

1 **The SEISCOPE Optimization Toolbox: A large-scale**
2 **nonlinear optimization library based on reverse**
3 **communication**

4 **Ludovic Métivier^{*,†}, Romain Brossier[†]**

5 ^{*} *LJK, Univ. Grenoble Alpes, CNRS, France. E-mail: ludovic.metivier@ujf-grenoble.fr.*

6 [†] *ISTerre , Univ. Grenoble Alpes, France. E-mail: romain.brossier@ujf-grenoble.fr*

7 Peer-reviewed code related to this article can be found at <http://software.seg.org/2016/0001>

8 Code number: 2016/0001

9 (December 2, 2015)

10 Running head: **The SEISCOPE Optimization Toolbox**

ABSTRACT

11 The SEISCOPE optimization toolbox is a set of FORTRAN 90 routines which implement first-
12 order methods (steepest-descent, nonlinear conjugate gradient) and second-order methods
13 (*l*-BFGS, truncated Newton), for the solution of large-scale nonlinear optimization prob-
14 lems. An efficient linesearch strategy ensures the robustness of these implementations. The
15 routines are proposed as black boxes easy to interface with any computational code where
16 such large-scale minimization problems have to be solved. Travel-time tomography, least-
17 squares migration, or full waveform inversion are examples of such problems in the context
18 of geophysics. Integrating the toolbox for solving this class of problems presents two ad-
19 vantages. First, it helps to separate the routines depending on the physics of the problem
from the ones related to the minimization itself, thanks to the reverse communication pro-
tocol. This enhances flexibility in code development and maintenance. Second, it allows

1 us to switch easily between different optimization algorithms. In particular, it reduces the
2 complexity related to the implementation of second-order methods. As the latter benefit
3 from faster convergence rates compared to first-order methods, significant improvements in
4 terms of computational efforts can be expected.

INTRODUCTION

1 Nonlinear optimization problems are ubiquitous in geophysics. Methods such as travel-
2 time tomography (Nolet, 1987), least-squares migration (Lambaré et al., 1992; Nemeth
3 et al., 1999), or full waveform inversion (FWI) (Virieux and Operto, 2009), are based on
4 the minimization of an objective function which measures the misfit between simulated
5 and observed data from seismic. The simulated data (travel-time, partial measurement of
6 the pressure and/or displacement velocity wavefields) depends on subsurface parameters
7 (P- or S-wave velocities, density, anisotropy parameters, reflectivity among others). The
8 minimization of the objective function consists in computing the subsurface parameters
9 which best explain the observed data. In the general case, the simulated data depends
10 nonlinearly on the subsurface parameters, yielding a nonlinear optimization problem.

11 Two general classes of methods exist to solve this type of problem: global and semi-
12 global methods (Sen and Stoffa, 1995; Spall, 2003) or local descent algorithms (Nocedal
13 and Wright, 2006; Bonnans et al., 2006). The methods belonging to the first class have
14 the capability to find the global minimum of an objective function from any given starting
15 point. They are said to be globally convergent. These methods proceed in a guided-random
16 sampling of the parameter space to find the global minimum of the objective function.
17 However, as soon as the number of discrete parameters exceeds a few tenths to hundreds, the
18 number of evaluations of the objective function required to find its global minimum becomes
19 excessively important for these methods to be applied in a reasonable time, even using
20 modern high performance computing (HPC) platforms. In geophysics, realistic applications
21 yield optimization problems for which the number of unknowns parameters is easily several
22 orders of magnitude higher (billions may be involved for 3D FWI applications), rendering

1 global approaches intractable.

2 Local descent algorithms only guarantee local convergence: from a given starting point,

3 the nearest local minimum of the objective function is found. Instead of proceeding to a

4 random sampling, these methods are based on a guided exploration of the model space

5 following the descent directions of the objective function. From an initial guess, an estima-

6 tion of the local minimum is found by successive updates following these directions. The

7 convergence and the monotonic decrease of the objective function is guaranteed by a proper

8 linesearch or trust-region strategy (Nocedal and Wright, 2006). This consists in giving an

9 appropriate scale to the update brought to the current model estimate. This guided ex-

10 ploration strategy is drastically cheaper than a random sampling of the model space: the

11 number of iterations required to converge to a minimum is far lower than the number of

12 evaluations of the objective function which would be required to obtain an acceptable sam-

13 pling of the model space when its dimension is large. This is the reason why local descent

14 methods are used for solving large scale optimization problems.

15 Among possible local descent methods, the simplest ones use only the first-order deriva-

16 tives of the objective function (the gradient) to define the descent direction. These methods

17 are referred to as first-order methods. The steepest-descent and nonlinear conjugate gradi-

18 ent algorithms are examples of these methods. More elaborated descent schemes account

19 for the curvature of the objective function through quasi-Newton algorithms. These al-

20 gorithms are based on approximations of the inverse Hessian operator (the matrix of the

21 second-order derivatives of the objective function) following three different (possibly com-

22 plementary) strategies.

23 The first strategy is referred to as the preconditioning strategy. In some situations,

1 an approximation of the inverse Hessian operator can be computed from analytic or semi-
2 analytic formulas. Diagonal approximations can be used. An example in the context of
3 FWI is the pseudo-Hessian preconditioning strategy proposed by Shin et al. (2001). This
4 strategy can be used to derive an approximate inverse of the Hessian operator, which is
5 integrated within first-order algorithms by multiplying the gradient by this approximation.
6 The second strategy is based on the computation of an approximate inverse Hessian op-
7 erator, from finite-differences of the gradient, such as in the *l*-BFGS algorithm (Nocedal,
8 1980). Instead of using an approximate inverse of the Hessian operator, the third strategy,
9 which is referred to as the truncated Newton method, solves incompletely the linear system
10 associated with the Newton equations, namely, the linear system relating the gradient to
11 the descent direction, through a conjugate gradient iterative solver (Nash, 2000).

12 In geophysics, the solution of minimization problems is often computed using basic type
13 methods such as the steepest descent or the nonlinear conjugate gradient. In addition,
14 linesearch algorithms are not always implemented, and as a consequence the convergence
15 and the monotonic decrease of the objective function is not guaranteed. Moreover, it is
16 not unusual to see computer codes in which the minimization procedure is embedded into
17 the routines related to the physical problem at stake. This type of implementation is
18 not flexible, as it implies that the optimization process has to be re-developed for each
19 application. Improving the optimization process from first-order to second-order methods
20 may also require significant modifications of the code.

21 The SEISCOPE optimization toolbox has been designed to propose a solution to these
22 issues. The first objective of the toolbox is to propose minimization routines which can be
23 easily interfaced with any computational code through a reverse communication protocol
24 (Dongarra et al., 1995). The principle is the following: the computation of the solution

1 of the optimization problem is performed in a specific routine of the code. This routine
2 is organized as a minimization loop. At each iteration of the loop, the minimization rou-
3 tine from the toolbox chosen by the user is called. This routine communicates what is the
4 quantity required at the current iteration: objective function, gradient, or Hessian-vector
5 product. These quantities are computed by the user in specific routines, external to the
6 minimization loop. The process continues until the convergence is reached. This implemen-
7 tation paradigm yields a complete separation between the code related to the physics of the
8 problem and the code related to the solution of the minimization problem. This ensures
9 a greater versatility as one can easily modify one of these parts while keeping the other
10 unchanged.

11 The second objective of the toolbox is to propose robust and efficient methods for
12 large-scale nonlinear optimization problems. Four different optimization schemes are im-
13 plemented: steepest-descent, nonlinear conjugate gradient, *l*-BFGS and truncated Newton
14 method. All these methods share the same linesearch strategy, which assures the conver-
15 gence and the monotonic decrease of the objective function, and allows to compare the
16 efficiency of the methods on a given application. In addition, the four methods can be
17 combined with preconditioning strategies: any information regarding the curvature of the
18 objective function can be easily incorporated to improve the convergence rate of the algo-
19 rithms. Finally, the toolbox offers the possibility to use second-order methods as well as
20 first-order methods. The complexity associated with the implementation of the *l*-BFGS
21 approximation is hidden to the user. The use of truncated Newton methods only requires
22 to implement Hessian-vector products.

23 It should be mentioned that a similar development in the C++ language has been provided
24 by members of the Rice project. The developments are based on RVL, the Rice Vector

1 library. The available methods are *l*-BFGS and the truncated Newton method implemented
2 in the framework of the trust-region globalization strategy, also known as the Steihaug-Toint
3 algorithm (Steihaug, 1983; Gould et al., 1999). However, to the best of our knowledge, the
4 implementation is not based on a reverse communication protocol.

5 We devote this article to the presentation of the SEISCOPE optimization toolbox. First,
6 we give an overview of the optimization methods which are implemented. No details on
7 convergence proofs and computation of convergence rates are given. Those can be found
8 in reference textbooks (Nocedal and Wright, 2006; Bonnans et al., 2006). We focus on
9 the general principles of these methods. We present a code sample in which the toolbox
10 is interfaced to illustrate the reverse communication protocol in a practical example. We
11 provide details on the implementation of the routines. In the numerical study section,
12 we investigate the use of the toolbox on two different applications. We minimize the bi-
13 dimensional Rosenbrock function using the different methods proposed in the toolbox. We
14 compare the convergence of these algorithms and the optimization path they follow. We
15 give a second illustration on a large scale nonlinear minimization problem related to a 2D
16 acoustic frequency-domain FWI case study. We consider synthetic data acquired on the
17 Marmousi 2 model. We compare the convergence of the different minimization algorithms.
18 These two examples emphasize the superiority of second-order optimization methods over
19 first-order methods. We give a conclusion and perspectives in the final section.

THEORY

¹ Generalities

The routines implemented in the SEISCOPE optimization toolbox are designed to solve unconstrained and bound constrained nonlinear minimization problems, under the general form

$$\min_{x \in \Omega} f(x), \quad (1)$$

where

$$\Omega = \prod_{i=1}^n [a_i, b_i] \subset \mathbb{R}^n, \quad n \in \mathbb{N}, \quad (2)$$

- ² and $f(x)$ is a sufficiently smooth function (at least twice differentiable) depending nonlinearly on the variable x .

All the routines considered in the toolbox belong to the class of local descent algorithms.

From an initial guess $x_0 \in \Omega$, a sequence of iterates is built following the recurrence

$$x_{k+1} = x_k + \alpha_k \Delta x_k, \quad (3)$$

- ⁴ where $\alpha_k \in \mathbb{R}_+^*$ is a steplength and $\Delta x_k \in \mathbb{R}^n$ is a descent direction. The recurrence (3) is
⁵ repeated until a certain stopping convergence criterion is reached.

- ⁶ The steplength α_k is computed through a linesearch process to satisfy the standard
⁷ Wolfe conditions*. Satisfying the Wolfe conditions ensures the convergence toward a local
⁸ minimum, provided $f(x)$ is bounded and sufficiently smooth (Nocedal and Wright, 2006).
⁹ The satisfaction of the bound constraints is ensured through the projection of each iterate
¹⁰ within the feasible domain Ω in the linesearch process.

*Standard is understood as opposed to the strong Wolfe conditions, which are more restrictive. This issue is important regarding the nonlinear conjugate gradient method implementation (see the nonlinear conjugate gradient subsection). The Wolfe conditions are described later in equations (13) and (14).

1 The different nonlinear minimization methods differ by the computation of the descent
2 direction Δx_k . First, a description of the computation of these quantities depending on the
3 minimization routines is given. More details on the linesearch and the bound constraints
4 enforcement algorithm are given after.

