



Limited-Memory BFGS Diagonal Preconditioners for a Data Assimilation Problem in Meteorology

F. VEERSÉ

INRIA Rhône-Alpes, Monbonnot Saint Martin, France

email: fabrice.veerse@imag.fr

D. AUROUX*

Ecole Normale Supérieure, Lyon, France

email: didier.auroux@ens-lyon.fr

M. FISHER

European Centre for Medium-Range Weather Forecasts, Reading, UK

Received December 17, 1999; Revised September 27, 2000

Abstract. This paper uses two simple variational data assimilation problems with the 1D viscous Burgers' equation on a periodic domain to investigate the impact of various diagonal-preconditioner update and scaling strategies, both on the limited-memory BFGS (Broyden, Fletcher, Goldfarb and Shanno) inverse Hessian approximation and on the minimization performance. These simple problems share some characteristics with the large-scale variational data assimilation problems commonly dealt with in meteorology and oceanography.

The update formulae studied are those proposed by Gilbert and Lemaréchal (*Math. Prog.*, vol. 45, pp. 407–435, 1989) and the quasi-Cauchy formula of Zhu et al. (*SIAM J. Optim.*, vol. 9, pp. 1192–1204, 1999). Which information should be used for updating the diagonal preconditioner, the one to be forgotten or the most recent one, is considered first. Then, following the former authors, a scaling of the diagonal preconditioner is introduced for the corresponding formulae in order to improve the minimization performance. The large negative impact of such a scaling on the quality of the L-BFGS inverse Hessian approximation led us to propose an alternate updating and scaling strategy, that provides a good inverse Hessian approximation and gives the best minimization performance for the problems considered. With this approach the quality of the inverse Hessian approximation improves steadily during the minimization process. Moreover, this quality and the L-BFGS minimization performance improves when the amount of stored information is increased.

Keywords: quasi-Newton, limited-memory BFGS, diagonal updating, quasi-Cauchy, large-scale minimization

1. Introduction

Numerical simulations in meteorology and oceanography consist essentially in integrating a discretized version of the system of partial derivative equations (PDEs) modeling the evolution of the atmosphere and/or the ocean. This system of PDEs comprises a dynamical core derived from Navier-Stokes equations using relevant approximations, an equation of state for the fluid of interest, an equation representing the first law of thermodynamics, and physical parameterizations modeling subgrid-scale processes (convection, radiation,

*Corresponding author.

precipitation, turbulence, surface drag, etc.). To perform this numerical integration, one needs to provide an initial state and possibly boundary conditions. It is the purpose of the so-called “data assimilation” discipline to merge the information coming from the observations with that from the physical laws governing the fluid evolution (available under the form of a numerical model), in order to infer the initial and boundary conditions that will lead to the best-quality simulation or prediction. The so-called “variational data assimilation (VDA) method” tries to achieve this using optimal-control techniques (fitting the model trajectory to the observations), weighting both sources of information by their respective error covariances. Because the dimension of the model state vector is usually large (10^5 – 10^9) and the relation between the model variables are complex, it is not possible in practice to handle the covariance matrix of forecast errors used to weight the information coming from the model (10^{10} – 10^{18} scalar components). Instead, the corresponding linear operator is modeled as the composition of operators that can be managed with current parallel computers. As a result the specified error covariance information is almost always climatologic and does not depend on the underlying dynamics. This is a major deficiency of current implementations of VDA. Recently Veersé proposed a method to specify dynamical forecast error covariances using limited-memory quasi-Newton operators (Veersé, 1999, to appear). For such a method to be efficient in practice, it is required that the limited-memory inverse Hessian approximation be of good quality. This motivated a study for assessing this quality using a simple model sharing some characteristics with atmospheric models. Not surprisingly the way the limited-memory BFGS (L-BFGS) diagonal preconditioner is specified became a touchstone in this study. This led us to assess some of the diagonal preconditioners proposed in the literature (Nocedal, 1980; Gilbert and Lemaréchal, 1989; Liu and Nocedal, 1989; Zhu et al., 1999) and to propose some alternatives.

L-BFGS implements a limited-memory version of the Broyden-Fletcher-Goldfarb-Shanno update formula for the inverse Hessian (Broyden, 1969; Dennis and Moré, 1977; Gilbert and Lemaréchal, 1989):

$$\mathbf{H}^+ = U(\mathbf{H}, \mathbf{y}, \mathbf{s}) = \left(\mathbf{I} - \frac{\mathbf{s} \otimes \mathbf{y}}{\langle \mathbf{y}, \mathbf{s} \rangle} \right) \mathbf{H} \left(\mathbf{I} - \frac{\mathbf{y} \otimes \mathbf{s}}{\langle \mathbf{y}, \mathbf{s} \rangle} \right) + \frac{\mathbf{s} \otimes \mathbf{s}}{\langle \mathbf{y}, \mathbf{s} \rangle} \quad (1)$$

where \mathbf{H}^+ is the updated inverse Hessian, $\mathbf{s} = \mathbf{x}^+ - \mathbf{x}$ is the difference between the new iterate and the previous one, and $\mathbf{y} = \mathbf{g}^+ - \mathbf{g}$ is the corresponding gradient increment. Here $\langle \cdot, \cdot \rangle$ is the scalar product with respect to which the gradient is defined and the minimization is to be performed; $\mathbf{u} \otimes \mathbf{v}$ is the linear operator that to a vector \mathbf{d} associate the vector $\langle \mathbf{v}, \mathbf{d} \rangle \mathbf{u}$. In the limited version (Nocedal, 1980; Gilbert and Lemaréchal, 1989; Liu and Nocedal, 1989) aimed at large-scale unconstrained minimization problems, one can afford to store say m couples of vectors (\mathbf{s}, \mathbf{y}) . The above update formula is used for the first m iterations. For the subsequent ones, the following algorithm is used:

$$\begin{aligned} \mathbf{H}_k^0 &= \mathbf{D}_k, \\ \mathbf{H}_k^{i+1} &= U(\mathbf{H}_k^i, \mathbf{y}_{k-m+i}, \mathbf{s}_{k-m+i}), \quad \text{for } 0 \leq i \leq m-1, \\ \mathbf{H}_k &= \mathbf{H}_k^m \end{aligned} \quad (2)$$

The starting matrix \mathbf{D}_k is diagonal and several formulations have been proposed and tested in Gilbert and Lemaréchal (1989) and by Liu and Nocedal (1989). Recently Zhu et al.

