

# MINIMIZATIONS IN THE CYCLE 44 OF ARPEGE/IFS.

YESSAD K. (METEO-FRANCE/CNRM/GMAP/ALGO)

March 6, 2017

## *Abstract:*

*This documentation describes the minimizations done in ARPEGE/IFS, for example in the 3DVAR and 4DVAR assimilation. First some algorithmic aspects are recalled about minimizers and solving of linear algebra systems. There is then a specific description of the minimizers employed in ARPEGE/IFS (M1QN3, N1CG1 and CONGRAD) and of the applications using these minimizers.*

## *Résumé:*

*Cette documentation décrit les minimisations faites dans ARPEGE/IFS, par exemple dans les assimilations 3DVAR et 4DVAR. On rappelle d'abord quelques notions sur les minimiseurs et sur la résolution des systèmes linéaires. On fournit ensuite une description spécifique des minimiseurs utilisés dans ARPEGE/IFS (M1QN3, N1CG1 et CONGRAD) et des applications utilisant ces minimiseurs.*

# Contents

<b>1</b>	<b>Introduction.</b>	<b>3</b>
<b>2</b>	<b>Some algorithms of minimization and of linear systems solving.</b>	<b>4</b>
2.1	Purely direct algorithms. . . . .	4
2.1.1	Gauss method. . . . .	4
2.1.2	Gauss-Jordan method. . . . .	4
2.1.3	<i>LU</i> factorisation method. . . . .	4
2.1.4	Cholesky method. . . . .	4
2.1.5	<i>QR</i> factorisation method. . . . .	4
2.2	Iterative algorithms for small linear systems. . . . .	4
2.2.1	Jacobi method. . . . .	4
2.2.2	Gauss-Seidel method. . . . .	5
2.2.3	Relaxation method. . . . .	5
2.3	Iterative and “mixed” algorithms for large linear systems and minimization problems. . . . .	5
2.3.1	General considerations and denotations. . . . .	5
2.3.2	Preconditioning. . . . .	6
2.3.3	Newton method. . . . .	6
2.3.4	Quasi-Newton method. . . . .	7
2.3.5	Limited memory quasi-Newton method. . . . .	8
2.3.6	Truncated Newton method. . . . .	8
2.3.7	Fletcher-Reeves conjugate gradient method. . . . .	8
2.3.8	Polak-Ribière conjugate gradient method. . . . .	9
2.3.9	Gradient method with optimal step (or steepest descent method). . . . .	9
2.3.10	Conjugate gradient method combined with a Lanczos algorithm. . . . .	9
2.3.11	When to stop the algorithm. . . . .	11
<b>3</b>	<b>Available minimizers in the “arpifs” and “algor” libraries: description and organigrammes.</b>	<b>12</b>
3.1	Minimizer <b>M1QN3</b> developed at INRIA. . . . .	12
3.2	Minimizer <b>M1QN3R</b> . . . . .	13
3.3	Minimizer <b>N1CG1</b> . . . . .	14
3.4	Minimizer <b>CONGRAD</b> . . . . .	15
<b>4</b>	<b>Configurations of use of minimizations, and simulators.</b>	<b>16</b>
4.1	Minimization for orography in climatological files (configuration 923). . . . .	16
4.2	Minimization in the 4DVAR assimilation. . . . .	16
4.3	Other applications using minimizers. . . . .	16
<b>5</b>	<b>Modules and namelists.</b>	<b>17</b>
<b>6</b>	<b>Bibliography.</b>	<b>18</b>
6.1	Publications. . . . .	18
6.2	Some internal notes and other ARPEGE notes. . . . .	18

# 1 Introduction.

ARPEGE/IFS needs to minimize a cost function for example in the 4DVAR assimilation but also to compute a spectral orography in the climatological files (configuration 923). If the cost function is purely quadratic the minimization is equivalent to solve a linear system of equations.

There are several ranges of ways to solve such problems, some apply to solve small linear systems, some other to solve small or big linear systems, some other to minimize functions which are not necessary quadratic ones:

- Direct algorithms to invert linear systems: they can be used to invert rather small systems (the space dimension has a range of several hundreds but not more):
  - Gauss method.
  - Gauss-Jordan method.
  - $LU$  factorisation (square matrices).
  - Cholesky  $LL^T$  factorisation (symmetric positive definite square matrices).
  - $QR$  factorisation (square matrices).
- Pure iterative algorithms to invert linear systems or minimize functions.
  - Jacobi (solves a small linear system).
  - Gauss-Seidel (solves a small linear system).
  - Relaxation (solves a small linear system).
  - Newton and quasi-Newton algorithm (minimization of a function which is not necessary a quadratic one, the dimension of the problem can be large).
  - Relaxation method (minimization, the dimension of the problem can be large).
  - Gradient method (minimization, the dimension of the problem can be large).
- Mixed iterative algorithms to invert linear systems or minimize functions. What we say by “mixed” is that these algorithms are direct algorithms (the maximum number of iterations is equal to the dimension of the problem), but one uses them with a number of iterations significantly lower than the possible maximum, so one uses them as iterative algorithms.
  - Conjugate gradient methods (minimization, the dimension of the problem can be large). There are several versions of this algorithm, some for purely quadratic functions, some other which can be also used for non quadratic functions.
  - Conjugate gradient methods combined with a Lanczos algorithm (minimization, the dimension of the problem can be large).

These algorithms are described in some mathematical books but for each of ones one will recall the principles of them and the range of use in order to have a synthetic overview of them. What is indeed our main concern in this documentation is the sub-list of iterative minimization algorithms which allow to minimize quadratic or non-quadratic functions and to solve big linear systems (the dimension of the problem to solve being in a range around tens or hundreds of thousands).

## 2 Some algorithms of minimization and of linear systems solving.

### 2.1 Purely direct algorithms.

These algorithms generally apply to solve small linear systems. The system which has to be solved writes:

$$\mathbf{A}\mathbf{X} = \mathbf{B} \quad (1)$$

where  $\mathbf{A}$  is a  $N * N$  matrix,  $\mathbf{X}$  is the unknown vector of dimension  $N$  and  $\mathbf{B}$  is the RHS known vector of dimension  $N$ .