In the following, the gradient of the objective function is denoted by $\nabla f(x)$, and the Hessian operator by $H(x)$. A preconditioner for the Hessian operator $H(x)$ is denoted by $P(x)$, such that

$$P(x) \simeq H^{-1}(x). \quad (4)$$

5 **Steepest descent**

The simplest optimization routine implemented is the steepest descent algorithm. This method uses the following descent direction

$$\Delta x_k = -P(x_k)\nabla f(x_k). \quad (5)$$

6 The convergence rate of this method is linear. When no preconditioner is available (P is
7 the identity matrix), the descent direction is simply the opposite of the gradient. In this
8 case, examples show that the number of iterations required to reach the convergence can
9 be extremely high.

10 **Nonlinear conjugate gradient**

The conjugate gradient algorithm is an iterative linear solver for symmetric positive definite systems. This method can be interpreted as a minimization algorithm for quadratic functions. The nonlinear conjugate gradient method is conceived as an extension of this algorithm to the minimization of general nonlinear functions. The model update is computed

as a linear combination of the opposite of the gradient and the descent direction computed at the previous iteration

$$\begin{cases} \Delta x_0 = -P(x_0)\nabla f(x_0), \\ \Delta x_k = -P(x_k)\nabla f(x_k) + \beta_k \Delta x_{k-1}, \quad k \geq 1, \end{cases} \quad (6)$$

where $\beta_k \in \mathbb{R}$ is a scalar parameter. Different formulations can be used to compute β_k , giving as many versions of the nonlinear conjugate gradient algorithm. Standard implementations use formulas from Fletcher and Reeves (1964), Hestenes and Stiefel (1952) or Polak and Ribi  re (1969). In the SEISCOPE optimization toolbox, an alternative formula, proposed by Dai and Yuan (1999), is used. In this particular case, the scalar β_k is computed as

$$\beta_k = \frac{\|\nabla f(x_k)\|^2}{(\nabla f(x_k) - \nabla f(x_{k-1}))^T \Delta x_{k-1}}. \quad (7)$$

- ¹ This formulation ensures the convergence toward the nearest local minimum as soon as the
- ² steplength satisfies the Wolfe conditions. This is not the case for the standard formulations,
- ³ which require satisfying the strong Wolfe conditions. This is the reason why the formu-
- ⁴ lation from Dai and Yuan (1999) is preferred. This allows us to use the same linesearch
- ⁵ strategy for all the optimization routines proposed in the toolbox, including the nonlinear
- ⁶ conjugate gradient. The convergence of the nonlinear conjugate gradient is linear, as well
- ⁷ as for the steepest descent algorithm. In some cases, the reduction of the number of itera-
- ⁸ tions can be significant, however this potential acceleration is case dependent and therefore
- ⁹ unpredictable.

¹⁰ Quasi-Newton l -BFGS method

Among the class of quasi-Newton methods, the l -BFGS algorithm has become increasingly popular because of its simplicity and efficiency (Nocedal, 1980). The l -BFGS approximation

relies on an approximate inverse Hessian operator Q_k computed through finite-differences of l previous values of the gradient. The resulting descent direction is computed as

$$\Delta x_k = -Q_k \nabla f(x_k). \quad (8)$$

1 The equation (8) is noteworthy: it is similar in form as the equation (5), where Q_k plays the
2 role of the preconditioner $P(x_k)$. The difference comes from the fact that the preconditioning
3 matrix Q_k is computed in a systematical way, only from the objective function gradient
4 values, and updated at each iteration following the l -BFGS formula. In preconditioned
5 steepest descent algorithm, the user is in charge to provide the preconditioning matrix
6 $P(x_k)$. The l -BFGS method achieves a superlinear convergence rate, and is most of the
7 time faster than nonlinear conjugate gradient and steepest descent algorithms.

8 In terms of practical implementation, the matrix Q_k is never explicitly built. The
9 product $Q_k \nabla f(x_k)$ is directly computed instead, following a two-loop recursion algorithm
10 (Nocedal and Wright, 2006). Interestingly, it is possible to incorporate an estimation of the
11 inverse Hessian operator through the preconditioner $P(x_k)$ into this two-loop recursion. This
12 means that the information from the preconditioner $P(x_k)$ and the l -BFGS formula can be
13 combined to approximate as accurately as possible the action of the inverse Hessian operator.
14 Surprisingly, this possibility is rarely mentioned in the literature, while it seems to be one of
15 the most efficient strategy for large-scale nonlinear optimization problems. In the context of
16 FWI, the combination of a diagonal preconditioner and the l -BFGS approximation can be
17 shown to yield the best computational efficiency among the different optimization strategies
18 available from the toolbox (see the numerical results on the Marmousi 2 case study). The
19 details of this approach are given in the Appendix.

¹ **Truncated Newton method**

Instead of using an approximate inverse of the Hessian operator, the truncated Newton method computes an inexact (truncated) solution of the linear system associated with the Newton equation (Nash, 2000)

$$H(x_k)\Delta x_k = -\nabla f(x_k). \quad (9)$$

- ² This approximate solution of the linear system is computed through a matrix-free conjugate
³ gradient algorithm: only the action of the Hessian operator on a given vector is required. In
⁴ the context of FWI, the gradient is usually computed through first-order adjoint methods.
⁵ Second-order adjoint strategies can be used to compute the action of the Hessian on a
⁶ given vector (Métivier et al., 2012). In both cases, the computational complexity is reduced
⁷ to the solution of the wave equation for different source terms. This allows us to use
⁸ truncated Newton methods for solving FWI problems at a reasonable computational cost.
⁹ This strategy has been investigated in Métivier et al. (2013, 2014); Castellanos et al. (2015)

The stopping criterion for the inner linear system is

$$\|H(x_k)\Delta x_k + \nabla f(x_k)\| \leq \eta_k \|\nabla f(x_k)\|, \quad (10)$$

- ¹⁰ where η_k is a forcing term. The equation (10) reveals that η_k controls the accuracy of the
¹¹ relative residuals of the linear system (9). For large-scale nonlinear applications, particular
¹² choices for η_k are recommended, as it is detailed by Eisenstat and Walker (1994). Since
¹³ the targeted applications for this toolbox are precisely large-scale nonlinear minimization
¹⁴ problems, the default implementation of the computation of η_k is based on this work: η_k
¹⁵ should decreases as soon as the function which is minimized is locally quadratic. The
¹⁶ measure of the local quadratic approximation of the misfit function is based on the gradient

1 current and previous values. For FWI application, this strategy reveals to be efficient
2 (Métivier et al., 2013, 2014; Castellanos et al., 2015). However, for simpler problems, a
3 constant and small value for η_k might produce better results. In numerical results section,
4 for the Rosenbrock experiment, η_k is set to 10^{-5} and remains constant throughout the
5 iterations.

In the context of the truncated Newton method, the preconditioner $P(x)$ is used naturally as a preconditioner for the linear system (9). This implies that the symmetry of (9) has to be preserved through the preconditioning operation. Therefore, $P(x)$ cannot be used simply as a left or right preconditioner. The preconditioner $P(x)$ has to be symmetric definite positive. In this case, it can be factorized as

$$P(x) = C(x)^T C(x), \quad (11)$$

and the following symmetric preconditioning strategy can be used

$$C(x_k)^T H(x_k) C(x_k) \Delta y_k = -C(x_k)^T \nabla f(x_k), \quad \Delta x_k = C(x_k) \Delta y_k. \quad (12)$$

6 The symmetry of (9) is preserved as $C(x_k)^T H(x_k) C(x_k)$ is symmetric by construction.
7 In terms of implementation, the formulation (12) may imply that a factorization of the
8 preconditioner $P(x)$ under the form (11) is required. However, a particular implementation
9 of the preconditioned conjugate gradient algorithm, such as the one proposed in Nocedal
10 and Wright (2006), allows us to solve (12) through the conjugate gradient method using
11 only matrix-vector products from $H(x)$ and $P(x)$. In practice, this preconditioned conjugate
12 gradient implementation differs from the non-preconditioned conjugate gradient implemen-
13 tation only by the multiplication of the residuals by the matrix $P(x)$. Thus, from the user
14 point of view, only the computation of the action of $P(x)$ on a given vector is required.

1 The truncated Newton method has been shown to converge even faster than the l -BFGS
2 method, achieving a nearly quadratic convergence rate close to the solution. However, the
3 computation cost of each iteration is increased by the (inexact) solution of the Newton
4 equation. In practice, l -BFGS and truncated Newton methods have been shown to be com-
5 petitive on a collection of benchmark problems from the nonlinear optimization community
6 (Nash and Nocedal, 1991). The same conclusions have been drawn for mono-parameter
7 FWI applications (Métivier et al., 2013, 2014; Castellanos et al., 2015), and this is what
8 can be observed in the numerical experiments proposed in this study.

9 **Linesearch and bound constraints**

Using a linesearch strategy is necessary to ensure the robustness of the optimization routines, that is guaranteeing a monotonic decrease of the objective function and the convergence toward a local minimum. The linesearch algorithm enforces the Wolfe conditions. These are the sufficient decrease condition

$$f(x_k + \alpha \Delta x_k) \leq f(x_k) + c_1 \alpha \nabla f(x_k)^T \Delta x_k, \quad (13)$$

and the curvature condition

$$\nabla f(x_k + \alpha \Delta x_k)^T \Delta x_k \geq c_2 \nabla f(x_k)^T \Delta x_k, \quad (14)$$

where c_1 and c_2 are coefficients such that

$$0 < c_1 < c_2 \leq 1. \quad (15)$$

In practice, the values

$$c_1 = 10^{-4}, \quad c_2 = 0.9, \quad (16)$$

1 are used, following Nocedal and Wright (2006). The first condition indicates that the
2 steplength α should be computed to ensure that the update in the direction Δx_k will
3 generate a sufficient decrease of the objective function. The second condition is used to
4 rule out too small values for α itself, which would lead to undesirable small updates of the
5 estimation x_k .

The bound constraints satisfaction is enforced through a projection method, similarly to the strategy proposed by Byrd et al. (1995). At each iteration of the linesearch, before testing the two Wolfe conditions, the updated model $x_k^* = x_k + \alpha\Delta x_k$ is projected into the feasible domain Ω , through the operator $Proj(x)$ defined as

$$Proj(x)_i = \begin{cases} x_i & \text{if } a_i \leq x_i \leq b_i \\ a_i & \text{if } x_i < a_i \\ b_i & \text{if } x_i > b_i, \end{cases} \quad (17)$$

6 where the subscript i denotes the i th component of the vectors of \mathbb{R}^n . This procedure ensures
7 that the estimated model always remains in the feasible domain Ω . Although there is no
8 proof of convergence for optimization algorithms using this bound constraints enforcement,
9 the method has yielded satisfactory results in practice. In some sense, it can be viewed as a
10 simplification of the method of projection on convex sets (POCS) proposed in the context
11 of seismic inversion by Baumstein (2013).

IMPLEMENTATION

12 The reverse communication protocol on which are based all the routines of the SEISCOPE
13 optimization toolbox requires a specific implementation of the different algorithms. The
14 principle of reverse communication from the user point of view is first introduced. A spe-
15 cific code example is used, where the minimization of a function is performed with the

1 preconditioned *l*-BFGS routine from the toolbox. Next, details on the implementation of
2 the algorithms of the toolbox are given. The general algorithmic structure of the different
3 optimization methods in the toolbox is presented, as well as the structure of the linesearch
4 algorithm. Note that in the current version, all the toolbox routines are implemented in
5 single precision, with the purpose to save computational time on large-scale applications.