(1999) proposed a variant, based on the quasi-Cauchy relation $\langle \mathbf{D}\mathbf{y}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{s} \rangle$ that the experiments of Gilbert and Lemaréchal (1989) showed it was beneficial to enforce (by rescaling the diagonal matrix \mathbf{D} before updating it). The M1QN3 minimization code from the INRIA MODULOPT library either specifies \mathbf{D}_k as the identity matrix multiplied by the Oren-Spedicato factor $\langle \mathbf{y}_{k-1}, \mathbf{s}_{k-1} \rangle / \langle \mathbf{y}_{k-1}, \mathbf{y}_{k-1} \rangle$, or updates \mathbf{D}_k using a scaled version of a “diagonalized” BFGS formula (Eq. (4.9) in Gilbert and Lemaréchal (1989)). In our study the focus is put on the latter case, as it is usually more efficient.

During the minimization only the multiplication of the approximate inverse Hessian matrix by a given vector is needed and is performed efficiently using a two-loop recursion proposed by Nocedal (1980), and the corresponding matrices are never formed. It is precisely this aspect that makes the limited-memory inverse BFGS algorithm suitable for large-scale VDA in meteorology and oceanography, as the size of the corresponding matrices may typically reach $10^5 \times 10^5$.

The paper is organized as follows. In the next section, the simple VDA problems used in this study are presented. Section 3 details the updating formulae for the diagonal preconditioner \mathbf{D}_k , that are assessed. The following two sections present the results from numerical experiments for a quadratic and a non-quadratic VDA problem respectively. A brief discussion based on the present case study is then proposed.

2. The variational data assimilation case study

4D variational data assimilation (4D-Var) may be expressed as the minimization of a cost function $\mathcal{J}(\mathbf{x})$ which measures the misfit of a model state vector \mathbf{x} to a set of observations \mathbf{y}_o and an a priori (background) estimate \mathbf{x}_b of the true vector \mathbf{x}_t . The cost function may be written as

$$\mathcal{J}_{\mathbf{x}}(\mathbf{x}) = (\mathbf{x} - \mathbf{x}_b)^T \mathbf{B}^{-1} (\mathbf{x} - \mathbf{x}_b) + (\mathcal{H}(\mathbf{x}) - \mathbf{y}_o)^T \mathbf{O}^{-1} (\mathcal{H}(\mathbf{x}) - \mathbf{y}_o). \quad (3)$$

The matrices \mathbf{B} and \mathbf{O} are the covariance matrices for random errors in $(\mathbf{x}_t - \mathbf{x}_b)$ and $(\mathcal{H}(\mathbf{x}_t) - \mathbf{y}_o)$ respectively. \mathcal{H} is an operator which, applied to \mathbf{x} produces an estimate of the observation vector. In 4D-Var, the observations are spread over a period of time from t_0 to $t_0 + T$, while the control vector represents the state of the atmosphere or the ocean at time t_0 . Therefore \mathcal{H} includes an integration of the model to the observation times and an interpolation to the observation locations.

A useful degree of preconditioning is achieved by rewriting the problem in terms of the transformed control vector $\chi = \mathbf{L}^{-1}(\mathbf{x} - \mathbf{x}_b)$, where $\mathbf{B} = \mathbf{L}\mathbf{L}^T$. Since the dimension of the control vector χ is typically much larger than the dimension of the observation vector \mathbf{y}_o , this transformation gives a Hessian matrix which has many eigenvalues equal to one. The transformed problem is then written as

$$\mathcal{J}_{\chi}(\chi) = \chi^T \chi + \delta^T \mathbf{O}^{-1} \delta \quad (4)$$

where

$$\delta = (\mathcal{H}(\mathbf{x}) - \mathbf{y}_o). \quad (5)$$

For the atmosphere and the ocean at basin scales the spatial error correlations of the background (a priori) estimate imply that the observations have little influence on the solution (called *analysis*) at the smallest scales, provided the observations are well separated. For this reason, and to reduce the computational burden of computing the analysis in operational numerical weather prediction, the minimization problem is usefully modified to its incremental version (Courtier et al., 1994):

$$\mathcal{J}_\chi^{\text{inc}}(\chi) = \chi^T \chi + \delta^T \mathbf{O}^{-1} \delta \quad (6)$$

where

$$\delta = (\mathcal{H}(\mathbf{x}_b) + \mathcal{H}'_S(\mathbf{S}\mathbf{x}_b) \cdot (\mathbf{x} - \mathbf{x}_b) - \mathbf{y}_o). \quad (7)$$

Here, $\chi = \mathbf{L}^{-1}\mathbf{S}(\mathbf{x} - \mathbf{x}_b)$ where \mathbf{S} is a simplification operator (e.g. a projection to a lower spatial resolution) and $\mathcal{H}'_S(\mathbf{S}\mathbf{x}_b)$ is the derivative of the simplified equivalent of \mathcal{H} evaluated at the simplified background state, $\mathbf{S}\mathbf{x}_b$.

In this simple case study, the evolution model is the viscous Burgers' equation over a one-dimensional cyclic domain:

$$\frac{\partial x}{\partial t} + \frac{1}{2} \frac{\partial x^2}{\partial s} - \nu \frac{\partial^2 x}{\partial s^2} = 0 \quad (8)$$

where s represents the distance in meters around the 45°N constant-latitude circle. The period of the domain is roughly 28.3×10^6 m. The diffusion coefficient is set to $10^5 \text{ m}^2 \text{ s}^{-1}$, as in the experiments by Fisher and Courtier (1995).

The corresponding numerical model used in the calculation of $\mathcal{H}(\mathbf{x})$ represents \mathbf{x} as a vector of discrete Fourier coefficients. The nonlinear term is calculated using the transform method on a non-aliasing grid. A conditionally stable (Matsuno) predictor-corrector scheme is used for the time integration.

For each element of the vector of observations, \mathbf{y}_o , the corresponding element of $\mathcal{H}(\mathbf{x})$ (resp. $\mathcal{H}'_S(\mathbf{S}\mathbf{x}_b) \cdot (\mathbf{x} - \mathbf{x}_b)$) is the result of a model (resp. linearized model) integration followed by an inverse Fourier transform, followed by linear interpolation from the two grid points closest to the observation.

