal{L}} \langle B - B_k, Y \rangle_W \\ &= \max_{\|Y\|_W \leq 1} \left\{ \langle \bar{B} - B_k, Y \rangle_W + \min_{X \in \mathcal{L}} \langle X, Y \rangle_W \right\} \\ &= \max_{\|Y\|_W \leq 1, Y \in \mathcal{L}^\perp} \langle \bar{B} - B_k, Y \rangle_W, \end{aligned}$$

where the second equality follows from the minimax theorem, and the last equality follows from the fact that $\min\{\langle X, Y \rangle_W : X \in \mathcal{L}\}$ equals zero if $Y \in \mathcal{L}^\perp$, and $-\infty$ otherwise. Using the Cauchy–Schwarz inequality, the first equality above holds only if $B_{k+1} - B_k = \alpha \bar{Y}$ for some $\alpha \in \mathbb{R}$, where \bar{Y} is the optimal solution to (3.4). This completes the proof, since \mathcal{L}^\perp is given in Lemma 3.3 and $W^{-1}s_k = y_k$. \square

We note that the dual optimal solution \bar{Y} is the point in \mathcal{L}^\perp making the *smallest* angle (in the W -inner product) with the point $\bar{B} - B_k$.

Similarly, we also have

Theorem 3.7. *Let y_k, s_k, B_k , and W satisfy the conditions in Corollary 2.5. Let $\bar{H} \in \mathbb{S}^n$ be any matrix satisfying $H y_k = s_k$ (say $\bar{H} = (s_k s_k^T) / \langle s_k, y_k \rangle$). The problem*

$$\begin{aligned} \min \quad & \langle H_k - \bar{H}, Z \rangle_W \\ \text{s. t.} \quad & \|Z\|_W \leq 1, \\ & Z = \lambda s_k^T + s_k \lambda^T, \lambda \in \mathbb{R}^n, \end{aligned} \tag{3.5}$$

is dual to the minimization problem (2.7) for the BFGS update. The BFGS update formula is given by

$$H_{k+1} = H_k + \alpha \bar{Z}$$

where \bar{Z} is the solution to (3.5) and α is chosen so that the secant equation $H_{k+1}y_k = s_k$ is satisfied.

4 Trace–determinant function minimization problems in quasi–Newton methods

The function $\psi : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$ defined by

$$\psi(X) = \text{tr } X - \ln \det X = \langle I, X \rangle - \ln \det X \tag{4.1}$$

is used by Byrd and Nocedal [3] in the convergence analysis of the BFGS update rule. Later, Fletcher [11] shows that the DFP and BFGS update rules can be obtained

from the minimization of function $\psi(X)$ subject to the same constraints as in the least squares minimization case. Here, we present a more geometric version of the proof of main Theorem 2.1 in [11].

Theorem 4.1. *Let the affine set $\{X \in \mathbb{S}^n : Xs = y\}$ contain a positive definite matrix. The solution \bar{X} to the minimization problem in \mathbb{S}^n*

$$\begin{aligned} \min \quad & \psi(X) = \langle I, X \rangle - \ln \det X \\ \text{s. t.} \quad & Xs = y \end{aligned} \tag{4.2}$$

satisfies

$$\bar{X}^{-1} = I + \frac{ss^T}{\langle y, s \rangle} - \frac{sy^T + ys^T}{\langle y, s \rangle} + \frac{\langle y, y \rangle}{\langle y, s \rangle^2} ss^T. \tag{4.3}$$

Proof. The first and second derivatives of ψ are given by

$$\nabla \psi(x) = I - X^{-1}, \quad \nabla^2 \psi(x) = X^{-1} \otimes X^{-1},$$

see equation (6.4) in the Appendix. Thus, ψ is strictly convex on the cone of positive definite matrices in \mathbb{S}^n (the semidefinite cone). It is also coercive on the same cone. Lemma 6.1 and Lemma 2.1 imply that the optimal solution \bar{X} satisfies the condition

$$I - \bar{X}^{-1} = s\lambda^T + \lambda s^T, \tag{4.4}$$

for some $\lambda \in \mathbb{R}^n$. Using the secant equation $\bar{X}^{-1}y = s$, we get