6 **Reverse communication protocol**

7 The principle of reverse communication is illustrated from the code sample presented in
8 Figure 1. In this example, the user minimizes an objective function $f(x)$ which corresponds
9 to the Rosenbrock function. In this example, x is a two-dimensional real vector such that
10 $x \in \mathbb{R}^2$. In general, x is a n -dimensional vector. It is implemented in an array structure.
11 The minimization is performed from an initial value x_0 which initializes x . At the end of the
12 minimization loop, x contains the solution of the minimization problem. In this example,
13 the minimization is performed with the preconditioned version of the *l*-BFGS algorithm
14 proposed in the toolbox by the routine PLBFGS.

15 The code is organized in a `do while` loop. The principle is the following: while the
16 convergence has not been reached, further iterations of the minimization process are needed.
17 Each iteration of this loop starts by a call to the optimization routine. The communication
18 with the optimization routine is achieved through the four characters chain `FLAG`. This
19 character chain is an input/output variable of the optimization routine. On return of each
20 call, it is updated to communicate with the user and to require the computation of particular
21 quantities. On first call, this variable is initialized by the user to the chain of characters
22 ‘‘INIT’’.

1 Following this principle, after each call to the optimization routine, the user tests the
2 possible values for the communicator and performs the subsequent actions. If **FLAG** has
3 been set to ‘‘GRAD’’, the computation of the objective function and its gradient at the
4 current value of **x** is required. In this case, the code calls the subroutine **Rosenbrock** with
5 the current **x** as an input. On return, the subroutine **Rosenbrock** computes the objective
6 function as well as its gradient at the current **x**, and store these quantities in the variables
7 **fcost** and **grad** respectively. The code reaches the end of the **if/end if** loop and returns
8 to the minimization function **PLBFGS** with the required values for the variables **fcost** and
9 **grad**.

10 If **FLAG** is set to ‘‘PREC’’, the optimization loop requires the preconditioner to be
11 applied. In this case, the user extracts the vector **y** from the variable **optim**. The data type
12 of **optim** is proper to the optimization toolbox and is described in the file **optim.h**. The
13 field of interest is **optim%q_plb**. The user then applies a preconditioner to **y** through the call
14 to the routine **apply_Preco**. On return of this routine, the quantity **Py** is transferred into
15 the data structure **optim** in the field **optim%q_plb**. The code reaches the end of the **if/end**
16 **if** loop and returns to the minimization function **PLBFGS** with the required values for the
17 variable **optim%q_plb**. A more efficient strategy in terms of high performance computing
18 would consist in performing the two copies involving the vector **optim%q_plb** directly in the
19 routine **apply_Preco** to preserve the locality of the data. However, the purpose here is to
20 illustrate the use of the toolbox rather than proposing an optimized implementation.

21 When **FLAG** is set to ‘‘NSTE’’, this is an indication that the linesearch process just
22 terminated and has yielded a new iterate x_{k+1} of the minimization sequence. This new
23 iterate is stored in the variable **x**. This information can be useful to track the history
24 of the computed approximation of the solution. In this example, the user calls a routine

1 `print_information` to this purpose.

If `FLAG` is ‘‘CONV’’, the default convergence criterion implemented in the optimization routine has been reached. This convergence criterion is based on the reduction of the relative objective function. The convergence is declared when the condition

$$\frac{f(x_k)}{f(x_0)} \leq \varepsilon, \quad (18)$$

2 is satisfied, where ε is a threshold parameter initialized by the user.

3 If `FLAG` is ‘‘FAIL’’, a failure in the linesearch process has been detected. The user
4 is in charge of defining the minimization loop and is responsible for defining the actual
5 stopping criterion. This stopping criterion can be based on the information provided by the
6 optimization routine. For instance, in the example provided, the minimization loop stops as
7 soon as `FLAG` is set to ‘‘CONV’’ or to ‘‘FAIL’’. However, it is also possible to ignore this
8 information, or to complete it with additional tests to define the actual stopping criterion.

9 On exit of the `do while` loop, `x` contains the solution reached when the convergence criterion
10 is satisfied.

11 The code sample presented in this example can be used as a template for the minimiza-
12 tion of any objective function $f(x)$ for which it is possible to compute the gradient. The
13 optimization algorithm used here is the preconditioned l -BFGS method, however, this code
14 sample can be straightforwardly extended to the use of the other methods from the toolbox.

15 To this purpose, the table 2 presents the different values that can be taken by the variable
16 `FLAG` and the corresponding actions required depending on the choice of the minimization
17 routine. For practical use, more details are given in the toolbox documentation.

18 In the following subsections, details on the implementation of the toolbox are given.
19 Readers only interested in the use of the toolbox may skip the two following subsections

1 and jump directly to the numerical examples section.

2 **General algorithmic structure of the optimization routines**

3 The diagram in Figure 2 presents the general structure of the toolbox optimization routines.

4 The workflow starts with a call of the routine from the user. The value of the communicator

5 **FLAG** is first tested.

6 If it is equal to “INIT”, this means that this is a first call from the user. In this case,

7 subsequent initializations have to be performed. An initialization routine (here `PSTD_init`)

8 is thus called. This routine is not described in details here; this is where memory allocation

9 and parameters settings (such as the scalar c_1 and c_2 for the linesearch) are performed.

10 The memory allocation consists in allocating different fields in the data structure `optim`.

11 Depending on the minimization routines, certain fields may be allocated or not. The l -BFGS

12 algorithm requires specific fields related to the l -BFGS approximation to be allocated for

13 instance. After this initialization has been performed, the communicator **FLAG** is set to

14 “GRAD” and the routine returns to the user.

15 If the communicator **FLAG** has any value different from “INIT”, the linesearch algorithm

16 is called. Returning from the linesearch, the algorithm asks if the linesearch has terminated

17 or not.

18 If the linesearch has terminated, this means that a new iterate has been computed. This

19 new iterate is stored in the variable `x`. In this situation, the algorithm tests if the stopping

20 criterion (18) is satisfied. If this is the case, the communicator **FLAG** is set to “CONV” and

21 the routine returns to the user. If not, this means that a new step has been taken in the

22 minimization process, while the convergence has not been reached yet. In this case, a new

1 descent direction is computed, and the communicator **FLAG** is set to “NSTE”. This value
2 of the communicator is useful for intermediate printing of the solution for instance. The
3 routine returns to the user.

4 If the linesearch process is not terminated, this means that a new value of the objective
5 function and its gradient are required to continue the linesearch process. In this case, the
6 steplength is updated, the communicator **FLAG** is set to “GRAD”, and the routine returns to
7 the user.

8 The steepest descent, nonlinear conjugate gradient and *l*-BFGS algorithms share exactly
9 this algorithmic structure in the toolbox implementation. The only difference is related to
10 the routine called for the initialization and the computation of the descent direction.

11 The structure has to be modified as soon as the computation of the descent direction also
12 requires a direct intervention of the user. This is the case for instance for the preconditioned
13 *l*-BFGS routine. In this situation, the computation of the descent direction is split into two
14 parts, corresponding to the two-loop recursion (see Appendix for more details). Between
15 these two loops, a call to the user is required for applying the preconditioner.

16 The same is true for truncated Newton methods (with or without preconditioning).
17 The computation of the descent direction requires additional communications with the user
18 for applying the Hessian operator and the preconditioner to particular vectors, in order to
19 compute the incomplete solution of the system (12).

20 From this overview of the data structure, one can see that the linesearch algorithm plays
21 a central role in the optimization routines implementation. More details on its implemen-
22 tation are given in the next subsection.

1 Linesearch implementation

2 A diagram summarizing the linesearch algorithmic structure is presented in Figure 3. The
3 structure is also based on a reverse communication paradigm. The calling program is in
4 charge of the computation of the quantities required by the linesearch algorithm: objective
5 function and gradient computed at particular points.

The first call to the linesearch subroutine is indicated by an internal linesearch flag. In this case, a first candidate for the next model update x_{k+1} is computed. This candidate is denoted by x_k^* , computed as

$$x_k^* = \text{Proj}(x_k + \alpha_k \Delta x_k). \quad (19)$$

6 The formula (19) means that the candidate is computed as the current value x_k updated
7 in the descent direction Δx_k with a step-length α_k and projected into the set Ω to satisfy
8 the bound constraints which could be potentially imposed.

9 At this stage, the Wolfe conditions have to be tested on x_k^* . However, this requires the
10 computation of the quantities $f(x_k^*)$ and $\nabla f(x_k^*)$. The linesearch routine thus goes back to
11 the calling program with a request for computing these quantities. In particular, it indicates
12 that the linesearch did not terminate.

13 When these quantities have been computed, the calling program goes back to the line-
14 search routine. This time the internal linesearch flag indicates it is not a first call to the
15 linesearch. As a consequence, the two Wolfe conditions are tested in a sequential order. If
16 they are satisfied, the linesearch algorithm terminates successfully with a new model up-
17 date x_{k+1} which is set to the candidate value x_k^* . The internal linesearch flag is set back to
18 “first_call”.

1 If one of the two Wolfe conditions is not satisfied, the steplength α_k has to be adjusted.
2 The modification of α_k is performed following the bracketing strategy proposed in Bonnans
3 et al. (2006). Once α_k has been modified, a new candidate x_k^* is computed, and the linesearch
4 algorithm goes back to the calling program with a request for the computation of $f(x_k^*)$ and
5 $\nabla f(x_k^*)$.

6 The optimization routines implemented in the toolbox are thus based on a two-level
7 reverse communication process. The first-level is between the user and the optimization
8 routine. The second level is embedded in the optimization routine itself and corresponds
9 to the linesearch implementation. The potential complexity of such an implementation is
10 transparent for the user. When asked to compute the gradient of the objective function,
11 the user does not have to know if this computation is requested for the linesearch process
12 or the evaluation of a new descent direction.

13 In terms of practical implementation, the linesearch algorithm requires to define an
14 initial step-length α_0 . In the current implementation, this step-length is simply set to
15 1. Besides, at each iteration k of the minimization, the new step-length α_k is initialized
16 to the previous value α_{k-1} . From our experience, this approach results in the following
17 behavior. At the first iteration of the minimization process, the linesearch algorithm may
18 require several internal iterations to adjust the step-length. However, for further iterations,
19 the step-length computed at the previous iteration often directly yields a candidate which
20 satisfies the Wolfe conditions. The additional computational cost related to the linesearch
21 is thus limited to the first iteration of the minimization. Note that this cost can be reduced
22 by the user by using an appropriate scaling of the objective function. Indeed, at the first
23 iteration, for most optimization methods, the descent direction is equal to the gradient.
24 The step-length adjustment from the linesearch thus corresponds to an automatic scaling of

1 the misfit function. This process can be accelerated by a proper scaling performed directly
2 by the user.

3 **Summary of the routines implemented in the toolbox**

4 Table 1 summarizes the name of the minimization routines implemented in the toolbox
5 together with their basic properties. Although four main minimization methods are pre-
6 sented, the toolbox proposes a total of six subroutines. This is due to the implementation
7 of preconditioning. While it is straightforward for steepest-descent and nonlinear conju-
8 gate gradient, the implementation of preconditioning for *l*-BFGS and truncated Newton
9 method requires substantial modifications. Our choice is thus to implement a single version
10 for steepest-descent and nonlinear conjugate gradient, where the user only has to apply
11 the identity operator if no preconditioner is available, while two distinct versions (with and
12 without preconditioning) are implemented for *l*-BFGS and truncated Newton method.

13 In theory, first-order methods have a linear convergence rate, while second-order meth-
14 ods benefits from superlinear convergence. While the Newton method converges with a
15 quadratic rate close from the solution, *l*-BFGS and the truncated Newton method only
16 benefit from a superlinear convergence rate as only an approximation of the inverse Hes-
17 sian operator is used. The six routines require the computation of the gradient. Only the
18 truncated Newton method (with and without preconditioning) requires the computation of
19 Hessian-vector products.