The observations are specified at regular time intervals at fixed locations which are chosen randomly, but with a probability distribution resembling the longitudinal density of radiosonde stations between 30°N and 60°N. The observed values are calculated by applying \mathcal{H} to \mathbf{x}_t and then adding a random error. The observation errors are normally distributed, uncorrelated and all observation errors have the same variance. The covariance matrix of observation error, \mathbf{O} , is therefore proportional to the identity matrix.

The specification of the preconditioning operator \mathbf{L} defining the background error covariance matrix \mathbf{B} is chosen to resemble that formerly used operationally at Météo-France and ECMWF (Courtier et al., 1993).

The scalar product $\langle \cdot, \cdot \rangle$ used for the definition of the gradient and to perform the minimization is given by

$$\langle \mathbf{x}_1, \mathbf{x}_2 \rangle = \mathbf{x}_1^T \mathbf{P} \mathbf{x}_2 \quad (9)$$

where \mathbf{P} is a diagonal matrix with diagonal elements equal to 2, except that corresponding to the constant term of the Fourier decomposition, which is set to 1.

3. Update formulae

This section details the four diagonal-preconditioner update formulae studied: three of them are taken from Gilbert and Lemaréchal (1989), and the remaining one is a generalized version of the quasi-Cauchy update formula (Zhu et al., 1999).

Inverse BFGS formula

The inverse BFGS diagonal-preconditioner update formula is given by Eq. (4.6) in Gilbert and Lemaréchal (1989). The i -th updated diagonal component is:

$$D_+^{(i)} = D^{(i)} + \left(\frac{1}{\langle \mathbf{y}, \mathbf{s} \rangle} + \frac{\langle \mathbf{D}\mathbf{y}, \mathbf{y} \rangle}{\langle \mathbf{y}, \mathbf{s} \rangle^2} \right) \langle \mathbf{s}, \mathbf{e}_i \rangle^2 - \frac{2D^{(i)} \langle \mathbf{y}, \mathbf{e}_i \rangle \langle \mathbf{s}, \mathbf{e}_i \rangle}{\langle \mathbf{y}, \mathbf{s} \rangle} \quad (10)$$

It is obtained as the i -th diagonal component of the matrix resulting from updating \mathbf{D} with the inverse BFGS formula. Here $(\mathbf{e}_i)_{1 \leq i \leq n}$ is an orthonormal basis of \mathbb{R}^n for the scalar product $\langle \cdot, \cdot \rangle$.

Direct BFGS formula

The direct BFGS diagonal-preconditioner update formula corresponds to Eq. (4.7) in Gilbert and Lemaréchal (1989):

$$D_+^{(i)} = \left(\frac{1}{D^{(i)}} + \frac{\langle \mathbf{y}, \mathbf{e}_i \rangle^2}{\langle \mathbf{y}, \mathbf{s} \rangle} - \frac{(\langle \mathbf{s}, \mathbf{e}_i \rangle / D^{(i)})^2}{\langle \mathbf{D}^{-1}\mathbf{s}, \mathbf{s} \rangle} \right)^{-1} \quad (11)$$

It results from taking the inverse of the diagonal of the matrix obtained by updating \mathbf{D}^{-1} with the direct BFGS formula.

Inverse DFP formula

The inverse DFP diagonal-preconditioner update formula is given by Eq. (4.8) in Gilbert and Lemaréchal (1989):

$$D_+^{(i)} = D^{(i)} + \frac{\langle \mathbf{s}, \mathbf{e}_i \rangle^2}{\langle \mathbf{y}, \mathbf{s} \rangle} - \frac{(D^{(i)} \langle \mathbf{y}, \mathbf{e}_i \rangle)^2}{\langle \mathbf{D}\mathbf{y}, \mathbf{y} \rangle} \quad (12)$$

It is the diagonal of the matrix obtained by updating \mathbf{D} with the inverse DFP formula.

Quasi-Cauchy formula

The quasi-Cauchy diagonal-preconditioner update formula is an extension of Eq. (9) in Zhu et al. (1999) to the case of a general metric defined by the scalar product $\langle \cdot, \cdot \rangle$. It may be written as

$$\mathbf{D}_+ = \begin{cases} \mathbf{D} & \text{if } \langle \mathbf{D}\mathbf{y}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{s} \rangle \\ (\mathbf{I} + \nu\mathbf{G})^{-2}\mathbf{D} & \text{if } \langle \mathbf{D}\mathbf{y}, \mathbf{y} \rangle \neq \langle \mathbf{y}, \mathbf{s} \rangle \end{cases} \quad (13)$$

where \mathbf{G} is the diagonal matrix whose i -th diagonal component is $\langle \mathbf{y}, \mathbf{e}_i \rangle^2$, and ν is the largest solution of $F(\nu) = \langle \mathbf{y}, \mathbf{s} \rangle$ with

$$F(\nu) = \langle (\mathbf{I} + \nu\mathbf{G})^{-2}\mathbf{D}\mathbf{y}, \mathbf{y} \rangle. \quad (14)$$

This diagonal-preconditioner update formula is obtained by solving the minimization problem

$$\min \langle \mathbf{w}, \mathbf{w} \rangle \quad \text{such that } \langle (\mathbf{D}^{1/2} + \Omega)^2\mathbf{y}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{s} \rangle \quad (15)$$

where Ω is a diagonal matrix whose nonzero components are given by the corresponding components of the vector \mathbf{w} .

Scaling

As shown in Gilbert and Lemaréchal (1989), a large number of iterations and function/gradient evaluations may be saved by scaling the diagonal matrix \mathbf{D} before updating it, that is multiplying it by $\langle \mathbf{y}, \mathbf{s} \rangle / \langle \mathbf{D}\mathbf{y}, \mathbf{y} \rangle$ so that the diagonal matrix to be updated satisfies the quasi-Cauchy relation.

The impact of this scaling on the quality of the L-BFGS inverse Hessian approximation will be assessed for the first three diagonal-preconditioner update formulae. This is irrelevant for the quasi-Cauchy diagonal-preconditioner update formula (13): scaling the diagonal to be updated would result in using the identity matrix multiplied by the Oren-Spedicato factor at each iteration.