#### 2.1.1 Gauss method.

The algorithm involves  $N_I$  iterations which leads to replace the system  $\mathbf{A}\mathbf{X} = \mathbf{B}$  by the system:

$$\mathbf{M}\mathbf{A}\mathbf{X} = \mathbf{M}\mathbf{B} \quad (2)$$

where the product  $[\mathbf{M}\mathbf{A}]$  is an upper triangular  $N * N$  matrix; this new system can be easily solved. This method can be used for any invertible matrix  $\mathbf{A}$ . This method requires  $N^3/3$  additions,  $N^3/3$  multiplications and  $N^2/2$  divisions.

#### 2.1.2 Gauss-Jordan method.

This is a variant of the Gauss method but  $\mathbf{M}$  is now computed in order to have a product  $[\mathbf{M}\mathbf{A}]$  which is diagonal. This method can be used for example when the storage of  $\mathbf{A}^{-1}$  is required.

#### 2.1.3 LU factorisation method.

Matrix  $\mathbf{A}$  is factorized into the matricial product  $\mathbf{L}\mathbf{U}$  where  $\mathbf{L}$  is a lower triangular matrix, the diagonal elements of which are all equal to 1, and  $\mathbf{U}$  is an upper triangular matrix. This is a particular case of Gauss method where  $\mathbf{M} = \mathbf{L}^{-1}$  but one assumes that  $\mathbf{L}$  is stored. This method can be used when the same matrix has to be inverted for a large number of unknown vectors  $\mathbf{X}$ : in this case the matrices  $\mathbf{U}$  and  $\mathbf{L}$  are first pre-computed and stored. This method can be used for any invertible matrix  $\mathbf{A}$  which a minimum amount of "good" properties which are not detailed in this documentation.

#### 2.1.4 Cholesky method.

This is a particular application of the Gauss method when  $\mathbf{A}$  is a definite positive symmetric matrix. In this case  $\mathbf{A}$  is factorised as  $\mathbf{A} = \mathbf{B}\mathbf{B}^T$  with the properties  $\mathbf{M} = \mathbf{B}^{-1}$  and  $\mathbf{M}\mathbf{A} = \mathbf{B}^T$ .  $\mathbf{B}$  is a real lower triangular matrix. This method requires  $N^3/6$  additions,  $N^3/6$  multiplications,  $N^2/2$  divisions and  $N$  square roots extractions.

#### 2.1.5 QR factorisation method.

This is a variant of the Gauss method but  $\mathbf{M}$  is chosen to be equal to  $\mathbf{Q}^{-1}$  and  $\mathbf{A}$  is factorised as  $\mathbf{Q}\mathbf{R}$  where:

- $\mathbf{Q}$  is a unitary matrix.
- $\mathbf{R}$  is an upper triangular matrix.
- $\mathbf{Q}^{-1}$  can be written as the product of  $N - 1$  Householder matrices:  $\mathbf{Q}^{-1} = \mathbf{H}_{N-1}\mathbf{H}_{N-2}\dots\mathbf{H}_1$ .

This method can be used for any matrix  $\mathbf{A}$ . If  $\mathbf{A}$  is an invertible matrix, the factorisation  $\mathbf{Q}\mathbf{R}$  is an unique one.

## 2.2 Iterative algorithms for small linear systems.

$\mathbf{D}$  is the diagonal part of  $\mathbf{A}$ ,  $-\mathbf{E}$  (resp.  $-\mathbf{F}$ ) is the lower (resp. upper) triangular part of  $\mathbf{A}$ . So  $\mathbf{A}$  writes  $\mathbf{D} - \mathbf{E} - \mathbf{F}$ . These iterative methods are used for small linear systems and use the decomposition of matrix  $\mathbf{A}$  into its diagonal part, its lower triangular part and its upper triangular part.

### 2.2.1 Jacobi method.

- $\mathbf{A}$  is written  $\mathbf{A} = \mathbf{D} - (\mathbf{E} + \mathbf{F})$ .
- An iteration writes:  $\mathbf{X}_{k+1} = \mathbf{D}^{-1}(\mathbf{E} + \mathbf{F})\mathbf{X}_k + \mathbf{D}^{-1}\mathbf{B}$ .

This method is convergent if the spectral radius (i.e. the maximum of the absolute values of the eigenvalues) of the matrix  $[\mathbf{D}^{-1}(\mathbf{E} + \mathbf{F})]$  is lower than 1.

### 2.2.2 Gauss-Seidel method.

- A is written  $\mathbf{A} = (\mathbf{D} - \mathbf{E}) - \mathbf{F}$ .
- An iteration writes:  $\mathbf{X}_{k+1} = (\mathbf{D} - \mathbf{E})^{-1}\mathbf{F}\mathbf{X}_k + (\mathbf{D} - \mathbf{E})^{-1}\mathbf{B}$ .

This method is convergent if the spectral radius of the matrix  $[(\mathbf{D} - \mathbf{E})^{-1}\mathbf{F}]$  is lower than 1.

### 2.2.3 Relaxation method.

- A is written:

$$\mathbf{A} = \left(\frac{1}{\omega}\mathbf{D} - \mathbf{E}\right) - \left(\frac{1-\omega}{\omega}\mathbf{D} + \mathbf{F}\right)$$

- An iteration writes:

$$\mathbf{X}_{k+1} = \left(\frac{1}{\omega}\mathbf{D} - \mathbf{E}\right)^{-1} \left(\frac{1-\omega}{\omega}\mathbf{D} + \mathbf{F}\right) \mathbf{X}_k + \left(\frac{1}{\omega}\mathbf{D} - \mathbf{E}\right)^{-1} \mathbf{B}$$

This method is convergent if the spectral radius of the matrix  $\left(\frac{1}{\omega}\mathbf{D} - \mathbf{E}\right)^{-1} \left(\frac{1-\omega}{\omega}\mathbf{D} + \mathbf{F}\right)$  is lower than 1. When matrix  $\mathbf{A}$  is a definite positive hermitian matrix, this method is convergent if  $0 < \omega < 2$  and the best convergence is obtained for a value of  $\omega$  which is always between 1 and 2.