20 Table 1 also summarizes the comparisons in terms of convergence speed (reduction of
21 the misfit function depending on the number of iteration) and computational efficiency (re-
22 duction of the misfit function depending on the number of gradient/Hessian-vector product

1 evaluations) observed on the two numerical examples presented in the following section.

NUMERICAL EXAMPLES

2 In this section, numerical examples of use of the SEISCOPE optimization toolbox are
3 provided. First, the bi-dimensional Rosenbrock objective function is considered. This
4 function serves as a benchmark in reference textbooks to calibrate optimization meth-
5 ods. Its particularity is to present a narrow valley of attraction where the function has
6 a very flat minimum. It is therefore a challenging problem for nonlinear optimization.
7 The interest of this example is to illustrate the properties of the different optimization
8 methods proposed in the toolbox. This example can be reproduced using the test case
9 implemented in the toolbox (see file `00README_for_GEOPHYSICS_submission`). The second
10 example shows an application of 2D acoustic frequency-domain FWI on the Marmousi
11 2 model. The Marmousi 2 model is a standard benchmark for today FWI algorithms
12 (Martin et al., 2006). The interest of this example is to illustrate the capability of the
13 SEISCOPE optimization toolbox to handle realistic scale optimization problems and to
14 perform comparisons between the different algorithms which are proposed. In this case
15 the number of discrete unknowns reaches almost 100 000. This example is not directly
16 proposed in the toolbox code attached to this submission. However, it could be repro-
17 duced using the open-source code from the SEISCOPE team TOY2DAC (dowloadable here
18 <http://seiscope2.osug.fr/TOY2DAC,82?lang=en>, accessed on 09/23/2015). The Marmousi
19 2 model can be downloaded here <http://www.agl.uh.edu/downloads/downloads.html>, ac-
20 cessed on 09/23/2015

1 The Rosenbrock function

The analytic expression for the bi-dimensional Rosenbrock function is

$$f(x, y) = (1 - x)^2 + 100(y - x^2)^2. \quad (20)$$

The gradient and the Hessian of the Rosenbrock function are, respectively,

$$\nabla f(x, y) = \begin{pmatrix} 2(x - 1) - 400x(y - x^2) \\ 200(y - x^2) \end{pmatrix}, \quad (21)$$

$$H(x, y) = \begin{pmatrix} 1200x^2 - 400y + 2 & -400x \\ -400x & 200 \end{pmatrix}. \quad (22)$$

2 The Rosenbrock function reaches its minimum value at the point $(1, 1)$. In Figure 4 the
3 valley of attraction appears in blue. The steepness of the objective function in this valley
4 is weak, which renders the convergence to the minimum $(1, 1)$ difficult for local descent
5 algorithms.

6 As an illustration of the performance of the algorithms proposed in the toolbox, the
7 different paths taken by steepest-descent, nonlinear conjugate gradient, l -BFGS and trun-
8 cated Newton methods are presented in Figure 5. The starting point (initial guess) is taken
9 as $(x_0, y_0) = (0.25, 0.25)$.

10 The memory parameter l , which controls the number of gradient stored to build the l -
11 BFGS approximation of the inverse Hessian, is set to 20. Changes of this value do not show
12 significant modification of the l -BFGS path. For realistic size case studies, 20 is already a
13 “large” value for this parameter. For strongly nonlinear misfit functions (for instance in the
14 presence of noisy data when the method is applied to nonlinear least-squares minimization),
15 it may be useful to take lower values such that $l = 3$ or $l = 5$, as the information carried

1 out by older iterates may be not relevant to the local approximation of the inverse Hessian
2 operator.

3 Starting from (x_0, y_0) , all the methods head rapidly toward the valley of attraction.
4 The path they take when they reach this area is however very different. The steepest-
5 descent algorithm performs the worst. It generates very small steps and therefore numerous
6 accumulation points appear on the path it takes. The convergence is reached only after
7 several thousand iterations. The nonlinear conjugate gradient performs better, but presents
8 a quite irregular convergence toward the minimum. In particular, the algorithm goes beyond
9 the minimum before heading back to it. The l -BFGS algorithm path is more regular,
10 however a zone of accumulation points is created near the minimum. Finally, the truncated
11 Newton method proposes the most efficient path, taking long updates in the attraction
12 valley before reaching the minimum around which smaller steps are taken.

13 This example is completed by the analysis of the convergence curves of the four algo-
14 rithms in terms of number of iterations and computational effort in Figure 6. Not surpris-
15 ingly, the steepest-descent algorithm convergence is very slow compared to the three other
16 methods. Although this does not appear in Figure 6, the convergence is reached only after
17 several thousand iterations. On the other hand, the three other methods converge in less
18 than 60 iterations. Among these three, the methods follow the hierarchy which could be
19 expected: the nonlinear conjugate gradient is the slowest (53 iterations), followed by the
20 l -BFGS method (29 iterations); the most efficient one is the truncated Newton method (18
21 iterations).

22 However, this observation has to be mitigated by the analysis of the behavior of the
23 methods in terms of computational effort. The latter is measured by the number of gra-

1 dient computation required to reach convergence. In the case of the truncated-Newton
2 method, the cost of a Hessian-vector multiplication is equivalent to the cost of a gradi-
3 ent computation. In terms of computational effort, the hierarchy between l -BFGS and the
4 truncated Newton method is inverted. This is understandable as even if the truncated New-
5 ton method converges in less iterations, each iteration of this algorithm has a significantly
6 higher cost related to the additional Hessian-vector products it has to perform.

7 The Rosenbrock case study is a good illustration of what can be expected from the dif-
8 ferent descent methods implemented in the SEISCOPE optimization toolbox. For difficult
9 problems, the information on the local curvature of the objective function helps to guarantee
10 a faster convergence speed. This information is embedded in the second-order derivatives
11 of the objective function. The two methods accounting at best for this information are the
12 l -BFGS method and the truncated Newton method. The latter has the capability to inte-
13 grate the information on the Hessian with a great accuracy, which decreases the number of
14 iteration required to reach the minimum, at the expense of an increased computational com-
15 plexity. The l -BFGS method proposes the better trade-off between computational efficiency
16 and accuracy of the Hessian estimation.

17 **FWI on the Marmousi case study**

18 *Full Waveform Inversion principle*

FWI is a seismic imaging method based on the minimization between observed data and synthetic data computed through the numerical solution of a wave propagation problem. The method allows us to retrieve high resolution quantitative estimates of subsurface parameters such as P-wave and S-wave velocities, density, or anisotropy parameters. For this

simple illustration, a 2D acoustic frequency-domain application with constant density is used. In this context, the wave propagation is described by the Helmholtz equation

$$-\frac{\omega^2}{v_P^2}u - \Delta u = s, \quad (23)$$

- ¹ where ω denotes the circular frequency, $v_P(x)$ the P-wave velocity, $u(x, \omega)$ the pressure
- ² wavefield, Δ the Laplacian operator, and $s(x, \omega)$ an explosive source term.

Given partial measurement of the pressure wavefield d_{obs} acquired at the surface, the FWI problem consists in the nonlinear least-squares minimization problem

$$\min_{v_P} \frac{1}{2} \|d_{cal}(v_P) - d_{obs}\|^2, \quad (24)$$

where $\|\cdot\|$ is the L^2 norm and $d_{cal}(v_P)$ is computed as

$$d_{cal}(v_P) = R u(x, \omega; v_P). \quad (25)$$

- ³ In equation (25), $u(x, \omega; v_P)$ is the solution of the Helmholtz equation (23) for the P-wave
- ⁴ velocity model v_P , and R is a restriction operator mapping this wavefield to the receiver
- ⁵ location where the measurements are performed.

⁶ In the following numerical examples, the **TOY2DAC** code from the **SEISCOPE** group
⁷ is used. This code is based on the solution of the equation (23) through a fourth-order
⁸ mixed staggered-grid finite differences discretization (Hustedt et al., 2004). The solution of
⁹ the Helmholtz equation in an infinite domain is simulated with Perfectly Matched Layers
¹⁰ (Bérenger, 1994). The associated linear system is solved by a LU factorization performed
¹¹ using the massively parallel solver MUMPS (Amestoy et al., 2000; MUMPS-team, 2011).

¹² The computation of the gradient and Hessian-vector products is performed following first-
¹³ order and second-order adjoint methods. Using these techniques the computation of these
¹⁴ quantities is reduced to the solution of wave propagation problems of type (23) with different

1 source terms (Plessix, 2006; Métivier et al., 2012, 2013). The TOY2DAC code is interfaced
2 with the SEISCOPE optimization toolbox.

3 *Comparison of the optimization routines on the Marmousi2 benchmark model*

4 The Marmousi 2 benchmark model and the initial model which is used are presented in
5 Figure 7. The initial model is obtained after applying a Gaussian smoothing to the exact
6 model with a 500 m correlation length. A fixed-spread acquisition system with 336 sources
7 and receivers positioned at 50 m depth, from $x = 0.15$ km to $x = 16.9$ km, equally spaced
8 each 50 m, is used. The synthetic observed data-set is constructed using six frequencies
9 from 3 Hz to 8 Hz, with 1 Hz sampling. No free-surface condition is implemented. The
10 synthetic observed data thus do not contain any surface multiples. No multi-scale strategy
11 is used: the data corresponding to the six frequencies is inverted simultaneously. After
12 discretization, the P-wave velocity model is described by approximately 100 000 discrete
13 parameters, which yields a large-scale nonlinear optimization problem.

14 For FWI, the computational efficiency of the different optimization methods can be
15 measured in terms of the number of gradient required to reach a certain level of accuracy.
16 Indeed, a complexity analysis of the method reveals that the average complexity of one
17 gradient computation is in $O(N^4)$ for 2D experiments ($O(N^6)$ for 3D experiments), where
18 N is the average number of discrete points in one direction. This approximation is ob-
19 tained assuming that the number of sources is in $O(N)$ (resp. $O(N^2)$ for the 3D case). In
20 comparison, the complexity if the operations performed within the different optimization
21 methods is in $O(N^2)$ (resp. $O(N^3)$ for the 3D case). Therefore the computation associated
22 with the optimization algorithm in itself is negligible compared to the computation cost

1 associated with gradient computation. It should be mentioned here that in the context
2 of adjoint strategies for computing Hessian-vector product, the computation cost of this
3 operation is the same as the one for the gradient, therefore this simplifies the comparison
4 between first-order method and l -BFGS with the truncated Newton method.

5 The results obtained using the four different optimization methods available in the
6 toolbox, steepest-descent, nonlinear conjugate gradient, l -BFGS, truncated Newton, are
7 presented in Figure 8. All the methods are combined with a diagonal preconditioner. To
8 compute this preconditioner, the diagonal elements of the Hessian operator are approxi-
9 mated through the pseudo-Hessian approach promoted by Shin et al. (2001). The precon-
10 ditioner is computed as the inverse of the diagonal matrix formed with these elements. A
11 similar comparison of optimization methods in the context of FWI has been performed by
12 Castellanos et al. (2015) on different models.

13 The number of gradients stored for computing the l -BFGS approximation is set to $l = 10$.
14 For the truncated Newton method, the maximum number of iterations for the incomplete
15 solution of the Newton equations is set to 10. In practice, this bound is not reached: the
16 stopping criterion from Eisenstat and Walker (1994) is satisfied at each nonlinear iteration
17 in less than 10 iterations. The purpose of this experiment is to compare the different
18 convergence rates of the methods, therefore the stopping criterion which is used is based
19 on the maximum number of nonlinear iterations to be performed, which is set to 20. The
20 convergence curves in function of the number of iterations and the number of gradient
21 estimations are presented in Figure 9.