4. The quadratic case

The quality of the L-BFGS inverse Hessian approximation is first studied for the incremental formulation of 4D-Var (Eq. (6)), an unconstrained quadratic minimization problem. The simplification operator consists in using a lower spectral truncation (number of terms retained in the discrete Fourier series). To enable a large number of computations, the dimension of the high-resolution fields \mathbf{x} is taken to be 258. The low-resolution fields $\mathbf{S}\mathbf{x}$ and the control variable have dimension 130. Other parameters are summarized in Table 1. Some values need some comments: ε_m is the machine epsilon, approximatively equal to 2.220×10^{-16} since all the studies are performed using IEEE 64-bit floating-point

Table 1. Parameters for the quadratic problem.

Parameter	Value
Number of (\mathbf{s}, \mathbf{y}) couples	5
Minimum l_∞ distance between successive iterates	ε_m
Wolfe's line-search parameters	$\alpha = 10^{-4}$ $\beta = 0.9$
Max. number of iterations	120
Max. number of function/gradient evaluations	144
Expected decrease at first iteration	$\mathcal{J}_x^{inc}(0)/2$
Minimum expected final/initial gradient norm ratio	$\sqrt{\varepsilon_m}$

arithmetic. The expected decrease at the first iteration is used to obtain an estimate of the step-size at the first iteration, taken equal to 2 times this value divided by the initial-gradient norm squared (defined by the scalar product $\langle \cdot, \cdot \rangle$ mentioned in Section 2). The minimum ratio of the final to the initial gradient norms is the convergence criterion used for the minimization.

The quality of the L-BFGS inverse Hessian approximation is assessed by computing the eigen-spectrum of

$$\mathbf{H}_{true}^{-1} - \mathbf{H}_{L-BFGS}^{-1} \quad (16)$$

where \mathbf{H}_{true}^{-1} is the true Hessian computed using a second-order adjoint method (Wang et al., 1992; Le Dimet et al., 1997), and \mathbf{H}_{L-BFGS}^{-1} is the L-BFGS Hessian approximation built as in Veersé (to appear).

In order to have a relative measure of this quality, the eigen-spectrum of

$$\mathbf{I} - \mathbf{H}_{true}^{-1} \mathbf{H}_{L-BFGS} \quad (17)$$

is also computed. Since \mathbf{H}_{true}^{-1} and \mathbf{H}_{L-BFGS} are symmetric to a high accuracy, the eigen-spectrum of $\mathbf{I} - \mathbf{H}_{L-BFGS} \mathbf{H}_{true}^{-1}$ is almost identical to the latter one, to which we thus restrict our attention.

Both eigen-spectra are computed using an Implicit Restarted Arnoldi method (Lehoucq et al., 1997).

Although we are interested primarily in the quality of the L-BFGS inverse Hessian approximation, the efficiency of the minimization algorithm is also of concern and will be measured by the number of iterations needed to achieve convergence and the corresponding number of simulations (evaluations of the cost function and its gradient).

4.1. The full-memory case

Before studying the limited-memory versions, the convergence of the full-memory BFGS was checked. This was done by using L-BFGS with a value of the storing index m greater

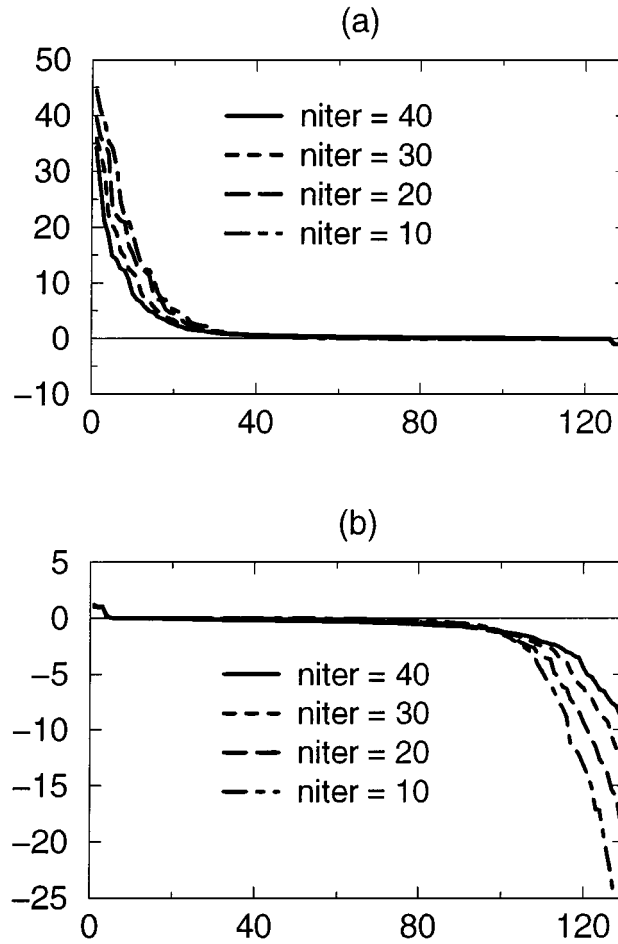


Figure 1. Eigen-spectra of operators for various iteration numbers of the minimization algorithm in the full-memory case. (a) Hessian difference operator $\mathbf{H}_{true}^{-1} - \mathbf{H}_{L-BFGS}^{-1}$. (b) Relative operator $\mathbf{I} - \mathbf{H}_{true}^{-1} \mathbf{H}_{L-BFGS}$.

than the total number of iterations required (equal to 40). Figure 1 shows the eigen-spectra of (16) and (17) for various iteration indices; the convergence is evidenced.

4.2. Choice of the (s, y) pair

A degree of freedom is given for the choice of the vectors \mathbf{s} and \mathbf{y} used to update the diagonal preconditioner. The L-BFGS algorithm (Eq. (2)) shows that the diagonal preconditioner should resemble \mathbf{H}_{k-m-1} as much as possible. This suggests to use the pair that is about to be dropped—hereafter called the *oldest pair*— $(\mathbf{s}_{k-m}, \mathbf{y}_{k-m})$ for updating the diagonal matrix \mathbf{D}_k in order to obtain \mathbf{D}_{k+1} . For the first m iterations the initial diagonal matrix (usually the

identity matrix scaled by the Oren-Spedicato factor) is to be used. This is consistent with the fact that all the curvature information from the last m iterations is completely retained in the (\mathbf{s}, \mathbf{y}) couples. However one of us advocated that the diagonal preconditioner should be updated with the newly computed couple of vectors $(\mathbf{s}_k, \mathbf{y}_k)$ —hereafter called *the newest pair*. This puts more weight onto the most recent curvature information, whose quality will improve as the minimization proceeds in the non-quadratic case.