## 2.3 Iterative and “mixed” algorithms for large linear systems and minimization problems.

### 2.3.1 General considerations and denotations.

One tries to do a synthesis of existing algorithms and to explain them briefly, preferably using the same denotations for all these algorithms. Literature provides a large number of books or papers explaining these algorithms, with often different denotations. For example, (Bonnans et al., 1997) gives (in French) a good overview of these algorithms, (IDNEWT) gives (in French) a synthetic overview of Newton and quasi-Newton algorithms, (Shewchuk, 1994) is more oriented towards conjugate gradient methods with good practical examples. The following denotations will be used:

- $N$  is the dimension of the problem.
- $N_I$  is the number of iterations required.
- $\mathbf{X}$  is the unknown vector to be found by an iterative algorithm; the successive iterations give  $\mathbf{X}_{iter=0}$ ,  $\mathbf{X}_{iter=1}$ , ...,  $\mathbf{X}_{iter=N_I}$ .  $\mathbf{X}$  has  $N$  components numbered from  $X_{i=1}$  to  $X_{i=N}$ .
- The function to be minimized is  $f(\mathbf{X})$  (this is a scalar function).
- $f$  is not necessary a quadratic function.  $f$  can be written as a sum of its second-order limited development and of a residual term, around a reference value of  $\mathbf{X}$  denoted by  $\mathbf{X}_{ref}$ ; if  $f(\mathbf{X}_{ref})$  is different from zero we can consider the function  $F$  defined by  $F(\mathbf{X}) = f(\mathbf{X}) - f(\mathbf{X}_{ref})$  instead of  $f$ ; if  $\mathbf{X}_{ref}$  is different from zero we can consider the function  $F$  defined by  $F(\mathbf{X} - \mathbf{X}_{ref}) = f(\mathbf{X})$ ; so we can stick to the case where  $f(0) = 0$  and where the second-order limited development is done around  $\mathbf{X}_{ref} = 0$ . In this case the second-order limited development writes:

$$f(\mathbf{X}) = J(\mathbf{X}) = -\langle \mathbf{B}, \mathbf{X} \rangle + 0.5 \langle \mathbf{A}\mathbf{X}, \mathbf{X} \rangle + residual$$

where:

- $\langle \mathbf{X}, \mathbf{Y} \rangle$  is the scalar product of vectors  $\mathbf{X}$  and  $\mathbf{Y}$ . It can be the “canonical” scalar product or some other scalar product.
- $-\mathbf{B}$  is the vector gradient of  $f$  at  $\mathbf{X} = 0$ ; the  $i$ -th component of  $\mathbf{B}$  contains  $-\frac{\partial f}{\partial X_i}(\mathbf{X} = 0)$ .
- Matrix  $\mathbf{A}$  contains the second-order derivatives  $\frac{\partial^2 f}{\partial X_i \partial X_j}(\mathbf{X} = 0)$  and is the Hessian operator of  $f$  at  $\mathbf{X} = 0$ .
- In the particular case where  $f$  is a quadratic function and where  $f(0) = 0$ , the residual “residual” becomes zero, so  $f$ , which is also denoted by  $J$  in this case, writes:

$$f(\mathbf{X}) = J(\mathbf{X}) = 0.5 \langle \mathbf{A}\mathbf{X}, \mathbf{X} \rangle - \langle \mathbf{B}, \mathbf{X} \rangle$$

The Hessian  $J''$  defined by matrix  $\mathbf{A}$  is no longer  $\mathbf{X}$ -dependent. Taking the gradient of this equation writes:

$$\mathbf{g} = [\nabla J](\mathbf{X}) = \mathbf{A}\mathbf{X} - \mathbf{B}$$

A necessary condition of minimum of  $f$  is that  $\mathbf{X}$  zeroes the gradient of the RHS, i.e.:

$$\mathbf{A}\mathbf{X} - \mathbf{B} = 0$$

The searched solution is still solution of a linear system. But the given solution can give a maximum of  $f$ , a minimum of  $f$  or a “saddle” value of  $f$  so if one wants to be sure that it gives a minimum of  $f$  one must assume some additional good properties for  $f$ . If the matrix  $\mathbf{A}$  is a positive definite one, the searched solution gives a minimum of  $f$ .

- Some methods stand to a minimization and can be used for non quadratic functions  $f$ , some other methods apply to quadratic functions and assume the equivalence of “minimizing  $f$ ” and “solving the linear system  $\mathbf{AX} = \mathbf{B}$ ”.
- When solving the linear system  $\mathbf{AX} = \mathbf{B}$ , at each iteration:
  - the residual term is  $\mathbf{r}_{iter} = \mathbf{AX}_{iter} - \mathbf{B}$ .
  - the increment is  $\mathbf{S}_{iter} = \mathbf{X}_{iter+1} - \mathbf{X}_{iter}$ .
  - this increment  $\mathbf{S}_{iter}$  generally writes as  $\rho_{iter}\mathbf{d}_{iter}$  where  $\mathbf{d}_{iter}$  is the direction and  $\rho_{iter}$  is the step.
  - the gradient increment is  $\mathbf{Y}_{iter} = \mathbf{g}_{iter+1} - \mathbf{g}_{iter}$ .

### 2.3.2 Preconditioning.

A good conditioning of matrix  $A$  (i.e. the ratio of the absolute values of the extreme eigenvalues  $r_{\text{cond}} = \frac{|\lambda_{\min}|}{|\lambda_{\max}|}$  is as close as possible to 1) is required to have a quick convergence of the iterative algorithm. Left preconditioning: in some cases it is more interesting to solve

$$\mathbf{H}^{-1}\mathbf{AX} = \mathbf{H}^{-1}\mathbf{B}$$

instead of:

$$\mathbf{AX} = \mathbf{B}$$

with a “good” matrix  $\mathbf{H}$  the properties of which are not always easy to match. We want at the same time that this matrix  $\mathbf{H}$  be:

- easily invertible, and if it is diagonal it would be even better. If  $\mathbf{A}$  is a positive definite symmetric matrix,  $\mathbf{H}$  must be also a positive definite symmetric matrix.
- its “biggest” eigenvalue (in absolute value) is close to  $\lambda_{\max}$ , and its “lowest” eigenvalue (in absolute value) is close to  $\lambda_{\min}$ .
- $\mathbf{H}^{-1}\mathbf{A}$  is not so far from identity and has a very good conditioning (considerably improved compared to the one of  $\mathbf{A}$ ).
- $\mathbf{H}$  can be constant or can be recomputed at each iteration.

Some algorithms to provide a good conditioning will be described later. All algorithms will be described from now on including matrix  $\mathbf{H}$ .

Remark: an alternative way of using a preconditioner is to write the linear system to solve as follows (right preconditioning):

$$\mathbf{AH}^{-1}[\mathbf{HX}] = \mathbf{B}$$

the unknown to find being  $\mathbf{HX}$  instead of  $\mathbf{X}$ .  $\mathbf{H}$  must have the “good properties” listed above.