$$\langle y - s, y \rangle = \langle (I - \bar{X}^{-1})y, y \rangle = \langle [s\lambda^T + \lambda s^T]y, y \rangle = 2\langle \lambda, y \rangle \langle y, s \rangle,$$

which yields $\langle \lambda, y \rangle = \langle y - s, y \rangle / (2\langle y, s \rangle)$. Substituting this in the equation

$$y - s = (I - \bar{X}^{-1})y = \langle \lambda, y \rangle s + \langle y, s \rangle \lambda$$

gives

$$\lambda = \frac{y - s}{\langle s, y \rangle} + \frac{\langle s - y, y \rangle}{2\langle s, y \rangle^2} s.$$

Substituting this in (4.4) and simplifying the result gives (4.3). \square

Corollary 4.2. *Let $\langle s_k, y_k \rangle > 0$ and $H_k \in \mathbb{S}^n$ be positive definite, where H_k is the approximation to the inverse Hessian at iteration k . The matrix H_{k+1} in the BFGS update formula (1.4) satisfies $H_{k+1} = \bar{B}^{-1}$ where \bar{B} is the optimal solution to the minimization problem in \mathbb{S}^n*

$$\begin{aligned} \min \quad & \psi(H_k^{1/2} B H_k^{1/2}) = \langle H_k, B \rangle - \ln \det B + \text{const} \\ \text{s. t.} \quad & Bs_k = y_k. \end{aligned} \tag{4.5}$$

Proof. The change of variables $X = H_k^{1/2} B H_k^{1/2}$, $y = H_k^{1/2} y_k$, and $s = H_k^{-1/2} s_k$ reduces the problem to the minimization problem (4.2) in Theorem 4.1. Substituting the values of X, y, s above in equation (4.3) and simplifying, we obtain

$$\bar{B}^{-1} = H_k - \frac{s_k y_k^T H_k + H_k y_k s_k^T}{\langle s_k, y_k \rangle} + \frac{s_k s_k^T}{\langle s_k, y_k \rangle} + \frac{\langle H_k y_k, y_k \rangle}{\langle s_k, y_k \rangle^2} s_k s_k^T.$$

The right hand side of this formula is identical to the one in (1.4). Consequently, $\bar{B} = H_{k+1}^{-1}$. This completes the proof. \square

Corollary 4.3. *Let $\langle s_k, y_k \rangle > 0$ and $B_k \in \mathbb{S}^n$ be positive definite, where B_k is the approximation to the Hessian at iteration k . The matrix B_{k+1} in the DFP update formula (1.1) satisfies $B_{k+1} = \bar{H}^{-1}$ where \bar{H} is the optimal solution to the minimization problem in \mathbb{S}^n*

$$\begin{aligned} \min \quad & \psi(B_k^{1/2} H B_k^{1/2}) = \langle B_k, H \rangle - \ln \det H + \text{const} \\ \text{s. t.} \quad & H y_k = s_k. \end{aligned}$$

Proof. This is similar to the proof of Corollary 4.2, using the change of variables $X = B_k^{1/2} H B_k^{1/2}$, $y = B_k^{1/2} s_k$, and $s = B_k^{-1/2} y_k$. \square

4.1 Dual of the trace–determinant function minimization problem

Here we work out the dual to the minimization problem (4.5) for the BFGS update. As in the case of the least squares minimization problems, this dual does not seem to be studied in the literature. As we will see, the optimal solution to the dual problem is directly related to the BFGS update matrix. Similar results also hold true for the DFP update formulas.

Theorem 4.4. *Let $X_0, Y_0 \in \mathbb{S}^n$. The following minimization problems are duals of each other,*

$$\begin{aligned} (P) \quad \min \quad & \langle X_0, X \rangle - \ln \det X \\ & X \in Y_0 + \mathcal{L} \\ (D) \quad \min \quad & \langle Y_0, Y \rangle - \ln \det Y \\ & Y \in X_0 + \mathcal{L}^\perp \end{aligned}$$

If (P) and (D) both have feasible positive definite solutions, then they both have optimal solutions and the strong duality theorem holds true. Furthermore, the optimal solutions of (P) and (D) are inverses of each other, that is, $\bar{Y} = (\bar{X})^{-1}$ where \bar{X} and \bar{Y} are the optimal solutions of (P) and (D), respectively.