To study the impact of both choices on the quality of the inverse Hessian approximation, the update formulae of Section 3 are implemented without scaling. The corresponding eigen-spectra are shown in figure 2. It is clear that the use of the newest pair provides an inverse Hessian approximation of far better quality than does that of the oldest pair. Moreover the panel (d) suggests that the direct BFGS and quasi-Cauchy update formulae are more accurate than the other two.

Table 2 shows the corresponding number of iterations and simulations (joint evaluations of the function and its gradient) needed to achieve convergence. Except for the direct BFGS update formula, using the newest pair increases the number of simulations, even if

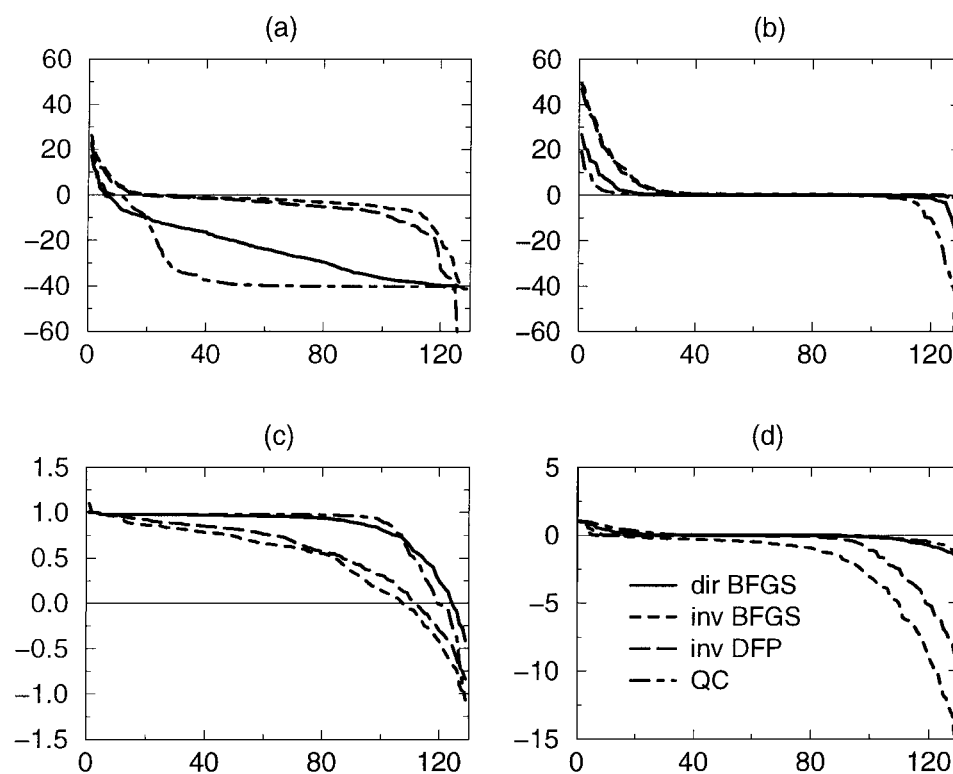


Figure 2. Eigen-spectra of the Hessian difference operator $\mathbf{H}_{true}^{-1} - \mathbf{H}_{L-BFGS}^{-1}$ for various update formulae using (a) the oldest and (b) newest pair respectively. (c) and (d) show the corresponding eigen-spectra of the relative operator $\mathbf{I} - \mathbf{H}_{true}^{-1} \mathbf{H}_{L-BFGS}$ for the oldest-pair and newest-pair cases respectively.

Table 2. Performance without scaling, using the oldest or newest pair for updating the diagonal preconditioner (# iterations/# simulations).

Formula	Oldest pair used	Newest pair used
Direct BFGS	77/78	40/52
Inverse BFGS	69/70	42/84
Inverse DFP	68/69	42/84
Quasi-Cauchy	96/103	120/130

it reduces the number of iterations for inverse BFGS and inverse DFP. The quasi-Cauchy update formula has the poorest performance on this quadratic problem. Using the newest pair with direct BFGS is beneficial both for approximating the inverse Hessian and for the minimization.

4.3. Impact of scaling

Gilbert and Lemaréchal (1989) show that scaling the diagonal matrix to force it to satisfy the quasi-Cauchy relation may lead to a better performance of the minimization. The M1QN3 minimization code of the MODULOPT library from INRIA implements the direct BFGS formula using the newest pair and scaling the diagonal before updating it. This leads to the equivalent scaled formula (4.9) in Gilbert and Lemaréchal (1989).

The impact of such a scaling for our case-study problem is now assessed for the first three diagonal-preconditioner update formulae of Section 3. This scaling is irrelevant for the quasi-Cauchy update formula, as already mentioned. Both options of scaling the diagonal matrix before or after updating it are studied. Figure 3 shows the eigenspectra of (16) and (17) for both cases. A comparison with the panels (b) and (d) of figure 2 clearly shows, for both options, a detrimental impact of the scaling onto the quality of the L-BFGS inverse Hessian approximation.

The impact in terms of minimization performance may be assessed from Table 3. Scaling leads to some reduction in the number of function and gradient evaluations needed to converge. Scaling the diagonal preconditioner after updating it is more efficient than doing it before, especially for the inverse BFGS formula. All three methods perform more or less equivalently when the scaling is done after updating, with a slight advantage for the direct BFGS diagonal-preconditioner update formula.

Table 3. Performance using the newest pair, when scaling the diagonal preconditioner before or after updating it (# iterations/# simulations).