Right and left preconditioning can be combined, the linear system to solve is rewritten as follows:

$$\mathbf{H}^{-1}\mathbf{AJ}^{-1}[\mathbf{JX}] = \mathbf{H}^{-1}\mathbf{B}$$

$\mathbf{H}$  and  $\mathbf{J}$  are two preconditioners.

One can find information about preconditioning for example in part 3 of (Barrett et al., 1994).

Let us list some examples of ways of computing preconditioning matrices.

- Quasi-Newton preconditioner: the inverse  $\mathbf{H}^{-1}$  can be directly computed with the same recurrence method as for the matrix  $\mathbf{W}$  used in the Quasi-Newton method. It can be considered as a good approximation of the inverse of  $\mathbf{A}$ . For more details see for example (Morales and Nocedal, 1999).
- Another example of preconditioner is given in the part 6.5 of ECMWF documentation (TDECDCAS).

### 2.3.3 Newton method.

This algorithm can be used for non-quadratic functions. One assumes that the function to minimize has a minimum in the “region of minimum”. At each iteration of the minimization algorithm, one will minimize the second-order limited development of  $f$  around the value  $\mathbf{X} = \mathbf{X}_{iter}$ . This is equivalent to solve a linear system which changes at each iteration.

- One starts from an arbitrary “guess”  $\mathbf{X}_0$  of the solution.
- Definition of  $\mathbf{A}_0 = \mathbf{J}''(\mathbf{X}_0)$ .
- Resolution of the linear system  $\mathbf{A}_0\mathbf{d}_0 = -\nabla J(\mathbf{X}_0)$  which provides  $\mathbf{d}_0$ .

- If  $\mathbf{d}_0 = 0$  the algorithm ends.
- If  $\mathbf{d}_0$  is not zero:
  - Definition of the step  $\rho_0$ : see below the general case for iteration “*iter*”.
  - Computation of:

$$\mathbf{X}_1 = \mathbf{X}_0 + \rho_0 \mathbf{d}_0$$

- The following iterations write:
  - Definition of  $\mathbf{A}_{iter} = J''(\mathbf{X}_{iter})$ .
  - Computation of  $\mathbf{g}_{iter} = \nabla J(\mathbf{X}_{iter})$ : that can be done by an analytic formula or by an iterative formula which writes  $\mathbf{g}_{iter} = \mathbf{g}_{iter-1} + \rho_{iter-1} \mathbf{A}_{iter-1} \mathbf{d}_{iter-1}$ .
  - Computation of  $\mathbf{d}_{iter}$ : resolution of the linear system  $\mathbf{A}_{iter} \mathbf{d}_{iter} = -\nabla J(\mathbf{X}_{iter})$  which provides  $\mathbf{d}_{iter}$ .
  - Computation of  $\rho_{iter}$ : several methods can be used, the two known ones are the following ones:
    - \* use of a trust region (this is a rather tricky algorithm, the description of which will not be provided in the current documentation; for more details see for example (IDNEWT)).
    - \* linear research along the descent direction, using a Wolfe algorithm. Roughly one searches for a value of  $\rho_{iter}$  (by dichotomy) which minimizes  $J(\mathbf{X}_{iter} + \rho_{iter} \mathbf{d}_{iter})$  (or a polynomial approximation of this function, often cubic). Moreover this value has to be bounded as follows:

$$\begin{aligned} J(\mathbf{X}_{iter} + \rho_{iter} \mathbf{d}_{iter}) &\leq J(\mathbf{X}_{iter}) + \alpha_1 \rho_{iter} < \nabla J(\mathbf{X}_{iter}); \mathbf{d}_{iter} > \\ &< \nabla J(\mathbf{X}_{iter} + \rho_{iter} \mathbf{d}_{iter}); \mathbf{d}_{iter} > \geq \alpha_2 < \nabla J(\mathbf{X}_{iter}); \mathbf{d}_{iter} > \end{aligned}$$

where  $0 < \alpha_1 < 0.5$  and  $\alpha_1 < \alpha_2 < 1$ .

- Computation of  $\mathbf{X}_{iter+1}$ :
$$\mathbf{X}_{iter+1} = \mathbf{X}_{iter} + \rho_{iter} \mathbf{d}_{iter}$$
- Algorithm stops when the gradient  $\nabla J(\mathbf{X}_{iter})$  has a “sufficiently” small norm (for more details see section (2.3.11)).
- This algorithm is a purely iterative one (in theory it does not give the exact solution in a finite number of iterations).
- The descent direction is generally not parallel to the gradient.
- Since  $\mathbf{A}$  depends on the iteration, the main shortcoming of this algorithm is that it requires the inversion of a big linear system at each iteration, which can be done in theory by some other methods described in this documentation (Fletcher-Reeves conjugate gradient method for example), but which can become very expensive (in time and memory), so the purely Newton method is generally replaced by cheaper algorithms which replace  $\mathbf{A}$  by an approximation easier to compute and to invert. That provides a class of quasi-Newton methods described below.

### 2.3.4 Quasi-Newton method.

This algorithm can be used for non-quadratic functions and is a simplification of the Newton method.  $\mathbf{A}$  is replaced by an approximation  $\mathbf{H}$  simpler to compute and to invert.  $\mathbf{H}$  is generally computed by a recurrence.  $\mathbf{H}$  must have some interesting properties:

- At each iteration,  $\mathbf{H}$  must be a good approximation of  $\mathbf{A}$ .
- $\mathbf{H}$  is easy to compute (generally by an iterative algorithm) and to invert; if  $\mathbf{H}$  is diagonal this is better.