Proof. Since the objective functions in (P) and (D) are coercive, both problems have optimal solutions and the strong duality theorem holds true. The problem (P) can be written as a minimax problem

$$\min_X \max_{Z \in \mathcal{L}^\perp} L(X, Z) := \langle X_0, X \rangle - \ln \det X + \langle X - Y_0, Z \rangle,$$

since $\max_{Z \in \mathcal{L}^\perp} \langle X - Y_0, Z \rangle = 0$ if $X - Y_0 \in \mathcal{L}$, and $+\infty$ otherwise. The dual problem with respect to $L(X, Z)$ is

$$\max_{Z \in \mathcal{L}^\perp} \min_X \{ \langle X_0, X \rangle - \ln \det X + \langle X - Y_0, Z \rangle \}.$$

The inner minimum is achieved at the point \hat{X} satisfying

$$(\hat{X})^{-1} = X_0 + Z. \quad (4.6)$$

Substituting this in L and simplifying, we arrive at $L(\hat{X}, Z) = n - \langle Y_0, Z \rangle + \ln \det(X_0 + Z)$. Thus, the dual problem is equivalent to

$$\min_{Z \in \mathcal{L}^\perp} \{ \langle Y_0, Z \rangle - \ln \det(X_0 + Z) \}.$$

With the change of variables $Y = X_0 + Z$, and using the description of \mathcal{L}^\perp in (2.2), we see that the above problem is equivalent to (D). It follows from (4.6) that $\bar{Y} = (\bar{X})^{-1}$. \square

The following corollary is immediate, using the choices $X_0 = H_k$, $Y_0 \in \mathbb{S}^n$ is any matrix satisfying the condition $Bs_k = y_k$ (such as $Y_0 = y_k y_k^T / \langle s_k, y_k \rangle$), and $\mathcal{L} = \{B \in \mathbb{S}^n : Bs_k = 0\}$.

Corollary 4.5. *Let s_k, y_k, H_k and H_{k+1} be defined as in Corollary 4.2. H_{k+1} in the BFGS update formula (1.4) is the optimal \bar{Y} to the minimization problem*

$$\begin{aligned} \min \quad & \langle Y_0, Y \rangle - \ln \det Y \\ & Y = H_k + \lambda s_k^T + s_k \lambda^T, \quad \lambda \in \mathbb{R}^n. \end{aligned} \quad (4.7)$$

Thus, we obtain here a new result that the BFGS update matrices B_{k+1} and $H_{k+1} := B_{k+1}^{-1}$ come from the primal dual problems (4.5) and (4.7), respectively.

A similar result holds for the DFP updates.

5 Quasi-Newton methods for sparse problems

In an unconstrained optimization problem, if the Hessian of the function f to be minimized has some sparsity pattern, then it is desirable that the approximations to the Hessian have the same pattern. For large-scale problems this becomes very important because at each iteration of the quasi-Newton method a linear system involving the approximation to the Hessian is solved.

Assume that the Hessian $H = H(x) := \nabla^2 f(x)$ has the same sparse structure at each point x . Define $S = \{(i, j) : H_{ij} = 0\}$, then the sparsity condition on the approximations becomes $B_{ij} = 0$ for all $(i, j) \in S$. These conditions can be included

in the minimization problems for updating rules, but solving the corresponding problems efficiently becomes an issue.

The form of the solutions to the least squares problems (2.6) and (2.7) allows us to conclude that they are positive definite under the condition $\langle s_k, y_k \rangle > 0$. However, adding sparsity constraints to the variational problems changes the structure of the optimal solution and positive definiteness can no longer be satisfied automatically. One way to overcome this difficulty is to approximate the obtained update with a positive definite matrix. Alternatively, the DFP and BFGS updates can be obtained from minimization of the trace–determinant function $\psi(X)$. The positive definiteness of the optimal solution is guaranteed by the nature of $\psi(X)$.