22 As can be observed in Figure 9, the nonlinear conjugate gradient does not yield any
23 improvement with respect to the steepest-descent method. This is not really a surprise,

1 given the known erratic behavior of this method. In this situation, the method seems less
2 efficient than the steepest-descent algorithm.

3 A second observation is the relatively good performance of the steepest-descent com-
4 pared to the very slow convergence observed for the Rosenbrock test case for this method.

5 Two reasons may be invoked to explain this observation. First, a preconditioner is used
6 in the Marmousi 2 case, which improves significantly the convergence rate of this method.
7 Second, the Rosenbrock case study is a pathological situation in which the steepest-descent
8 has the highest difficulty to converge. For the large-scale Marmousi 2 application, this seems
9 to be no longer true.

10 As in the Rosenbrock case study, second-order methods outperform first-order methods.

11 The truncated Newton method provides the fastest convergence rate in terms of nonlinear
12 iterations. In terms of computational efficiency, the l -BFGS method provides the best
13 performance.

14 The analysis of the four P-wave estimations obtained after 20 iterations of each of the
15 algorithm is in agreement with the objective function level attained by the four methods
16 (Fig. 8). The steepest-descent and nonlinear conjugate gradient estimations are compa-
17 rable. The shallow structures are well recovered. The low velocity anomalies accounting
18 for the presence of gas located at ($z = 2$ km, $x = 6$ km) and ($z = 2$ km, $x = 21$ km) are
19 efficiently reconstructed. However, the deeper structures below $z = 4.5$ km are not well de-
20 lineated. Conversely, the l -BFGS and truncated Newton methods are able to better refocus
21 the energy on this deeper part of the model. The deep reflectors are correctly reconstructed.
22 This is particularly visible in the truncated Newton reconstruction. The inverse Hessian
23 operator acts as an efficient deblurring filter in the context of FWI (Pratt et al., 1998).

1 Although the l -BFGS algorithm is more efficient than the truncated Newton method in
2 terms of computational cost for this mono-parameter FWI case study, the situation may
3 change in the context of multi-parameter FWI. In this context, non-negligible trade-offs
4 between different classes of parameters are expected (Operto et al., 2013). Removing these
5 trade-offs require to account accurately for the inverse Hessian operator at the first stages
6 of the inversion (Métivier et al., 2014). As the l -BFGS algorithm builds a progressive
7 estimation of this operator along the iterations of the minimization loop, the estimation
8 in the first iterations depends strongly on the prior information injected by the user. If
9 this information is poor, as it is often the case, it is expected that the trade-offs between
10 parameter classes contaminate the solution at the early stages of the inversion. Removing
11 these trade-offs in the next iterations is an extremely difficult task (see for instance the
12 multi-parameter FWI for ground penetrating radar data performed by Lavoué et al. (2014)).
13 In contrast, the truncated Newton method is based on an approximation of the inverse
14 Hessian operator whose accuracy does not depend on the convergence history. Therefore, it
15 is expected that the truncated Newton method performs better than the l -BFGS algorithm
16 in a multi-parameter FWI context.

17 The comparison between the different minimization algorithms on the Marmousi case
18 study should be ended up with the following comments. The settings of the experiments
19 have been designed such that the initial model is close from the exact one. This could favor
20 second-order algorithm over first-order algorithm because this is the configuration in which
21 the superlinear convergence of second-order algorithms is more likely to be observed (in the
22 vicinity of the minimum). In real case applications, the initial model can be poorer and the
23 difference between first-order and second-order algorithm may be less obvious, at least in
24 the early stages of the inversion. In any case, if the initial solution is very poor, all the four

1 minimization strategies will fail to produce a reliable subsurface model estimation. What
2 can be observed is an acceleration of the convergence rate when the solution approaches the
3 minimum of the misfit function. The use of Tikhonov regularization also has an impact on
4 the differences between first-order and second-order methods. This regularization technique,
5 mandatory in case of noisy data for instance, results in adding a multiple of the identity
6 to the Hessian operator. To see the impact of this regarding optimization method, one can
7 consider the limit case, for which the Hessian operator tends toward the identity operator. In
8 this case, first-order methods, relying on this approximation, become equivalent to second-
9 order methods. Therefore, using strong Tikhonov regularization weights will have the effect
10 of reducing the differences between first-order and second-order methods regarding their
11 convergence rate.

CONCLUSION

12 The SEISCOPE optimization toolbox is a set of **FORTRAN 90** routines for the solution of
13 large-scale nonlinear minimization problems. This type of problems is ubiquitous in geo-
14 physics. The toolbox implements four different optimization schemes: the steepest-descent,
15 the nonlinear conjugate gradient, the l -BFGS method and the truncated Newton methods.
16 All these routines are implemented with a linesearch strategy ensuring the robustness of the
17 methods: monotonic decrease of the objective function and guaranteed convergence toward
18 the nearest local minimum. Bound constraints can also be activated.

19 The use of these routines within a computational code is completely uncoupled from the
20 physical problem through a reverse communication protocol. Within this framework, the
21 minimization of the objective function is brought back to the definition of a minimization
22 loop controlled by the user in which the solver is called at each iteration. Depending on

1 the return values of a communication flag, the user performs the action required by the
2 solver at a given point, namely: computation of the objective function and its gradient,
3 application of a preconditioner, computation of a Hessian-vector product. This allows for
4 a flexible implementation and an easier use of second-order optimization methods.

5 The general structure of the minimization algorithms of the toolbox is based on a two-
6 level reverse communication principle. The first level consists in implementing the commu-
7 nication between the user code and the minimization code. The second level is embedded
8 in the minimization code itself, and implements the communication between a linesearch
9 algorithm and the minimization code. The minimization code transfers the requirement
10 from the linesearch (computation of the objective function and its gradient) directly to the
11 user. This imbrication of reverse communication levels is thus transparent for the user.

12 The example of use of the toolbox provided in the numerical results part show the
13 benefit one can expect from second-order methods. The case of the Rosenbrock function is
14 pathological, however it illustrates how powerful can be the introduction of the information
15 carried out by the second-order derivatives into the optimization. The Marmousi 2 case
16 for FWI also illustrates the benefit of introducing an appropriate estimation of the inverse
17 Hessian operator within the minimization to speed-up the convergence of the algorithm and
18 save considerable computation time.

19 Perspectives of development of the toolbox to adapt it to very-large scale HPC applica-
20 tions is currently considered. In its actual implementation, the method requires to collect
21 one model and one gradient on a single node of a supercomputer. For very large-scale ap-
22 plications, such as the one involved in 3D FWI, these quantities are usually distributed and
23 this should require additional global communications at each iteration of the minimization

1 loop. In some cases, the ability to store these quantities on a single node may even be
2 questionable. Therefore, a completely distributed version of the toolbox may be derived
3 to relax this requirement. One possible implementation is to externalize all the operations
4 related to scalar products and vector additions involved in the minimization procedures
5 through the reverse communication protocol. This could slightly complicate the interface
6 of the toolbox with other computer codes, requiring more actions to be performed within
7 the minimization loop. However, it would represent a very efficient way of performing the
8 minimization on very large-scale problems for which all the vector quantities have to be
9 distributed over the whole cluster.

10 On a longer term, the set of routines currently included in the toolbox could be extended
11 to routines dedicated to the solution of constrained optimization problems. Based on the
12 same architecture, sequential quadratic programming solvers could be designed, for han-
13 dling problems with not only bound constraints, but also linear and nonlinear equality and
14 inequality constraints. This type of problems arise in geophysics as soon as regularization
15 methods are considered, beyond simple additive techniques. Another longer term extension
16 consists in integrating trust-region algorithms to propose an alternative to the linesearch
17 technique currently implemented. In particular, the combination of the truncated Newton
18 method and the trust-region method, known as the Steihaug algorithm (Steihaug, 1983),
19 is reputed to benefit from better convergence properties than the linesearch version of the
20 truncated Newton method. These extensions may be the topic of future investigations.

ACKNOWLEDGEMENTS

21 This study was partially funded by the SEISCOPE consortium (<http://seiscope2.osug.fr>),
22 sponsored by BP, CGG, CHEVRON, EXXON-MOBIL, JGI, PETROBRAS, SAUDI ARAMCO,

1 SCHLUMBERGER, SHELL, SINOPEC, STATOIL, TOTAL and WOODSIDE. This study
 2 was granted access to the HPC resources of the Froggy platform of the CIMENT infras-
 3 tructure (<https://ciment.ujf-grenoble.fr>), which is supported by the Rhône-Alpes region
 4 (GRANT CPER07_13 CIRA), the OSUG@2020 labex (reference ANR10 LABX56) and the
 5 Equip@Meso project (reference ANR-10-EQPX-29-01) of the programme Investissements
 6 d’Avenir supervised by the Agence Nationale pour la Recherche, and the HPC resources
 7 of CINES/IDRIS under the allocation 046091 made by GENCI. The authors would like to
 8 thank all the actors of the SEISCOPE consortium research group for their active and con-
 9 stant support. Special thanks are due to Stéphane Operto for suggesting us this article and
 10 his (always) careful proof-reading. Special thanks also go to Joe Dellinger, Philippe Thierry,
 11 Esteban Diaz and two anonymous reviewers for their suggestions, which really helped us to
 12 improve our work.

APPENDIX

The l -BFGS algorithm is based on the computation of the descent direction Δx_k following the equation (8). In (8), Q_k is defined by

$$\begin{aligned}
 Q_k = & (V_{k-1}^T \dots V_{k-l}^T) Q_k^0 (V_{k-l} \dots V_{k-1}) \\
 & + \rho_{k-l} (V_{k-1}^T \dots V_{k-l+1}^T) s_{k-l} s_{k-l}^T (V_{k-l+1} \dots V_{k-1}) \\
 & + \rho_{k-l+1} (V_{k-1}^T \dots V_{k-l+2}^T) s_{k-l+1} s_{k-l+1}^T (V_{k-l+2} \dots V_{k-1}) \\
 & + \dots \\
 & + \rho_{k-1} s_{k-1} s_{k-1}^T,
 \end{aligned} \tag{26}$$

where the pairs s_k, y_k are

$$s_k = x_{k+1} - x_k, \quad y_k = \nabla f(x_{k+1}) - \nabla f(x_k), \tag{27}$$

the scalar ρ_k are

$$\rho_k = \frac{1}{y_k^T s_k}, \quad (28)$$

and the matrices V_k are defined by

$$V_k = I - \rho_k y_k s_k^T. \quad (29)$$

- ¹ These formulas are directly extracted from Nocedal and Wright (2006). As is also explained
- ² in this reference textbook, the computation of the quantity $\Delta x_k = -Q_k \nabla f(x_k)$, can be
- ³ performed through the double recursion algorithm presented below.

Data: $\rho_i, s_i, y_i, i = k-l, \dots, k-1, H_k^0, \nabla f(x_k)$

Result: $\Delta x_k = -Q_k \nabla f(x_k)$

$q = -\nabla f(x_k);$

for $i = k-l, \dots, k-1$ **do**

$\alpha_i = \rho_i s_i^T \Delta x_k;$

$q = q - \alpha_i y_i;$

end

$\Delta x_k = H_k^0 q;$

for $i = k-l, \dots, k-1$ **do**

$\beta = \rho_i y_i^T \Delta x_k;$

$\Delta x_k = \Delta x_k + (\alpha_i - \beta_i) s_i;$

end

Algorithm 1: Two-loop recursion algorithm

- 1 In Algorithm 1, the matrix H_k^0 is an estimation of the inverse Hessian matrix at the k^{th}
- 2 iteration of the l -BFGS minimization. In most applications, this estimation is not updated
- 3 throughout the l -BFGS iterations, and kept equal to H_0^0 . However, nothing prevents from
- 4 performing this update. In practice, in the FWI application presented in this study, this
- 5 strategy reveals to be efficient. The inverse Hessian estimation which is used is a diagonal
- 6 approximation based on the pseudo-Hessian strategy (Shin et al., 2001). The computation
- 7 cost of this approximation is cheap, and therefore can be recalculated at each iteration. The
- 8 implementation which is used in the toolbox allows for these updates to be perform. When
- 9 the routine PLBFGS is used, the preconditioning operation requested by the communicator
- 10 flag set to ‘‘PREC’’ corresponds to the multiplication $\Delta x_k = H_k^0 q$. The user may thus
- 11 update the matrix H_k^0 at each iteration k before performing this multiplication. Surprisingly,
- 12 this possibility is rarely applied (to the best of our knowledge) while it is explicitly specified

¹ in Nocedal and Wright (2006).