The alternative way is to compute directly a good approximation  $\mathbf{W}$  of  $\mathbf{A}^{-1}$ . In the literature one finds some methods to compute such operators ( $\mathbf{W}$  can be also used as preconditioners for steepest descent or conjugate gradient methods), the most known is the Broyden-Fletcher-Goldfarb-Shanno (abbreviated into BFGS) algorithm. At each iteration:

- the  $N$ -lines and one column matrix  $\mathbf{S}_{iter}$  containing the vector  $(\mathbf{X}_{iter} - \mathbf{X}_{iter-1})$  is stored.
- the  $N$ -lines and one column matrix  $\mathbf{Y}_{iter}$  containing the vector  $(\mathbf{g}_{iter} - \mathbf{g}_{iter-1})$  is stored.
- a guess  $\mathbf{W}_0$  is computed (that can be the identity matrix).
- a recurrence formula gives  $\mathbf{W}_1$  to  $\mathbf{W}_{iter}$ :

$$\mathbf{W}_{iter} = \left( I - \frac{1}{\mathbf{Y}_{iter}^T \mathbf{S}_{iter}} \mathbf{S}_{iter} \mathbf{Y}_{iter}^T \right) \mathbf{W}_{iter-1} \left( I - \frac{1}{\mathbf{Y}_{iter}^T \mathbf{S}_{iter}} \mathbf{S}_{iter} \mathbf{Y}_{iter}^T \right) + \frac{1}{\mathbf{Y}_{iter}^T \mathbf{S}_{iter}} \mathbf{S}_{iter} \mathbf{S}_{iter}^T$$

Since  $\mathbf{W}$  is never stored ( $N * N$  matrix, too big in memory), all the  $\mathbf{S}$  and  $\mathbf{Y}$  have to be stored, and all the recurrence from  $\mathbf{W}_0$  has to be performed to each iteration (for  $\mathbf{S}$  and  $\mathbf{Y}$  the total amount of data stored at iteration “*iter*” is  $(2iter) * N$ ). The previous iterative algorithm is applied to vector  $-\nabla J(\mathbf{X}_{iter})$  and gives the vector  $-\mathbf{W}_{iter} \nabla J(\mathbf{X}_{iter})$ . The quasi-Newton method follows the same scheme as the Newton method, excepted that:

$$\mathbf{d}_{iter} = -\mathbf{A}_{iter}^{-1} \nabla J(\mathbf{X}_{iter})$$

is replaced by:

$$\mathbf{d}_{iter} = -\mathbf{W}_{iter} \nabla J(\mathbf{X}_{iter})$$

### 2.3.5 Limited memory quasi-Newton method.

When the number of iterations  $N_I$  necessary to find the solution becomes important, the number of pairs  $(\mathbf{S}, \mathbf{Y})$  becomes important to be stored, and some algorithms limiting the storage to a “good” subset of pairs  $(\mathbf{S}, \mathbf{Y})$  have been developed.  $N_M$  is the maximal number of pairs kept (for iterations number  $N_M + 1$  to  $N_I$ ). The choice of the iterations kept is generally the  $N_M$  last pairs but it can be something else. Practically a value of  $N_M$  around 5 is a good choice. One now assumes that, for the iteration  $iter$ , the selected pairs are  $(\mathbf{S}_{(iter, m_1)}, \mathbf{Y}_{(iter, m_1)})$  to  $(\mathbf{S}_{(iter, m_{N_M})}, \mathbf{Y}_{(iter, m_{N_M})})$ .  $m_{N_M}$  is generally (but not necessary) equal to  $iter$ ;  $m_1$  is generally (but not necessary) equal to  $iter + 1 - N_M$ ; the list of  $m$  selected can depend of  $iter$ .

For the iteration  $iter$ :

- one defines a guess  $\mathbf{W}_{(iter, m=m_1)}$  which is generally a diagonal matrix  $\mathbf{D}_{iter}$ ; it can be the identity matrix (this is generally the case for  $\mathbf{D}_0$ ) but some authors have developed better algorithms to compute  $\mathbf{D}_{iter}$ ; for example the scalar initial scaling mode (“SIS” algorithm) yields:

$$\mathbf{D}_{iter} = \frac{\mathbf{Y}_{iter-1}^T \mathbf{S}_{iter-1}}{\mathbf{Y}_{iter-1}^T \mathbf{Y}_{iter-1}} \mathbf{I}$$

- the recurrence formula giving  $\mathbf{W}_{(iter, m=m_i)}$  knowing  $\mathbf{W}_{(iter, m=m_{i-1})}$  writes:

$$\mathbf{W}_{(iter, m=m_i)} = \mathbf{F}_{(iter, m=m_i)} \mathbf{W}_{(iter, m=m_{i-1})} \mathbf{F}_{(iter, m=m_i)} + \frac{1}{\mathbf{Y}_{(iter, m=m_i)}^T \mathbf{S}_{(iter, m=m_i)}} \mathbf{S}_{(iter, m=m_i)} \mathbf{S}_{(iter, m=m_i)}^T$$

where  $\mathbf{F}$  is:

$$\mathbf{F}_{(iter, m=m_i)} = \left( \mathbf{I} - \frac{1}{\mathbf{Y}_{(iter, m=m_i)}^T \mathbf{S}_{(iter, m=m_i)}} \mathbf{S}_{(iter, m=m_i)} \mathbf{Y}_{(iter, m=m_i)}^T \right)$$

More details are given for example in (IDQN3).

### 2.3.6 Truncated Newton method.

Looks like the Newton method. The resolution of the linear system  $\mathbf{A}_{iter} \mathbf{d}_{iter} = -\nabla J(\mathbf{X}_{iter})$  is still done (without approximation of the Hessian) by an iterative algorithm but the convergence criterion is coarse for lower values of  $iter$  (so the resolution can be done quickly with a small amount of iterations); when  $iter$  increases the convergence criterion becomes more severe (in the purely Newton method the convergence criterion does not depend on  $iter$ ). To sum up, the computation of  $\mathbf{d}_{iter}$  is relatively cheap for lower values of  $iter$  and more and more expensive when  $iter$  increases.

### 2.3.7 Fletcher-Reeves conjugate gradient method.

This algorithm can be used for quadratic functions. One assumes that the function is a positive definite quadratic one. The aim is to minimize:

$$f(\mathbf{X}) = J(\mathbf{X}) = 0.5 \langle \mathbf{A}\mathbf{X}, \mathbf{X} \rangle - \langle \mathbf{B}, \mathbf{X} \rangle$$

and this equivalent to solve the following linear system:

$$\mathbf{H}^{-1} \mathbf{A}\mathbf{X} = \mathbf{H}^{-1} \mathbf{B}$$

- One starts from an arbitrary “guess”  $\mathbf{X}_0$  of the solution.
- Definition of  $\mathbf{H}_0$ .
- Computation of  $\mathbf{d}_0 = -\mathbf{H}_0^{-1} \nabla J(\mathbf{X}_0)$ .
- If  $\mathbf{d}_0 = 0$  the algorithm ends.
- If  $\mathbf{d}_0$  is not zero:

– Definition of:

$$\rho_0 = \frac{\langle \nabla J(\mathbf{X}_0); -\mathbf{d}_0 \rangle}{\langle \mathbf{A}\mathbf{d}_0; \mathbf{d}_0 \rangle} = \frac{\langle \mathbf{H}_0^{-1} \nabla J(\mathbf{X}_0); \nabla J(\mathbf{X}_0) \rangle}{\langle \mathbf{A}\mathbf{d}_0; \mathbf{d}_0 \rangle}$$

– Computation of:

$$\mathbf{X}_1 = \mathbf{X}_0 + \rho_0 \mathbf{d}_0$$

- The following iterations write:

– Definition of  $\mathbf{H}_{iter}$  and computation of  $\mathbf{H}_{iter}^{-1}$ .