Formula	Scaling before updating	Scaling after updating
Direct BFGS	47/49	47/49
Inverse BFGS	55/60	51/53
Inverse DFP	52/53	50/52

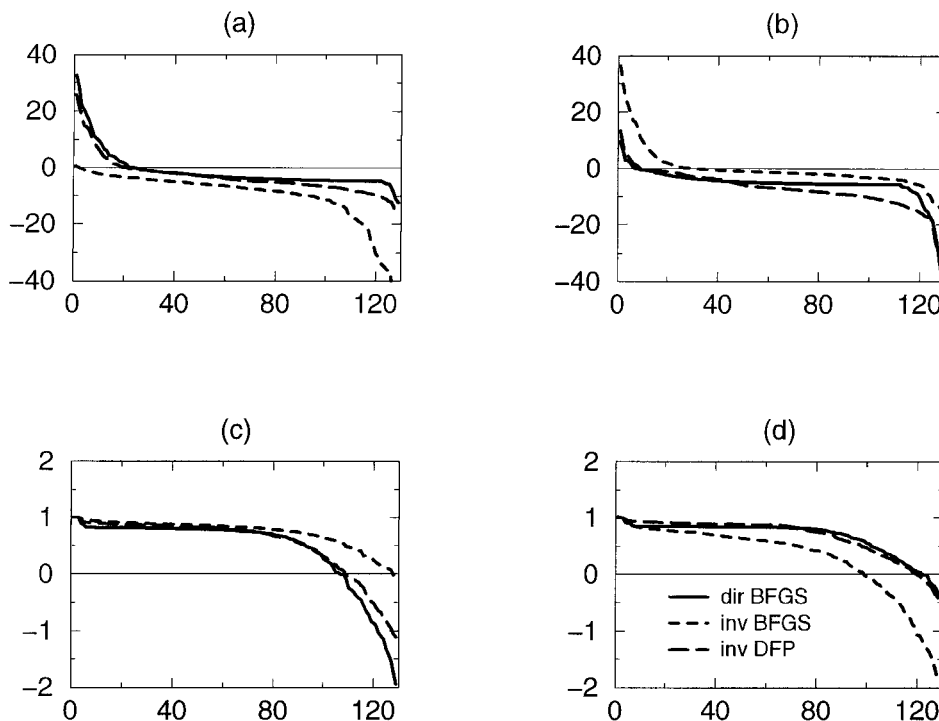


Figure 3. Eigen-spectra of the Hessian difference operator $\mathbf{H}_{true}^{-1} - \mathbf{H}_{L-BFGS}^{-1}$ when scaling the diagonal preconditioner (a) before updating it and (b) after updating it. (c) and (d) show the eigen-spectra of the relative operator $\mathbf{I} - \mathbf{H}_{true}^{-1} \mathbf{H}_{L-BFGS}$ corresponding to (a) and (b) respectively.

4.4. A new approach

Let us summarize our findings up to now. To obtain a good approximation of the inverse Hessian the newest pair should be used to update the diagonal preconditioner. But this reduces the minimization performance, except for direct BFGS. It is recommended to scale the diagonal preconditioner after updating, as this always leads to some improvements for the minimization. However the quality of the L-BFGS inverse Hessian approximation is then largely damaged.

It is natural therefore to consider a new approach where a scaled diagonal preconditioner is used for the minimization but the original (unscaled) one is updated. Figure 4 shows the corresponding eigenspectra of (16) and (17). As could be expected, an inverse Hessian quality similar to that of panels (b) and (d) in figure 2 is recovered. The possible differences between both figures may be explained by the different sequences of iterates generated.

Table 4 gives the corresponding impact on the minimization performance, in terms of iterations and simulations. This approach gives a further improvement in terms of simulations required for all three update formulae. All three perform nicely, with a slight advantage for the direct BFGS one.

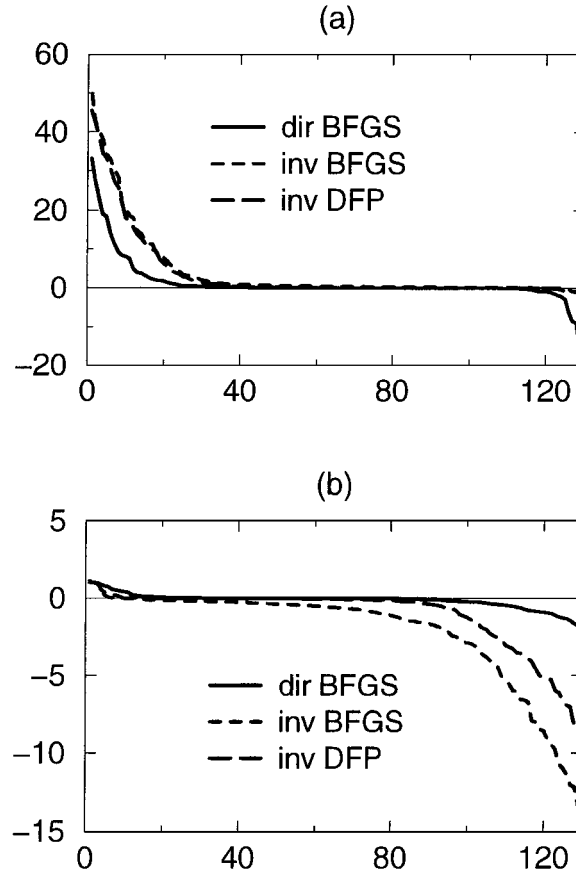


Figure 4. Eigen-spectra of (a) the Hessian difference operator $\mathbf{H}_{true}^{-1} - \mathbf{H}_{L-BFGS}^{-1}$ and (b) the relative operator $\mathbf{I} - \mathbf{H}_{true}^{-1}\mathbf{H}_{L-BFGS}$ when a scaled diagonal preconditioner is used for the minimization but the unscaled one is updated.

4.5. Impact of m

The focus is now put on how the quality of the L-BFGS inverse Hessian approximation and the minimization performance are affected by a change in the number m of (\mathbf{s}, \mathbf{y}) couples used. This is studied using our best configuration, namely updating the unscaled diagonal preconditioner with direct BFGS using the newest pair but using its scaled version for the minimization. Figure 5 shows the corresponding eigenspectra. There is not much difference on the quality of the L-BFGS inverse Hessian approximation for different values of m . However the approximation tends to improve with increasing values of the storing index.

The corresponding minimization performance are given in Table 5. The performance increases with increasing values of the storing index and the usual optimum value between 3 and 20 is not found with this diagonal-preconditioner update strategy. This is related to the corresponding increase in the quality of the L-BFGS inverse Hessian approximation.

Table 4. Performance using the newest pair, when using a scaled diagonal preconditioner for the minimization but updating the original unscaled one.