REFERENCES

- 1 Amestoy, P., I. S. Duff, and J. Y. L'Excellent, 2000, Multifrontal parallel distributed symmetric and
2 unsymmetric solvers: Computer Methods in Applied Mechanics and Engineering, **184**, 501–520.
- 3 Baumstein, A., 2013, POCS-based geophysical constraints in multi-parameter Full Wavefield In-
4 version: Presented at the 75th EAGE Conference & Exhibition incorporating SPE EUROPEC
5 2013.
- 6 Bérenger, J.-P., 1994, A perfectly matched layer for absorption of electromagnetic waves: Journal
7 of Computational Physics, **114**, 185–200.
- 8 Bonnans, J. F., J. C. Gilbert, C. Lemaréchal, and C. A. Sagastizábal, 2006, Numerical optimization,
9 theoretical and practical aspects: Springer series, Universitext.
- 10 Byrd, R. H., P. Lu, and J. Nocedal, 1995, A limited memory algorithm for bound constrained
11 optimization: SIAM Journal on Scientific and Statistical Computing, **16**, 1190–1208.
- 12 Castellanos, C., L. Métivier, S. Operto, R. Brossier, and J. Virieux, 2015, Fast full waveform inversion
13 with source encoding and second-order optimization methods: Geophysical Journal International,
14 **200**(2), 720–744.
- 15 Dai, Y. and Y. Yuan, 1999, A nonlinear conjugate gradient method with a strong global convergence
16 property: SIAM Journal on Optimization, **10**, 177–182.
- 17 Dongarra, J., V. Eijkhout, and A. Kalhan, 1995, Reverse communication interface for linear algebra
18 templates for iterative methods: Technical report, University of Tennessee.
- 19 Eisenstat, S. C. and H. F. Walker, 1994, Choosing the forcing terms in an inexact Newton method:
20 SIAM Journal on Scientific Computing, **17**, 16–32.
- 21 Fletcher, R. and C. M. Reeves, 1964, Function minimization by conjugate gradient: Computer
22 Journal, **7**, 149–154.
- 23 Gould, N. I. M., S. Lucidi, M. Roma, and P. Toint, 1999, Solving the trust-region subproblem using
24 the lanczos method: Siam Journal on Optimization, **9**, 504–525.
- 25 Hestenes, M. R. and E. Stiefel, 1952, Methods of conjugate gradient for solving linear systems:
26 Journal of Research of the National Bureau of Standards, **49**, no. 6, 409–436.

-
- 1 Hustedt, B., S. Operto, and J. Virieux, 2004, Mixed-grid and staggered-grid finite difference methods
2 for frequency domain acoustic wave modelling: Geophysical Journal International, **157**, 1269–
3 1296.
- 4 Lambaré, G., J. Virieux, R. Madariaga, and S. Jin, 1992, Iterative asymptotic inversion in the
5 acoustic approximation: Geophysics, **57**, 1138–1154.
- 6 Lavoué, F., R. Brossier, L. Métivier, S. Garambois, and J. Virieux, 2014, Two-dimensional permit-
7 tivity and conductivity imaging by full waveform inversion of multioffset GPR data: a frequency-
8 domain quasi-Newton approach: Geophysical Journal International, **197**, 248–268.
- 9 Martin, G. S., R. Wiley, and K. J. Marfurt, 2006, Marmousi2: An elastic upgrade for Marmousi:
10 The Leading Edge, **25**, 156–166.
- 11 Métivier, L., F. Breteau, R. Brossier, S. Operto, and J. Virieux, 2014, Full waveform inver-
12 sion and the truncated Newton method: quantitative imaging of complex subsurface structures:
13 Geophysical Prospecting, **62**, no. 6, 1353–1375.
- 14 Métivier, L., R. Brossier, J. Virieux, and S. Operto, 2012, Toward Gauss-Newton and exact Newton
15 optimization for full waveform inversion: Presented at the 74th Annual International Meeting,
16 EAGE.
- 17 ———, 2013, Full Waveform Inversion and the truncated Newton method: SIAM Journal On Scien-
18 tific Computing, **35(2)**, B401–B437.
- 19 MUMPS-team, 2011, MUMPS - Multifrontal Massively Parallel Solver users' guide - version
20 4.10.0 (may 10, 2011). ENSEEIHT-ENS Lyon, <http://www.enseeiht.fr/apo/MUMPS/> or
21 <http://graal.ens-lyon.fr/MUMPS>.
- 22 Nash, S. and J. Nocedal, 1991, A Numerical Study of the Limited Memory BFGS Method and
23 the Truncated-Newton Method for Large Scale Optimization: SIAM Journal on Optimization, **1**,
24 358–372.
- 25 Nash, S. G., 2000, A survey of truncated Newton methods: Journal of Computational and Applied
26 Mathematics, **124**, 45–59.
- 27 Nemeth, T., C. Wu, and G. T. Schuster, 1999, Least-squares migration of incomplete reflection data:

-
- 1 Geophysics, **64**, 208–221.
- 2 Nocedal, J., 1980, Updating Quasi-Newton Matrices With Limited Storage: Mathematics of Com-
- 3 putation, **35**, 773–782.
- 4 Nocedal, J. and S. J. Wright, 2006, Numerical optimization: Springer, 2nd edition.
- 5 Nolet, G., 1987, Seismic tomography with applications in global seismology and exploration geo-
- 6 physics: D. Reidel publishing Company.
- 7 Operto, S., R. Brossier, Y. Gholami, L. Métivier, V. Prieux, A. Ribodetti, and J. Virieux, 2013, A
- 8 guided tour of multiparameter full waveform inversion for multicomponent data: from theory to
- 9 practice: The Leading Edge, **Special section Full Waveform Inversion**, 1040–1054.
- 10 Plessix, R. E., 2006, A review of the adjoint-state method for computing the gradient of a functional
- 11 with geophysical applications: Geophysical Journal International, **167**, 495–503.
- 12 Polak, E. and G. Ribi  re, 1969, Note sur la convergence de m  thodes de directions conjugu  es:
- 13 Revue Fran  aise d’Informatique et de Recherche Op  rationnelle, **16**, 35–43.
- 14 Pratt, R. G., C. Shin, and G. J. Hicks, 1998, Gauss-Newton and full Newton methods in frequency-
- 15 space seismic waveform inversion: Geophysical Journal International, **133**, 341–362.
- 16 Sen, M. K. and P. L. Stoffa, 1995, Global optimization methods in geophysical inversion: Elsevier
- 17 Science Publishing Co.
- 18 Shin, C., S. Jang, and D. J. Min, 2001, Improved amplitude preservation for prestack depth migration
- 19 by inverse scattering theory: Geophysical Prospecting, **49**, 592–606.
- 20 Spall, J. C., 2003, Introduction to Stochastic Search and Optimization: Estimation, simulation and
- 21 control: Wiley-Interscience Series in Discrete Mathematics and Optimization. 1st edition.
- 22 Steihaug, T., 1983, The conjugate gradient method and trust regions in large scale optimization:
- 23 SIAM Journal on Numerical Analysis, **20**, 626–637.
- 24 Virieux, J. and S. Operto, 2009, An overview of full waveform inversion in exploration geophysics:
- 25 Geophysics, **74**, WCC1–WCC26.

TABLES

	First-order methods		Second-order methods			
	PSTD	PNLCG	LBFGS	PLBFGS	TRN	PTRN
Method	Steep. descent	Nonlin. CG	<i>l</i> -BFGS	<i>l</i> -BFGS	trunc. Newt.	trunc. Newt.
Preconditioning	Yes	Yes	No	Yes	No	Yes
Conv. rate	lin.	lin.	suplin.	suplin.	suplin.	suplin.
$\nabla f(m)$ required	Yes	Yes	Yes	Yes	Yes	Yes
$H(m)v$ required	No	No	No	No	Yes	Yes
Rank test 1	4,4	3,3	2,1	N/A	1,2	N/A
Rank test 2	3,3	4,4	N/A	2,1	N/A	1,2

Table 1: Summary of the minimization routines implemented in the toolbox. The ranks in the two last lines correspond to the performance on the numerical examples 1 (Rosenbrock function) and 2 (FWI on the Marmousi mode) in terms of convergence rate (normal font), and number of computed gradient (bold font).

FLAG values	Action required /Meaning
‘‘INIT’’	This flag is set only by the user, prior to the minimization loop. On reception of this flag, the solver performs the necessary initializations (parameter settings, memory allocations).
‘‘CONV’’	The convergence criterion $\frac{f(x_k)}{f(x_0)} < \varepsilon$ is satisfied. The variable x contains the solution x_k computed at this stage.
‘‘FAIL’’	The linesearch has terminated on a failure.
‘‘NSTE’’	The linesearch has terminated successfully. The variable x contains the new iterate x_{k+1} .
‘‘GRAD’’	Compute $f(x)$ and store it in fcost . Compute $\nabla f(x)$ and store it in grad . For PNLCG and PSTD, apply the preconditioner to $\nabla f(x)$ and store it in grad_preco .
‘‘PREC’’	Only for PLBGS and PTRN routines. If the routine called is PLBFGS, then apply the preconditioner to the variable optim%q_plb and store it in the same variable. If the routine called is PTRN, then apply the preconditioner to the variable optim%residual and store it in the variable optim%residual_preco
‘‘HESS’’	Only for TRN and PTRN routines. Apply the Hessian operator to the variable optim%d and store it in optim%Hd

Table 2: Summary of the different values taken by the communication variable **FLAG** and their meaning.

FIGURES CAPTION

- 1 • FIGURE 1: Example of a reverse communication loop
- 2 • FIGURE 2: General structure of the toolbox optimization routines
- 3 • FIGURE 3: Schematic view of the linesearch implementations
- 4 • FIGURE 4: Bi-dimensional Rosenbrock function map. The two white crosses highlight the
5 starting point for the optimization process located at $(0.25, 0.25)$ and the minimum of the
6 Rosenbrock function, at the point $(1, 1)$. The valley of attraction appears as a broad blue
7 channel where the function has a very flat minimum
- 8 • FIGURE 5: Convergence paths taken by the four different optimization methods: steepest
9 descent (a), nonlinear conjugate gradient (b), *l*-BFGS (c), truncated Newton (d). The starting
10 point is $(0.25, 0.25)$. The convergence is reached at $(1, 1)$.
- 11 • FIGURE 6: Rosenbrock function case study. Convergence rate of the four optimization
12 methods in terms of iterations (a) and computational cost in terms of the number of gradient
13 evaluations (b). For the truncated Newton method, the computation cost of a Hessian-vector
14 product is assimilated to the computation cost of a gradient.
- 15 • FIGURE 7: Marmousi 2 synthetic case study. Exact model (a). Initial model obtained after
16 smoothing the exact model with a Gaussian filter (b). The correlation length for the Gaussian
17 filter is set to 500 m. The number of discrete points is equal to 141×681 , which yields a
18 minimization problem involving approximately 96,000 unknowns.
- 19 • FIGURE 8: Marmousi case study. Models reconstructed by the four different optimization
20 methods. Steepest-descent (a), nonlinear conjugate gradient (b), *l*-BFGS (c), truncated New-
21 ton (d).
- 22 • FIGURE 9: Marmousi case study. Convergence curves of the four different optimization
23 methods in terms of iterations (a), and computational cost in terms of the number of gradient
24 evaluations (b). For the truncated Newton method, the computation cost of a Hessian-vector

1 product is assimilated to the computation cost of a gradient. This is consistent with the
2 implementation of second-order adjoint formulas for the computation of these quantities.