– Computation of  $\mathbf{g}_{iter} = \nabla J(\mathbf{X}_{iter})$ : that can be done by an analytic formula or by an iterative formula which writes  $\mathbf{g}_{iter} = \mathbf{g}_{iter-1} + \rho_{iter-1} \mathbf{A}\mathbf{d}_{iter-1}$ .



- Computation of  $\mathbf{d}_{iter}$ :

$$\mathbf{d}_{iter} = -\mathbf{H}_{iter}^{-1} \nabla J(\mathbf{X}_{iter}) + \frac{\langle \mathbf{H}_{iter}^{-1} \nabla J(\mathbf{X}_{iter}); \nabla J(\mathbf{X}_{iter}) \rangle}{\langle \mathbf{H}_{iter-1}^{-1} \nabla J(\mathbf{X}_{iter-1}); \nabla J(\mathbf{X}_{iter-1}) \rangle} \mathbf{d}_{iter-1}$$

- Computation of  $\rho_{iter}$ :

$$\rho_{iter} = \frac{\langle \mathbf{H}_{iter}^{-1} \nabla J(\mathbf{X}_{iter}); \nabla J(\mathbf{X}_{iter}) \rangle}{\langle \mathbf{A} \mathbf{d}_{iter}; \mathbf{d}_{iter} \rangle} = \frac{\langle \nabla J(\mathbf{X}_{iter}); -\mathbf{d}_{iter} \rangle}{\langle \mathbf{A} \mathbf{d}_{iter}; \mathbf{d}_{iter} \rangle}$$

- Computation of  $\mathbf{X}_{iter+1}$ :

$$\mathbf{X}_{iter+1} = \mathbf{X}_{iter} + \rho_{iter} \mathbf{d}_{iter}$$

- Algorithm stops when the gradient  $\nabla J(\mathbf{X}_{iter})$  has a “sufficiently” small norm (for more details see section (2.3.11)).
- This algorithm converges in a finite number of iterations, the maximum possible number is  $N$ .
- Two different iterations of  $\nabla J(\mathbf{X})$  (and also two different iterations of  $\mathbf{d}_{iter}$ ) are “orthogonal” relative to the scalar product defined with matrix  $\mathbf{H}^{-1}$ .
- The descent direction is not parallel to the gradient.

### 2.3.8 Polak-Ribière conjugate gradient method.

For purely quadratic functions this algorithm is identical to the Fletcher-Reeves conjugate gradient method. This algorithm is an extension of the Fletcher-Reeves one for non-quadratic functions. When the function to minimize is not a quadratic one the differences with the Fletcher-Reeves algorithm are:

- The new way of computing the direction  $\mathbf{d}_{iter}$  is:

$$\mathbf{d}_{iter} = -\mathbf{H}_{iter}^{-1} \nabla J(\mathbf{X}_{iter}) + \frac{\langle \mathbf{H}_{iter}^{-1} \nabla J(\mathbf{X}_{iter}); \nabla J(\mathbf{X}_{iter}) - \nabla J(\mathbf{X}_{iter-1}) \rangle}{\langle \mathbf{H}_{iter-1}^{-1} \nabla J(\mathbf{X}_{iter-1}); \nabla J(\mathbf{X}_{iter-1}) \rangle} \mathbf{d}_{iter-1}$$

- This algorithm does not converge in a finite number of iterations.
- Two different iterations of  $\nabla J(\mathbf{X})$  are no longer necessarily “orthogonal” relative to the scalar product defined with matrix  $\mathbf{H}^{-1}$ .

### 2.3.9 Gradient method with optimal step (or steepest descent method).

This algorithm can be used for quadratic functions. One assumes that the function is a positive definite quadratic one. This algorithm looks like the Fletcher-Reeves conjugate gradient method excepted on the following points:

- Computation of  $\mathbf{d}_{iter}$ :

$$\mathbf{d}_{iter} = -\mathbf{H}_{iter}^{-1} \nabla J(\mathbf{X}_{iter})$$

- The descent direction is parallel to the gradient and orthogonal to the iso- $J$ .
- This algorithm does not converge in a finite number of iterations.
- Two different iterations of  $\nabla J(\mathbf{X})$  are no longer necessary “orthogonal” relative to the scalar product defined with matrix  $\mathbf{H}^{-1}$ .

### 2.3.10 Conjugate gradient method combined with a Lanczos algorithm.

This algorithm is a combination of a Fletcher-Reeves conjugate gradient algorithm with a Lanczos algorithm. The main difference with the Fletcher-Reeves algorithm is that  $\mathbf{X}_{iter+1}$  is not computed using directly formula  $\tilde{\mathbf{X}}_{iter+1} = \tilde{\mathbf{X}}_{iter} + \rho_{iter} \mathbf{d}_{iter}$  but using a more tricky algorithm involving a matrix  $\mathbf{T}$  of dimensions  $(iter + 1) * (iter + 1)$ , the eigenvalues of which converge towards the eigenvalues of  $\mathbf{A}$  when  $iter$  increases.

This algorithm can be used for quadratic functions. One assumes that the function is a positive definite quadratic one. The aim is to minimize:

$$f(\mathbf{X}) = J(\mathbf{X}) = 0.5 \langle \mathbf{A} \mathbf{X}, \mathbf{X} \rangle - \langle \mathbf{B}, \mathbf{X} \rangle$$

and this equivalent to solve the following linear system:

$$\mathbf{H}^{-1} \mathbf{A} \mathbf{X} = \mathbf{H}^{-1} \mathbf{B}$$

- One starts from an arbitrary “guess”  $\mathbf{X}_0$  of the solution.
- Definition of  $\mathbf{H}_0$ .
- Computation of  $\mathbf{d}_0 = -\mathbf{H}_0^{-1} \nabla J(\mathbf{X}_0)$ .