Devising efficient quasi–Newton methods for sparse unconstrained optimization problems is a very important issue. Limited memory quasi–Newton methods of Nocedal [17] are popular, and the particular structure of the sparsity pattern can be exploited, but the research for better algorithms is ongoing. Recently, there has been exciting developments in the use of l_1 –norm minimization methods and convex optimization for obtaining sparse solutions for some problems, see [4],[9], [21]. It is not clear at present what impact these developments will have on quasi–Newton methods.

In this section, we will be content with giving short and transparent solutions to two variational results in this area as well as highlighting some related algorithmic issues.

Theorem 5.1. *(Toint [20]) Let $S \subset \{(i, j) : 1 \leq i, j \leq n\}$. Consider the minimization problem in the vector space \mathbb{S}^n*

$$\begin{aligned} \min \quad & \frac{1}{2} \|X\|^2 \\ \text{s. t.} \quad & Xs = y, \\ & X_{ij} = 0, \quad (i, j) \in S. \end{aligned} \tag{5.1}$$

Define $\mathcal{L} := \{X : X_{ij} = 0, (i, j) \in S\}$. The optimal solution to (5.1) is given by

$$\bar{X} = \Pi_{\mathcal{L}}(\lambda s^T + s \lambda^T), \tag{5.2}$$

where λ is the solution of the linear equations in \mathbb{R}^n ,

$$Q\lambda = y,$$

where $Q^T = [s_1, s_2, \dots, s_n]$, $s_i := S_i s$, $S_i := \Pi_{\mathcal{L}}(s e_i^T + e_i s^T)$, $i = 1, \dots, n$.

Proof. Define $\mathcal{M} := \{X : Xs = 0\}$. Lemma 6.1 implies that $\bar{X} \in (\mathcal{M} \cap \mathcal{L})^\perp = \mathcal{M}^\perp + \mathcal{L}^\perp$. Write

$$\bar{X} = (\lambda s^T + s \lambda^T) + \Lambda$$

with $\lambda s^T + s \lambda^T \in \mathcal{M}^\perp$ (see Lemma 2.1), and $\Lambda \in \mathcal{L}^\perp$. Since $\bar{X} \in \mathcal{L}$, we see that

$$\bar{X} = \Pi_{\mathcal{L}}(\lambda s^T + s\lambda^T).$$

We have $y = \bar{X}s = \Pi_{\mathcal{L}}(\lambda s^T + s\lambda^T)s$, and

$$\begin{aligned} y_i &= \langle [\Pi_{\mathcal{L}}(\lambda s^T + s\lambda^T)] s, e_i \rangle = \left\langle \Pi_{\mathcal{L}}(\lambda s^T + s\lambda^T), \frac{e_i s^T + s e_i^T}{2} \right\rangle \\ &= \left\langle \frac{\lambda s^T + s\lambda^T}{2}, \Pi_{\mathcal{L}}(e_i s^T + s e_i^T) \right\rangle = \left\langle \frac{\lambda s^T + s\lambda^T}{2}, S_i \right\rangle \\ &= \langle S_i s, \lambda \rangle = \langle s_i, \lambda \rangle. \end{aligned}$$

□

The projection operator $\Pi_{\mathcal{L}}$ is the simple, so-called ‘‘gangster operator’’ defined by

$$\mathcal{G}(H)_{ij} = \begin{cases} 0 & (i, j) \in S, \\ H_{ij} & (i, j) \in S^\perp \end{cases}$$

with S^\perp denoting the complement of S , because it shoots ‘‘holes’’ at the entries $(i, j) \in S$ of matrix H .

We remark that the matrix Q is symmetric, since

$$\begin{aligned} Q_{ij} &= (s_i)_j = \langle s_i, e_j \rangle = \langle \Pi_{\mathcal{L}}(s e_i^T + e_i s^T) s, e_j \rangle = \left\langle \Pi_{\mathcal{L}}(s e_i^T + e_i s^T), \frac{s e_j + e_j s^T}{2} \right\rangle \\ &= \left\langle \frac{s e_i^T + e_i s^T}{2}, \Pi_{\mathcal{L}}(s e_j + e_j s^T) \right\rangle = Q_{ji}, \end{aligned}$$

and has the same sparsity pattern S : We have

$$Q_{ij} = \left\langle \Pi_{\mathcal{L}}(s e_i^T + e_i s^T), \frac{s e_j + e_j s^T}{2} \right\rangle = \langle \Pi_{\mathcal{L}}(s e_i^T + e_i s^T) e_j, s \rangle,$$

and it is easy to show that $\Pi_{\mathcal{L}}(s e_i^T + e_i s^T) e_j = 0$ if $(i, j) \in S$.