FIGURES

```
do while ((FLAG.ne.'CONV').and.(FLAG.ne.'FAIL'))  
    call PLBFGS(n,x,fcost,grad,grad-preco,optim,FLAG)  
    if (FLAG.eq.'GRAD') then  
        !compute objective function and its gradient at point x  
        call Rosenbrock(x,fcost,grad)  
    endif  
    if (FLAG.eq.'PREC') then  
        !apply preconditioning to optim%q_plb  
        y(:)=optim%q_plb(:)  
        call apply-preco(y,Py)  
        optim%q_plb(:)=Py(:)  
    endif  
    if (FLAG.eq.'NSTE') then  
        !a new step has been taken  
        call print-information(x)  
    endif  
enddo
```

Figure 1: Example of a reverse communication loop.

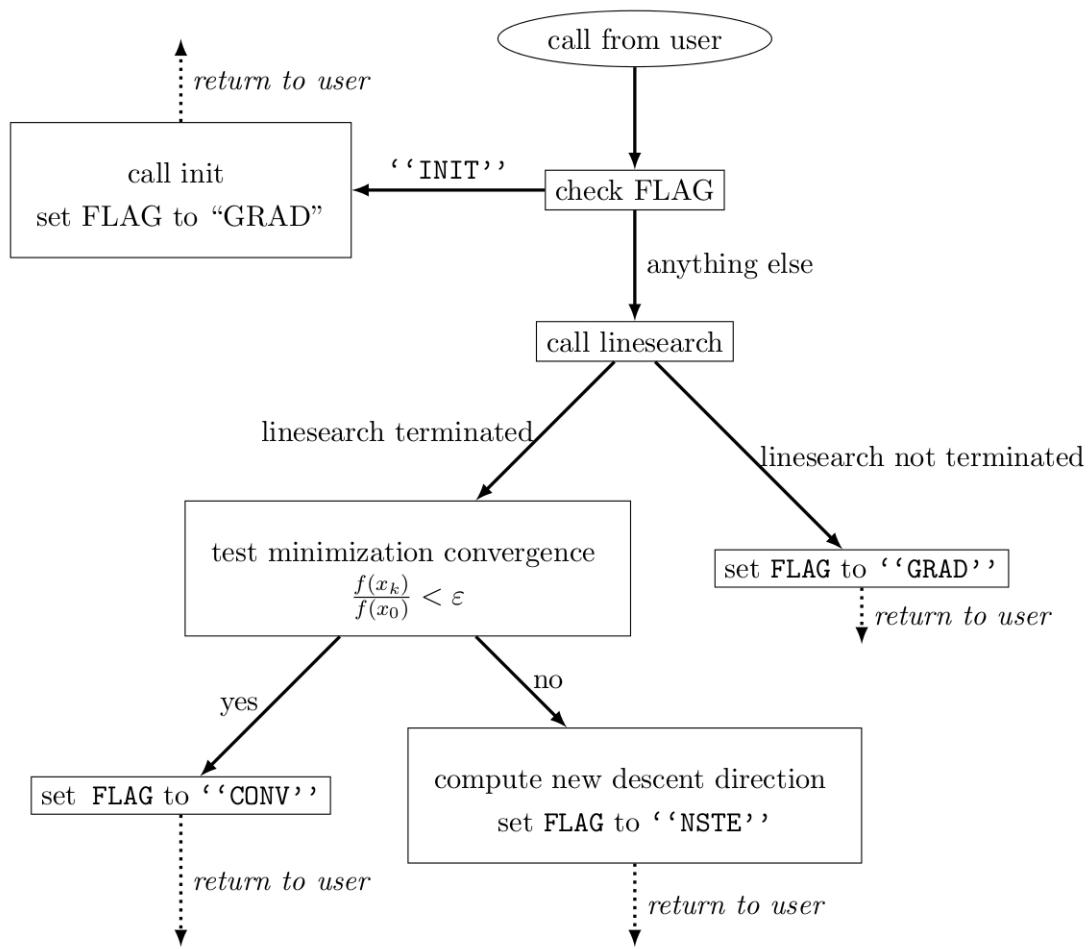


Figure 2: General structure of the toolbox optimization routines.

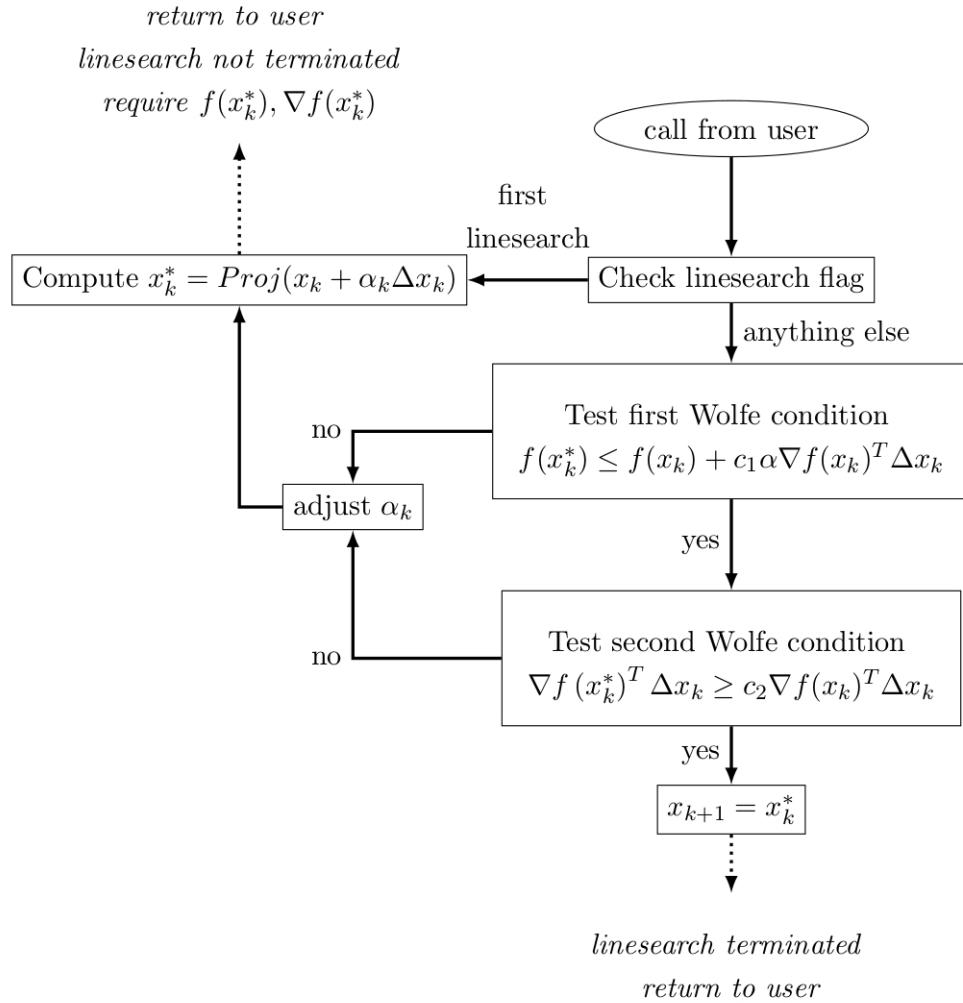


Figure 3: Schematic view of the linesearch implementation.

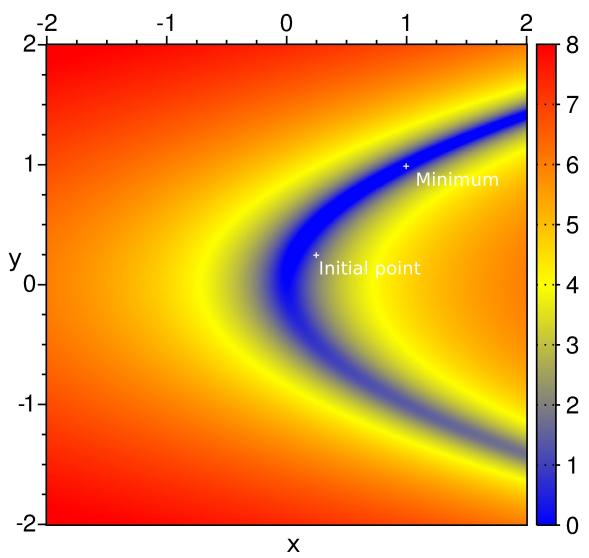


Figure 4: Bi-dimensional Rosenbrock function map. The two white crosses highlight the starting point for the optimization process located at $(0.25, 0.25)$ and the minimum of the Rosenbrock function, at the point $(1, 1)$. The valley of attraction appears as a broad blue channel where the function has a very flat minimum.

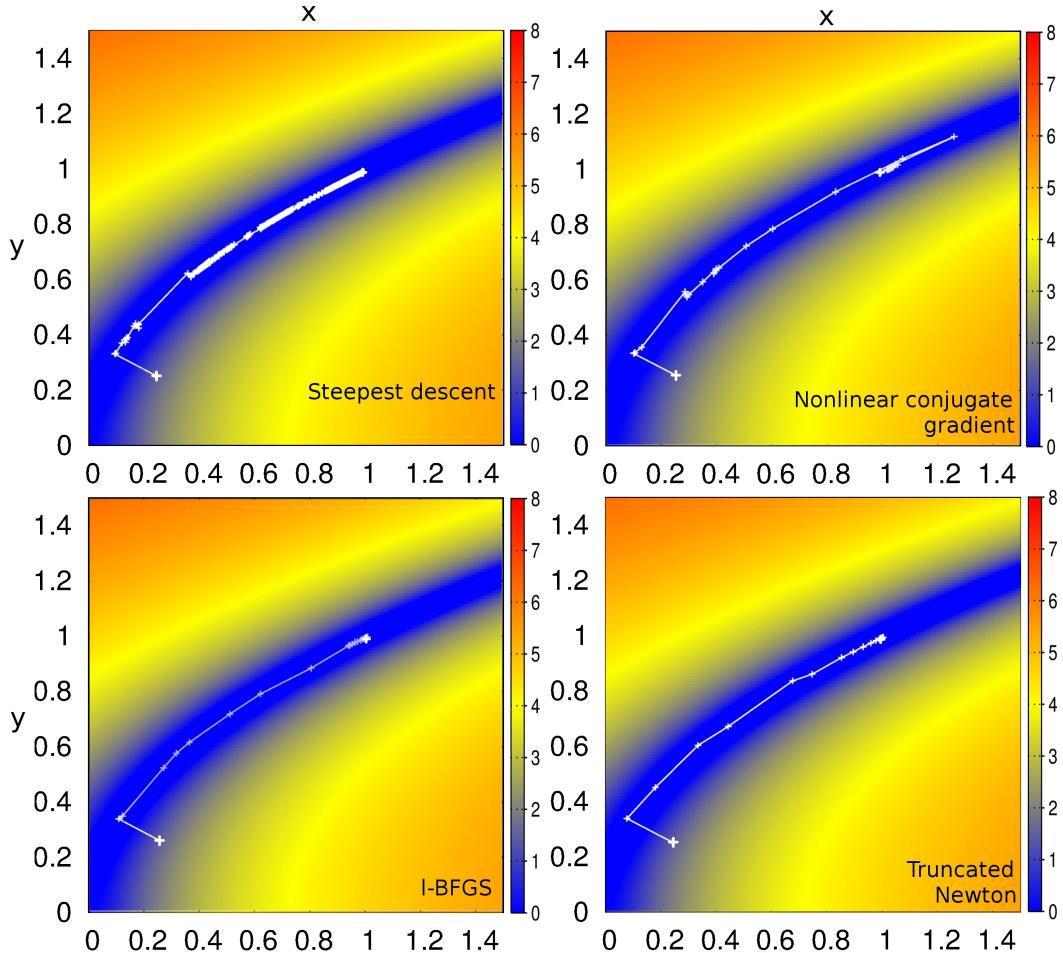


Figure 5: Convergence paths taken by the four different optimization methods: steepest descent (a), nonlinear conjugate gradient (b), l -BFGS (c), truncated Newton (d). The starting point is $(0.25, 0.25)$. The convergence is reached at $(1, 1)$.