Formula	Iterations/Simulations
Direct BFGS	40/43
Inverse BFGS	44/46
Inverse DFP	43/46

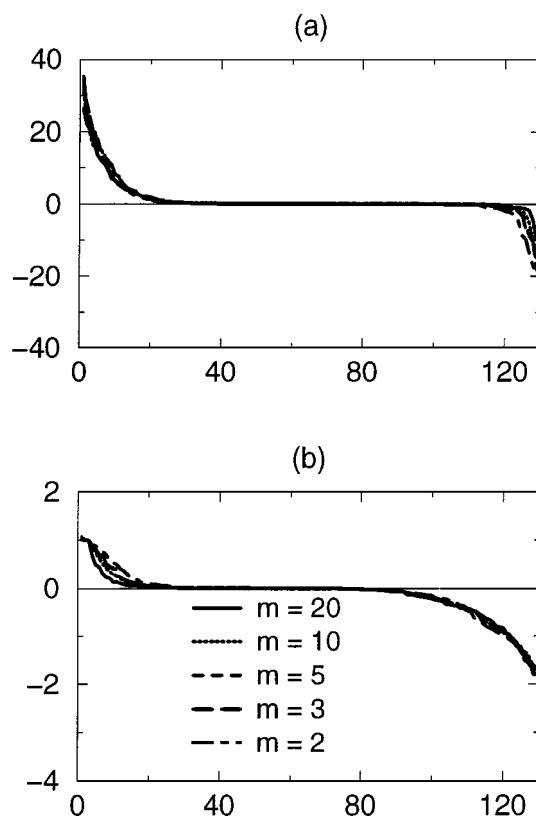


Figure 5. Eigen-spectra of (a) the Hessian difference operator $\mathbf{H}_{true}^{-1} - \mathbf{H}_{L-BFGS}^{-1}$ and (b) the relative operator $\mathbf{I} - \mathbf{H}_{true}^{-1} \mathbf{H}_{L-BFGS}$ for different values of the storing index m . The diagonal-preconditioner update formula is direct BFGS with the new approach.

4.6. Evolution during the minimization

In variational data assimilation, an approximation of the error covariances of the sought initial condition of the model is provided by the inverse Hessian at the minimum (Veersé, 1999, to appear). However due to the corresponding computational burden for realistic VDA problems in meteorology and oceanography the minimization is usually stopped

Table 5. Minimization performance for various storing indices. The diagonal preconditioner is updated using the direct BFGS formula with the new approach.

m	Iterations/Simulations
2	44/48
3	42/46
5	40/43
10	36/38
20	35/37

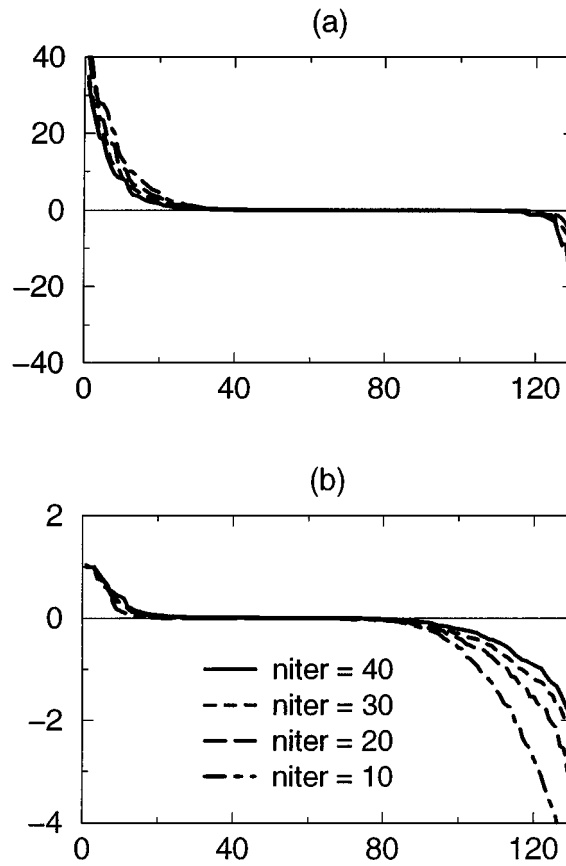


Figure 6. Eigen-spectra of (a) the Hessian difference operator $\mathbf{H}^{-1} - \mathbf{H}_{L-BFGS}^{-1}$ and (b) the relative operator $\mathbf{I} - \mathbf{H}_{true}^{-1} \mathbf{H}_{L-BFGS}$ for different values of the iteration index. The diagonal preconditioner is updated using direct BFGS with the new approach.

before reaching convergence, after a few tens of iterations have been performed. For this reason, it is interesting to see how the quality of the L-BFGS inverse Hessian approximation evolves with increasing iteration indices. Figure 6 shows the evolution of the corresponding eigen-spectra, using our best-case diagonal-preconditioner update formula with $m = 5$ (\mathbf{s}, \mathbf{y}) couples. Clearly, the quality of the approximation improves as the minimization proceeds. This is a natural evolution since the dimension of the subspace explored during the minimization increases whereas the true Hessian is constant for the present quadratic problem.

5. The non-quadratic problem

All the experiments of the previous section have been performed again using the standard formulation of 4D-Var (4), a non-quadratic minimization problem. The dimension of the control variable is 258, identical to the model phase space dimension. Now the Hessian depends on the point where it is evaluated. The quality of the various L-BFGS inverse Hessian approximations was assessed with respect to the Hessians computed with second-order adjoint techniques at the corresponding computed optimal points. The parameters used for the minimization are again those of Table 1, except the maximum numbers of iterations and simulations allowed have been increased to 200 and 250 respectively. The results do not differ qualitatively from those of the quadratic case and the plots of the corresponding eigen-spectra (not shown) are similar to those of the previous section.

A significant difference however is the failure of the first three diagonal-preconditioner update methods when the newest pair is used but no scaling is applied. This occurs at the second iteration during the model integration used for the cost-function and gradient computations. A likely explanation is the generation of an iterate during the minimization that after a few time steps leads to a violation of the stability condition of the model, which then explodes numerically.

The minimization performance may be assessed from Table 6. The standard (non-quadratic) 4D-Var cost function (Eq. (4)) requires a few more iterations and simulations to be minimized than its incremental (quadratic) approximation (Eq. (6)). The main conclusion is the same as for the quadratic case: using the direct BFGS diagonal-preconditioner update formula with the new scaling approach provides both the best performance for the minimization and the best inverse Hessian approximation.

Table 6. Minimization performance for the non-quadratic cost function (# iterations/# simulations).