We also give a short proof of a variational problem of Fletcher.

Theorem 5.2. (Fletcher [12]) *Let B_k be positive definite, and $H_k = B_k^{-1}$. The solution \bar{B} to the minimization problem in \mathbb{S}^n*

$$\begin{aligned} \min \quad & \psi_{H_k}(B) := \langle H_k, B \rangle - \ln \det(B) \\ \text{s. t.} \quad & B s_k = y_k, \\ & B_{ij} = 0, \quad (i, j) \in S \end{aligned} \tag{5.3}$$

is characterized by the existence of λ such that

$$\mathcal{G}(\bar{H}) = \mathcal{G}(H_k + \lambda s^T + s\lambda^T),$$

where $\bar{H} = \bar{B}^{-1}$.

Proof. As in the proof of Theorem 5.1, define $\mathcal{L} := \{B : B_{ij} = 0, (i, j) \in S\}$ and $\mathcal{M} := \{B : Bs_k = 0\}$. The optimal solution \bar{B} satisfies

$$\nabla_B \psi_{H_k}(\bar{B}) = H_k - \bar{B}^{-1} \in \mathcal{M}^\perp + \mathcal{L}^\perp,$$

that is,

$$\bar{B}^{-1} = H_k + \lambda s_k^T + s_k \lambda^T + \Lambda,$$

where $\Lambda \in \mathcal{L}^\perp$. The theorem is proved since $\mathcal{G}(\Lambda) = 0$. \square

We see again that the sparsity constraints complicate the problem and it becomes hard to obtain an explicit formula for the update. Fletcher [12] proposes a numerical method to solve the above nonlinear system. Although, he provides numerical results, no convergence analysis is given. One way of obtaining complexity results for this problem is to use Newton's method. The function ψ_{H_k} is a self-concordant function and the complexity of Newton's method can be estimated for such functions, see Nesterov and Nemirovskii [16]. At each step of Newton's method for minimizing (5.3), we need to solve the least squares problem in \mathbb{S}^n , where X is known and positive definite,

$$\begin{aligned} \min \quad & \frac{1}{2} \|X^{-1/2} [(\Delta X) - (X - XH_kX)] X^{-1/2}\|^2 \\ \text{s. t.} \quad & (\Delta X)_{s_k} = 0, \\ & (\Delta X)_{ij} = 0, \quad (i, j) \in S, \end{aligned} \tag{5.4}$$

see (6.3) and (6.4) in the Appendix. Since the problem (5.4) is the weighted projection of $X - XH_kX$ onto the intersection of two linear subspaces, obtaining an explicit solution is not easy. Next, we discuss a way of attacking this problem using the method of alternating projections.

5.1 The method of alternating projections

We have already encountered two least squares problems (5.1) and (5.4) which involve projecting a point onto the intersection of two linear subspaces. In these problems, projecting onto each subspace is trivial, but projecting onto the intersection is not. In the first problem, the problem was reduced to solving a small system of equations, but in general this is not possible.

The iterative *method of alternating projections* finds the projection of a point to the intersection of a finite number of closed linear subspaces by projecting onto each subspace in succession. The method and its convergence, in the case of two subspaces, seems to be due to von Neumann, see Lemma 22 in [22], and also [1], pp. 375–380.

Theorem 5.3. *Let $\mathcal{A}_1, \dots, \mathcal{A}_r$ be closed linear spaces in a Hilbert space H , Π_1, \dots, Π_r the orthogonal projection operators onto $\mathcal{A}_1, \dots, \mathcal{A}_r$, respectively, and Π the orthogonal projection operator onto $\cap_1^r \mathcal{A}_i$. Then for each $x \in H$,*

$$\lim_{n \rightarrow \infty} (\pi_r \pi_{r-1} \dots \pi_1)^n(x) = \pi(x), \quad (5.5)$$

and

$$\|(\pi_r \pi_{r-1} \dots \pi_1)^n(x) - \pi(x)\| \leq c^n \|x\|. \quad (5.6)$$

Here,

$$c := \left[1 - \prod_{i=1}^{r-1} (1 - c_i^2) \right]^{1/2}$$

and c_i is the cosine of the angle between M_i and $\bigcap_{j=i+1}^r M_j$.