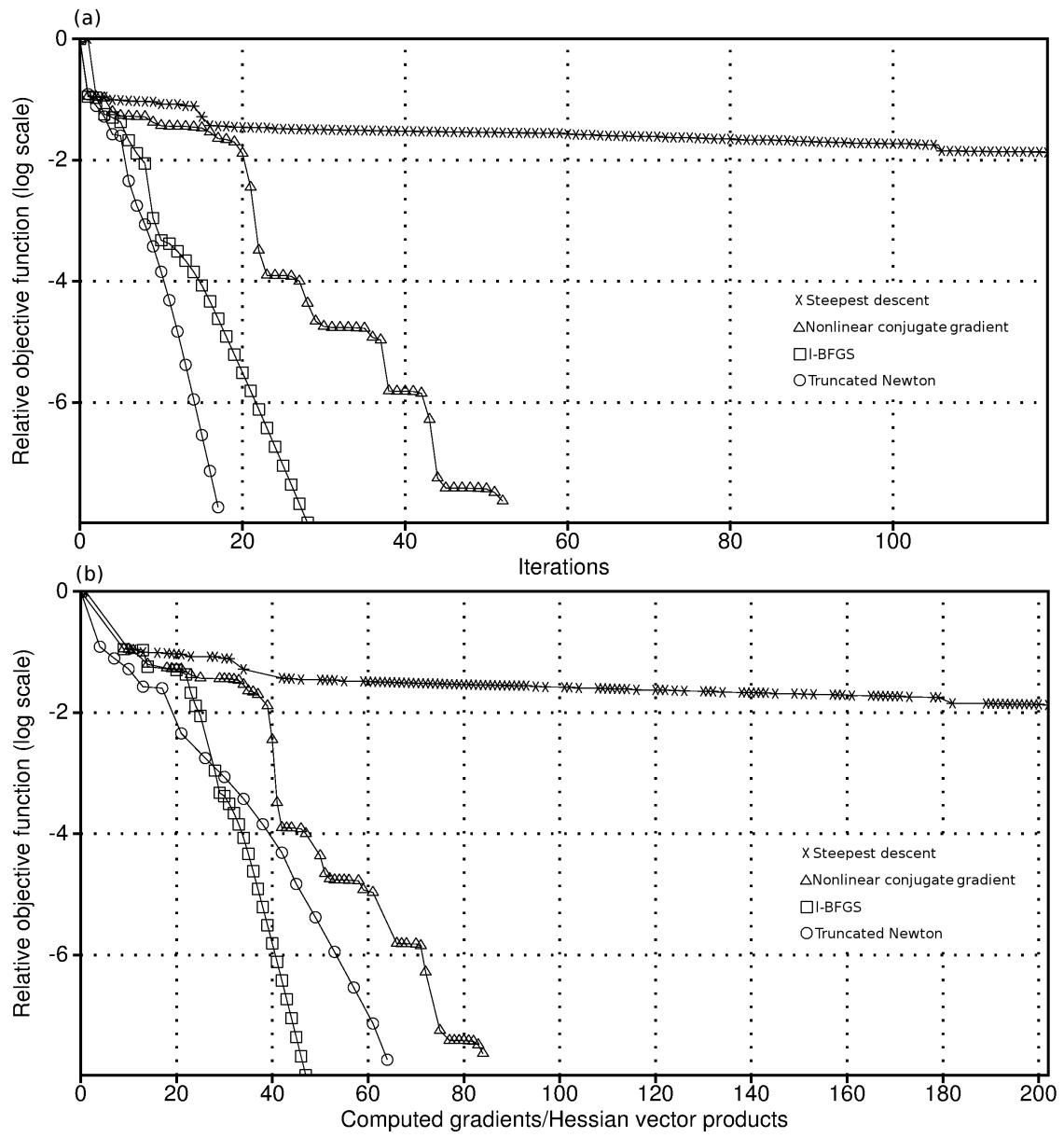


Figure 6: Rosenbrock function case study. Convergence rate of the four optimization methods in terms of iterations (a) and computational cost in terms of the number of gradient evaluations (b). For the truncated Newton method, the computation cost of a Hessian-vector product is assimilated to the computation cost of a gradient.

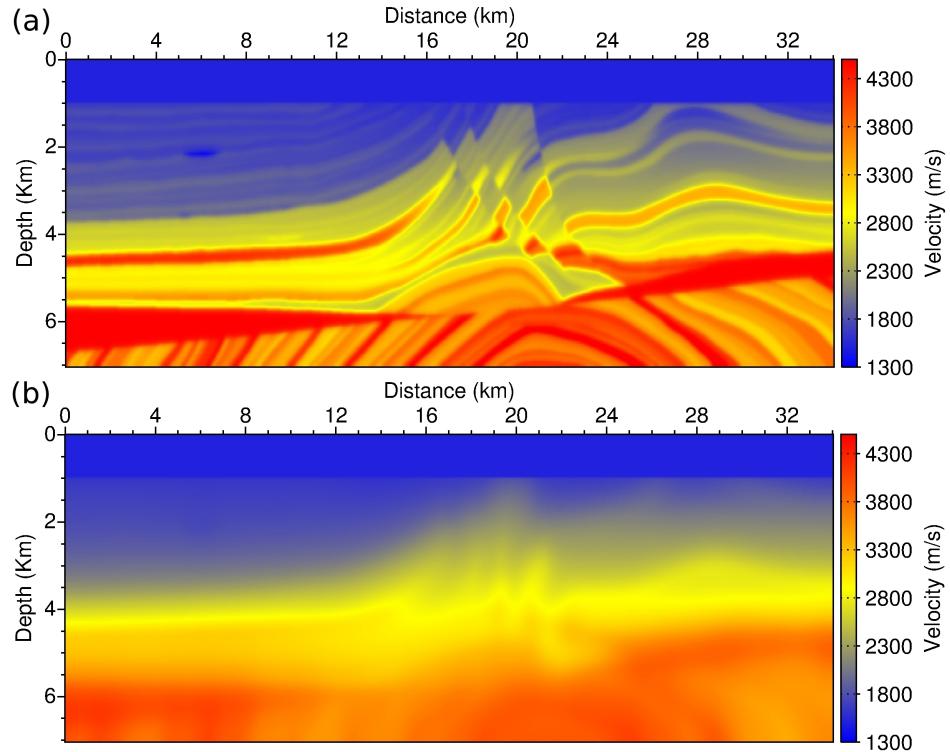


Figure 7: Marmousi 2 synthetic case study. Exact model (a). Initial model obtained after smoothing the exact model with a Gaussian filter (b). The correlation length for the Gaussian filter is set to 500 m. The number of discrete points is equal to 141×681 , which yields a minimization problem involving approximately 96,000 unknowns.

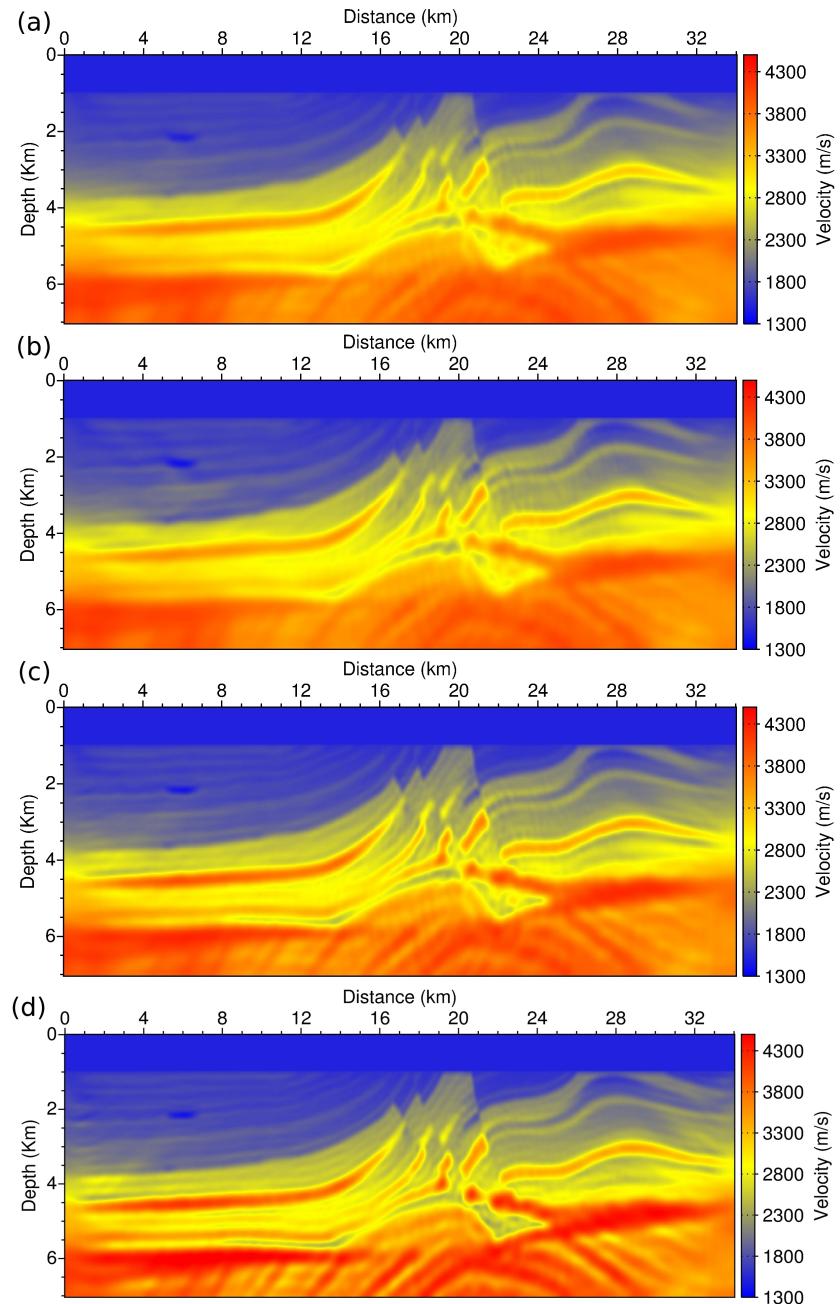


Figure 8: Marmousi case study. Models reconstructed by the four different optimization methods. Steepest-descent (a), nonlinear conjugate gradient (b), l -BFGS (c), truncated Newton (d).

Figure 9: Marmousi case study. Convergence curves of the four different optimization methods in terms of iterations (a), and computational cost in terms of the number of gradient evaluations (b). For the truncated Newton method, the computation cost of a Hessian-vector product is assimilated to the computation cost of a gradient. This is consistent with the implementation of second-order adjoint formulas for the computation of these quantities.