- If  $\mathbf{d}_0 = 0$  the algorithm ends.
- If  $\mathbf{d}_0$  is not zero:

- Definition of:

$$\rho_0 = \frac{\langle \nabla J(\mathbf{X}_0); -\mathbf{d}_0 \rangle}{\langle \mathbf{A}\mathbf{d}_0; \mathbf{d}_0 \rangle} = \frac{\langle \mathbf{H}_0^{-1} \nabla J(\mathbf{X}_0); \nabla J(\mathbf{X}_0) \rangle}{\langle \mathbf{A}\mathbf{d}_0; \mathbf{d}_0 \rangle}$$

(identical to Fletcher-Reeves conjugate gradient algorithm).

- Definition of:

$$\mathbf{q}_0 = \frac{1}{\langle \mathbf{H}_0^{-1} \nabla J(\mathbf{X}_0); \nabla J(\mathbf{X}_0) \rangle} \nabla J(\mathbf{X}_0)$$

- Definition of  $\mathbf{Q}_1$  which is the  $N$ -lines and one column matrix containing vector  $\mathbf{q}_0$ .
- Definition of:  $\delta_1 = \frac{1}{\rho_0}$ .
- Definition of the  $1 * 1$  matrix  $\mathbf{T}_1 = [\delta_0]$ .
- Definition of the  $1 * 1$  matrix:

$$\mathbf{Z}_1 = -\mathbf{T}_1^{-1} \mathbf{Q}_1^T \mathbf{H}_0^{-1} \nabla J(\mathbf{X}_0)$$

which is simply the scalar  $-\rho_0$ .

- Computation of:

$$\mathbf{X}_1 = \mathbf{X}_0 + \mathbf{H}_0^{-1} \mathbf{Q}_1 \mathbf{Z}_1$$

which is equivalent to:

$$\mathbf{X}_1 = \mathbf{X}_0 + \rho_0 \mathbf{d}_0$$

since  $\mathbf{Z}_1 = \rho_0$  and  $\mathbf{H}_0^{-1} \mathbf{Q}_1 = -\mathbf{d}_0$ .

- The following iterations write:

- Definition of  $\mathbf{H}_{iter}$  and computation of  $\mathbf{H}_{iter}^{-1}$ .
- Computation of  $\mathbf{g}_{iter} = \nabla J(\mathbf{X}_{iter})$ : that can be done by an analytic formula or by an iterative formula which writes  $\mathbf{g}_{iter} = \mathbf{g}_{iter-1} + \rho_{iter-1} \mathbf{A}\mathbf{d}_{iter-1}$  (identical to Fletcher-Reeves conjugate gradient algorithm).
- Computation of  $\mathbf{d}_{iter}$ :

$$\mathbf{d}_{iter} = -\mathbf{H}_{iter}^{-1} \nabla J(\mathbf{X}_{iter}) + \frac{\langle \mathbf{H}_{iter}^{-1} \nabla J(\mathbf{X}_{iter}); \nabla J(\mathbf{X}_{iter}) \rangle}{\langle \mathbf{H}_{iter-1}^{-1} \nabla J(\mathbf{X}_{iter-1}); \nabla J(\mathbf{X}_{iter-1}) \rangle} \mathbf{d}_{iter-1}$$

(identical to Fletcher-Reeves conjugate gradient algorithm).

- Computation of  $\rho_{iter}$ :

$$\rho_{iter} = \frac{\langle \mathbf{H}_{iter}^{-1} \nabla J(\mathbf{X}_{iter}); \nabla J(\mathbf{X}_{iter}) \rangle}{\langle \mathbf{A}\mathbf{d}_{iter}; \mathbf{d}_{iter} \rangle} = \frac{\langle \nabla J(\mathbf{X}_{iter}); -\mathbf{d}_{iter} \rangle}{\langle \mathbf{A}\mathbf{d}_{iter}; \mathbf{d}_{iter} \rangle}$$

(identical to Fletcher-Reeves conjugate gradient algorithm).

- Definition of:

$$\mathbf{q}_{iter} = \frac{1}{\langle \mathbf{H}_{iter}^{-1} \nabla J(\mathbf{X}_{iter}); \nabla J(\mathbf{X}_{iter}) \rangle} \nabla J(\mathbf{X}_{iter})$$

- Definition of  $\mathbf{Q}_{iter+1}$  which is the  $N$ -lines and  $iter + 1$ -columns matrix containing vectors  $\mathbf{q}_0$  to  $\mathbf{q}_{iter}$ .
- Definition of:

$$\delta_{iter+1} = \frac{\langle \mathbf{H}_{iter}^{-1} \nabla J(\mathbf{X}_{iter}); \nabla J(\mathbf{X}_{iter}) \rangle^2}{\langle \mathbf{H}_{iter-1}^{-1} \nabla J(\mathbf{X}_{iter-1}); \nabla J(\mathbf{X}_{iter-1}) \rangle^2} \frac{1}{\rho_{iter-1}} + \frac{1}{\rho_{iter}}$$

- Definition of:

$$\gamma_{iter} = -\frac{\langle \mathbf{H}_{iter}^{-1} \nabla J(\mathbf{X}_{iter}); \nabla J(\mathbf{X}_{iter}) \rangle}{\langle \mathbf{H}_{iter-1}^{-1} \nabla J(\mathbf{X}_{iter-1}); \nabla J(\mathbf{X}_{iter-1}) \rangle} \frac{1}{\rho_{iter-1}}$$

- Definition of the  $(iter + 1) * (iter + 1)$  matrix  $\mathbf{T}_{iter+1}$ :

- \*  $\mathbf{T}_{iter+1}$  is a symmetric tri-diagonal matrix.
- \* The main diagonal of  $\mathbf{T}_{iter+1}$  contains  $\delta_1$  to  $\delta_{iter+1}$ .
- \* The first lower and upper diagonal of  $\mathbf{T}_{iter+1}$  contains  $\gamma_1$  to  $\gamma_{iter}$ .

Computation of its inverse  $\mathbf{T}_{iter+1}^{-1}$ .