Theorem 5.4. For $r = 2$, the estimate (5.6) in Theorem 5.3 can be replaced by

$$\|(\pi_2 \pi_1)^n(x) - \pi(x)\| \leq c^{2n-1} \|x\|.$$

Furthermore, $\|(\pi_2 \pi_1)^n - \pi\| = c^{2n-1}$.

The proof of the theorems can be found in Chapter 9 of Deutsch [8]. Thus, the problem of projecting a point onto the intersection of linear spaces can be solved if the projection onto each space is easy. Theorem 5.4 gives a sharp convergence rate. In the general case of an arbitrary number of linear subspaces, a sharp theoretical rate of convergence is given in [23].

6 Appendix

In this appendix, we collect for completeness several results used in the main body of the paper.

6.1 Minimizing a function on an affine subspace

Lemma 6.1. Let $A = a + L \subseteq E$ be an affine set in a Euclidean space E where $a \in A$ and L is a linear subspace of E . Let $f : A \rightarrow \mathbb{R}$ be a differentiable function, and consider the minimization problem $\min\{f(x) : x \in A\}$. If a vector $\bar{x} \in A$ is a local minimizer of f , then

$$\nabla f(\bar{x}) \in L^\perp. \quad (6.1)$$

If f is convex, then (6.1) is a sufficient condition for \bar{x} to be a global minimizer of f on A .

Proof. Let $x \in A$ be an arbitrary point of A . For $|t|$ small enough, we have

$$f(\bar{x}) \leq f(\bar{x} + t(x - \bar{x})) = f(\bar{x}) + t\langle \nabla f(\bar{x}), x - \bar{x} \rangle + o(t),$$

where the inequality follows from \bar{x} 's being a local minimizer of f , and the equality follows from Taylor's formula. Thus, we have

$$t\langle \nabla f(\bar{x}), x - \bar{x} \rangle + o(t) \geq 0,$$

for all t small enough. For $t > 0$, dividing both sides by t and letting $t \rightarrow 0$ gives $\langle \nabla f(\bar{x}), x - \bar{x} \rangle \geq 0$. For $t < 0$, the same procedure leads to $\langle \nabla f(\bar{x}), x - \bar{x} \rangle \leq 0$. Since an arbitrary point in L can be represented as $x - \bar{x}$ for some $x \in A$, we obtain equation (6.1).

If f is convex and x is any point in A , we have $f(x) \geq f(\bar{x}) + \langle \nabla f(\bar{x}), x - \bar{x} \rangle = f(\bar{x})$, where the inequality follows from the convexity of f and the equality follows from (6.1). \square

Corollary 6.2. *The orthogonal projection of an $n \times n$ matrix A onto the subspace of symmetric matrices is $(A + A^T)/2$, that is, $\Pi_{\mathbb{S}^n}(A) = (A + A^T)/2$.*

Proof. The orthogonal projection $\Pi_{\mathbb{S}^n}(A)$ is the optimal solution to the problem

$$\min \left\{ f(X) := \frac{1}{2} \|X - A\|^2 : X \in \mathbb{S}^n \right\}$$

on the vector space $\mathbb{R}^{n \times n}$. Since $\nabla f(X) = X - A$ and $\nabla^2 f(X) = I$, f is convex. By Lemma 6.1 in the Appendix, it suffices to verify that $\nabla f(X^*) = (A + A^T)/2 - A = (A^T - A)/2$ is orthogonal to every symmetric matrix. If X is symmetric, we have

$$\langle A^T - A, X \rangle = \text{tr}((AX - A^T X)) = \text{tr}(AX) - \text{tr}((XA)^T) = 0,$$

since $\text{tr}(A) = \text{tr}(A^T)$ and $\text{tr}(AB) = \text{tr}(BA)$. \square