Formula	No scaling oldest pair	No scaling newest pair	Scaling before newest pair	Scaling after newest pair	New approach newest pair
Direct BFGS	78/79	Failed	52/56	47/49	43/45
Inverse BFGS	76/77	Failed	63/67	69/71	48/52
Inverse DFP	74/75	Failed	55/57	56/58	48/51
Quasi-Cauchy	98/139 ¹	188/248 ²			

6. Discussion

A simple variational data assimilation problem was used as a case-study to assess the impact of various strategies for scaling and updating the L-BFGS diagonal preconditioner, both on the quality of the L-BFGS inverse Hessian approximation and on the minimization performance. The former was evaluated from comparison with the Hessian provided by second-order adjoint techniques, using eigen-decompositions. This approach is feasible only with relatively small-size problems, due to the computational burden of computing these eigen-decompositions.

The minimization performance was measured in terms of number of iterations and simulations required to achieve convergence. Both points of view lead to a few constataions:

- Using the newest (\mathbf{s}, \mathbf{y}) pair to update the diagonal preconditioner gives a better inverse Hessian approximation and, except for quasi-Cauchy, requires less simulations. One should be reminded however that the corresponding computations failed for the non-quadratic problem in the absence of scaling.
- The quasi-Cauchy diagonal-preconditioner update formula was first implemented using Newton-Raphson's unidimensional root-finding algorithm, but there were some failures related to the difficulty of specifying a good problem-independent first estimate of the root. As in Zhu et al. (1999), it was finally implemented using bisection. Quasi-Cauchy performs worse in this case study than the formulae proposed in Gilbert and Lemaréchal (1989). This suggests that the latter was able to accumulate some useful information on the inverse Hessian.
- As in Gilbert and Lemaréchal (1989) scaling the diagonal preconditioner so that it satisfies the quasi-Cauchy relation improves the performance of the minimization, especially when this scaling is done after updating it. However the inverse Hessian approximation is largely damaged by such a scaling.

This fact led us to propose a new approach, where a scaled version of the diagonal preconditioner is used for the minimization but the original (unscaled) one is updated, using the newest pair. This approach allows a good approximation of the inverse Hessian while improving further the minimization performance. It has also some "natural" properties: increasing the storage leads to better minimization performance and inverse Hessian approximation, and the latter improves steadily during the minimization process. These results were obtained both for the quadratic and for the non-quadratic variational data assimilation problems.

However the improvements of the L-BFGS inverse-Hessian quality and the reduction of simulations needed to achieve convergence may well be specific to the problems studied. The proposed approach was also assessed for a large number of unconstrained problems. An improvement was obtained for the MODULOPT and MINPACK-2 problems, but the tests with the CUTE library evidenced a lack of robustness of the method (Veersé and Auroux, 2000).

Our original concern was the quality of the L-BFGS inverse Hessian approximation. Not surprisingly the diagonal-preconditioners that provide good-quality inverse Hessian approximations are also those that lead to good minimization performances, when the

proposed scaling approach is used. The methodology employed in the present study may thus be worthwhile when designing future diagonal-preconditioner update formulae.

Acknowledgments

The first author would like to dedicate this paper to Prof. P. Fabrie of Université de Bordeaux who taught him a few years ago the basics of optimization theory.

References

- C. G. Broyden, "A new double-rank minimization algorithm," *Notices American Math. Soc.* vol. 16, p. 670, 1969.
- P. Courtier, E. Andersson, W. Heckley, G. Kelly, J. Pailleux, F. Rabier, J.-N. Thépaut, P. Uden, D. Vasiljevic, C. Cardinali, J. Eyre, M. Hamrud, J. Haseler, A. Hollingsworth, A. McNally, and A. Stoffelen, "Variational assimilation at ECMWF," ECMWF Research Dept. Tech. Memo. vol. 194, p. 84, 1993.
- P. Courtier, J.-N. Thépaut, and A. Hollingsworth, "A strategy for operational implementation of 4D-Var, using an incremental approach," *Quart. J. Roy. Meteor. Soc.* vol. 120, pp. 1367–1387, 1994.
- J. E. Dennis and J. J. Moré, "Quasi-Newton methods, motivation and theory," *SIAM Rev.* vol. 19, pp. 46–89, 1977.
- J.-Ch. Gilbert and C. Lemaréchal, "Some numerical experiments with variable storage quasi-Newton algorithms," *Math. Prog.* vol. 45, pp. 407–435, 1989.
- F.-X. Le Dimet, H. E. Ngodock, and B. Luong, "Sensitivity analysis in variational data assimilation," *J. Meteorol. Soc. Japan* vol. 75-1B, pp. 245–255, 1997.
- R. B. Lehoucq, D. C. Sorensen, and C. Yang, *ARPACK User's Guide: Solution of Large-Scale Eigenvalue Problems with Implicit Restarted Arnoldi Methods*, 1997, Available from <http://www.caam.rice.edu/software/ARPACK/>
- D. C. Liu and J. Nocedal, "On the limited memory BFGS method for large scale optimization," *Math. Prog.* vol. 45, pp. 503–528, 1989.
- J. Nocedal, "Updating quasi-Newton matrices with limited-storage," *Math. Comput.* vol. 35, pp. 773–782, 1980.
- F. Veersé, "Variable-storage quasi-Newton operators as inverse forecast/analysis error covariance matrices in variational data assimilation," INRIA Research Report 3685, p. 27, 1999.
- F. Veersé, "Variable-storage quasi-Newton operators for modeling error covariances," in *Proceedings of the Third WMO International Symposium on Assimilation of Observations in Meteorology and Oceanography, 7–11 June 1999*, Quebec City, Canada, World Meteorological Organization, Geneva, to appear.
- F. Veersé and D. Auroux, "Some numerical experiments on scaling and updating L-BFGS diagonal preconditioners," INRIA Research Report 3858, p. 25, 2000.
- Z. Wang, I. M. Navon, F.-X. Le Dimet, and X. Zou, "The second order adjoint analysis: Theory and applications," *Meteorol. Atmos. Phys.* vol. 50, pp. 3–20, 1992.
- P. Wolfe, "Convergence conditions for ascent methods," *SIAM Rev.* vol. 11, no. 2, pp. 226–235, 1969.
- M. Zhu, J.-L. Nazareth, and H. Wolkowicz, "The quasi-Cauchy relation and diagonal updating," *SIAM J. Optim.* vol. 9, pp. 1192–1204, 1999.