6.2 Newton's method for minimizing a function on an affine subspace

Consider minimizing a strongly convex function ($\nabla^2 f(x)$ positive definite) on an affine subspace \mathcal{A} ,

$$\min \{ f(x) : x \in \mathcal{A} \}.$$

Let $\mathcal{A} = a + L$ where L is a linear subspace. Write $L = \mathcal{R}(A)$ where $A : V \rightarrow \mathbb{S}^n$ is a linear operator on a vector space V . The problem can be written as an unconstrained minimization problem on V

$$\min g(y) := f(a + Ay).$$

The Newton direction Δy minimizes the quadratic approximation of $g(y)$ given by

$$g(y) + \langle \nabla g(y), \Delta y \rangle + \frac{1}{2} \langle \nabla^2 g(y) \Delta y, \Delta y \rangle,$$

thus, Δy solves the linear system

$$\nabla^2 g(y) \Delta y + \nabla g(y) = 0. \tag{6.2}$$

Using the Taylor's expansion,

$$\begin{aligned} g(y + \Delta y) &= f(a + A(y + \Delta y)) \\ &= f(a + Ay) + \langle \nabla f(a + Ay), A\Delta y \rangle + \frac{1}{2} \langle \nabla^2 f(a + Ay) A\Delta y, A\Delta y \rangle + \dots, \end{aligned}$$

we see that (6.2) becomes

$$A^* \nabla^2 f(a + Ay) A\Delta y + A^* \nabla f(a + Ay) = 0.$$

In terms of the original variables, we obtain

$$\nabla^2 f(x) \Delta x + \nabla f(x) \in \mathcal{N}(A^*) = \mathcal{R}(A)^\perp = L^\perp.$$

Using (6.1), this is equivalent to stating that Δx solves the problem

$$\min \{ f(x) + \langle \nabla f(x), \Delta x \rangle + \frac{1}{2} \langle \nabla^2 f(x) \Delta x, \Delta x \rangle : \Delta x \in L \},$$

which is in fact the least squares minimization problem

$$\min \left\{ \frac{1}{2} \| (\nabla^2 f(x))^{1/2} [\Delta x + (\nabla^2 f(x))^{-1} \nabla f(x)] \|^2 : \Delta x \in L \right\}. \quad (6.3)$$

Consequently, Newton's direction for minimizing a function on an affine set is precisely the projection of the unconstrained Newton's direction $-(\nabla^2 f(x))^{-1} \nabla f(x)$ onto the linear subspace L with respect to the inner product with weight $(\nabla^2 f(x))^{1/2}$.

6.3 The derivatives of the log-determinant function

We expand the Taylor series of the function $f(X) = \ln \det X$ at $X \in \mathbb{S}^n$ positive definite and $D \in \mathbb{S}^n$ as follows.

$$\begin{aligned} \Delta f &:= f(X + tD) - f(X) = \ln \det(X + tD) - \ln \det X \\ &= \ln \det(X^{1/2}(I + tX^{-1/2}DX^{-1/2}X^{1/2})) - \ln \det X \\ &= \ln \det(I + tX^{-1/2}DX^{-1/2}). \end{aligned}$$

Define $\widehat{D} := X^{-1/2}DX^{-1/2}$. Writing the orthogonal decomposition of \widehat{D} in the form $\widehat{D} = Q\Lambda Q^T$ where Q is an $n \times n$ orthogonal matrix (whose columns are the eigenvectors of \widehat{D}), and Λ is an $n \times n$ diagonal matrix whose elements are the eigenvalues of \widehat{D} , we have

$$\begin{aligned} \Delta f &= \ln \det(I + t\Lambda) = \ln \prod_{i=1}^n (1 + t\lambda_i) = \sum_{i=1}^n \ln(1 + t\lambda_i) \\ &= t \operatorname{tr} \Lambda - \frac{t^2}{2} \operatorname{tr} \Lambda^2 + o(t^3) = t \operatorname{tr}(\widehat{D}) - \frac{t^2}{2} \operatorname{tr}(\widehat{D})^2 + o(t^3) \\ &= t \langle X^{-1}, D \rangle - \frac{t^2}{2} \langle (X^{-1}DX^{-1}), D \rangle + o(t^3), \end{aligned}$$

where the operator \otimes is defined as $(X^{-1} \otimes X^{-1})D := X^{-1}DX^{-1}$, and where we used $\ln(1 + t\alpha) = t\alpha - \frac{1}{2}(t\alpha)^2 + o(t^3)$ in the fourth equation. We conclude that

$$\nabla f(X) = X^{-1}, \quad \nabla^2 f(X) = -X^{-1} \otimes X^{-1}. \quad (6.4)$$

References

- [1] N. Aronszajn. Theory of reproducing kernels. *Trans. Amer. Math. Soc.*, 68:337–404, 1950.
- [2] C. G. Broyden. The convergence of an algorithm for solving sparse nonlinear systems. *Math. Comp.*, 25:285–294, 1971.
- [3] R. H. Byrd and J. Nocedal. A tool for the analysis of quasi-Newton methods with application to unconstrained minimization. *SIAM J. Numer. Anal.*, 26(3):727–739, 1989.
- [4] E. J. Candès, J. K. Romberg, and T. Tao. Stable signal recovery from incomplete and inaccurate measurements. *Comm. Pure Appl. Math.*, 59(8):1207–1223, 2006.
- [5] R. Courant and D. Hilbert. *Methods of mathematical physics. Vol. I*. Interscience Publishers, Inc., New York, N.Y., 1953.
- [6] W. C. Davidon. Variable metric method for minimization. *SIAM J. Optim.*, 1(1):1–17, 1991.
- [7] J. E. Dennis, Jr. and J. J. Moré. Quasi-Newton methods, motivation and theory. *SIAM Rev.*, 19(1):46–89, 1977.
- [8] F. Deutsch. *Best approximation in inner product spaces*. CMS Books in Mathematics/Ouvrages de Mathématiques de la SMC, 7. Springer-Verlag, New York, 2001.
- [9] D. L. Donoho. For most large underdetermined systems of linear equations the minimal l_1 -norm solution is also the sparsest solution. *Comm. Pure Appl. Math.*, 59(6):797–829, 2006.
- [10] R. Fletcher. A new approach to variable metric methods. *Comput. J.*, 13:317–322, 1970.
- [11] R. Fletcher. A new variational result for quasi-Newton formulae. *SIAM J. Optim.*, 1(1):18–21, 1991.
- [12] R. Fletcher. An optimal positive definite update for sparse Hessian matrices. *SIAM J. Optim.*, 5(1):192–218, 1995.

- [13] R. Fletcher and M. J. D. Powell. A rapidly convergent descent method for minimization. *Comput. J.*, 6:163–168, 1963/1964.
- [14] D. Goldfarb. A family of variable-metric methods derived by variational means. *Math. Comp.*, 24:23–26, 1970.
- [15] J. Greenstadt. Variations on variable-metric methods. (With discussion). *Math. Comp.*, 24:1–22, 1970.
- [16] Y. E. Nesterov and A. S. Nemirovskii. *Interior-point polynomial algorithms in convex programming*, volume 13 of *SIAM Studies in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1994.
- [17] J. Nocedal. Updating quasi-Newton matrices with limited storage. *Math. Comp.*, 35(151):773–782, 1980.
- [18] J. Nocedal. Theory of algorithms for unconstrained optimization. In *Acta numerica, 1992*, *Acta Numer.*, pages 199–242. Cambridge Univ. Press, Cambridge, 1992.
- [19] D. F. Shanno. Conditioning of quasi-Newton methods for function minimization. *Math. Comp.*, 24:647–656, 1970.
- [20] Ph. L. Toint. On sparse and symmetric matrix updating subject to a linear equation. *Math. Comp.*, 31(no 140):954–961, 1977.
- [21] J. A. Tropp. Greed is good: algorithmic results for sparse approximation. *IEEE Trans. Inform. Theory*, 50(10):2231–2242, 2004.
- [22] J. von Neumann. On rings of operators. Reduction theory. *Ann. of Math. (2)*, 50:401–485, 1949.
- [23] J. Xu and L. Zikatanov. The method of alternating projections and the method of subspace corrections in Hilbert space. *J. Amer. Math. Soc.*, 15(3):573–597 (electronic), 2002.