– Definition of the  $iter + 1$ -lines and one column matrix:

$$\mathbf{Z}_{iter+1} = -\mathbf{T}_{iter+1}^{-1} \mathbf{Q}_{iter+1}^T \mathbf{H}_{iter}^{-1} \nabla J(\mathbf{X}_0)$$

– Computation of:

$$\mathbf{X}_{iter+1} = \mathbf{X}_0 + \mathbf{H}_{iter}^{-1} \mathbf{Q}_{iter+1} \mathbf{Z}_{iter+1}$$

- Algorithm stops when the gradient  $\nabla J(\mathbf{X}_{iter})$  has a “sufficiently” small norm (for more details see section (2.3.11)).
- This algorithm converges in a finite number of iterations, the maximum possible number is  $N$ .
- Two different iterations of  $\nabla J(\mathbf{X})$  (and also two different iterations of  $\mathbf{d}_{iter}$ ) are “orthogonal” relative to the scalar product defined with matrix  $\mathbf{H}^{-1}$ .
- The descent direction is not parallel to the gradient.
- Matrix  $\mathbf{T}$  has some good properties:
  - Eigenvalues of  $\mathbf{T}_{iter}$  converge towards the eigenvalues of  $\mathbf{A}$  when  $iter$  increases.
  - Writing  $\mathbf{T}_{iter} = \mathbf{P}_{iter}^{-1} \mathbf{D}_{iter} \mathbf{P}_{iter}$  where  $\mathbf{D}_{iter}$  is diagonal and  $\mathbf{P}_{iter}$  contains the eigenvectors of  $\mathbf{T}_{iter}$ , the product  $\mathbf{Q}_{iter} \mathbf{P}_{iter}$  converges towards the eigenvectors of  $\mathbf{A}$ .

### 2.3.11 When to stop the algorithm.

This topic is discussed in chapter 4.2 of (Barrett et al., 1994). The stopping criterion is generally:

$$\frac{\langle \mathbf{g}_{iter}, \mathbf{g}_{iter} \rangle}{\langle \mathbf{g}_0, \mathbf{g}_0 \rangle} < \epsilon \quad (3)$$

An additional criterion can be used:

$$\langle \mathbf{X}_{iter} - \mathbf{X}_{iter-1}, \mathbf{X}_{iter} - \mathbf{X}_{iter-1} \rangle < \epsilon_2 \quad (4)$$

Some additional constraints can be added in some of the minimizers, for example:

- In **N1CG1**: the cost function must not increase, in the contrary the algorithm stops.
- In **CONGRAD**: the eigenvalues of the matrix  $\mathbf{T}$  obtained by the Lanczos algorithm must be all  $\geq 0$ , in the contrary the algorithm stops at the last iteration where all the eigenvalues are  $\geq 0$ .
- In **CONGRAD** one checks that matrix  $\mathbf{A}$  is positive definite, in the contrary the algorithm stops.
- In **CONGRAD** the algorithm stops when the ratio final gradient / initial gradient is below a predefined threshold, or when the final gradient is below a predefined threshold.

### 3 Available minimizers in the “arpifs” and “algor” libraries: description and organigrammes.

Codes are in library “algor” excepted for **CONGRAD** which is in the main library “arpifs”.

#### 3.1 Minimizer **M1QN3** developed at INRIA.

**M1QN3** follows a limited memory quasi-Newton method. It can be used for non-quadratic functions.

##### \* Organigramme of **M1QN3**:

```
M1QN3 ->
* ALLOCATE_CTLVEC
* M1QN3A ->
  - ALLOCATE_CTLVEC
  - [ YSTBL (not described in detail) ] [commented]
  - [PROSCA]
  - DD ->
    * [PROSCA]
    * [CTONB]
    * [CTCAB]
  - DDS ->
    * [ YSTBL (not described in detail) ] [commented]
    * [PROSCA]
    * [CTONB]
    * [CTCAB]
  - [SIMUL]
  - MLISO ->
    * [PROSCA]
    * [SIMUL]
    * ECUBE
  - DEALLOCATE_CTLVEC
* [PROSCA]
* DEALLOCATE_CTLVEC
```

##### \* Comments about routines involved:

- **M1QN3**: is the main routine; calls **M1QN3A** after having structured the available memory.
- **M1QN3A**: is the body of the minimizer.
- **ALLOCATE\_CTLVEC** (resp. **DEALLOCATE\_CTLVEC**) (in module-routine **CONTROL\_VECTORS**): allocation (resp. deallocation) of some “CONTROL\_VECTOR”-type arrays.
- [ **PROSCA** ]: dummy argument generic name of a routine doing a scalar product.
- [ **CTONB** ] and [ **CTCAB** ]: dummy argument generic name of routines doing a basis change and its inverse.
- **DD**: computes the descent direction; manages multiplications by the inverse of the preconditioner.
- **DDS**: does the same as **DD** but is used in the case when the pairs (**Y**;**S**) are not completely stored in memory but read on a buffer via the routine **YSTBL**.
- **MLISO**, **ECUBE**: subroutines realizing the line-search.
- [ **SIMUL** ]: dummy argument generic name of the simulator.
- **YSTBL**: is used for very large problems (currently commented).

\* **Additional remarks**: A more detailed documentation of **M1QN3** can be found in (IDQN3). Preconditioning is a left one and can be done by a quasi-Newton preconditioner.

## 3.2 Minimizer M1QN3R.

**M1QN3R** follows a limited memory quasi-Newton method. It can be used for non-quadratic functions. This is a version of **M1QN3** without the “CONTROL\_VECTOR” structure.

### \* Organigramme of M1QN3R:

```
M1QN3R ->
* MUPDTSR
* M1QN3AR ->
  - [PROSCAR]
  - DDR ->
    * [PROSCAR]
    * [CTONBR]
    * [CTCABR]
  - DDSR ->
    * YSTBLR (not described in detail)
    * [PROSCAR]
    * [CTONBR]
    * [CTCABR]
  - [SIMULR]
  - MLISOR ->
    * [PROSCAR]
    * [SIMULR]
    * ECUBER
  - YSTBLR (not described in detail)
* [PROSCAR]
```

### \* Comments about routines involved:

- **M1QN3R**: is the main routine; calls **M1QN3AR** after having structured the available memory.
- **M1QN3AR**: is the body of the minimizer.
- **MUPDTSR**: returns the number of updates to form the approximate Hessian.
- [ **PROSCAR** ]: dummy argument generic name of a routine doing a scalar product.
- [ **CTONBR** ] and [ **CTCABR** ]: dummy argument generic name of routines doing a basis change and its inverse.
- **DDR**: computes the descent direction; manages multiplications by the inverse of the preconditioner.
- **DDSR**: does the same as **DDR** but is used in the case when the pairs (**Y**;**S**) are not completely stored in memory but read on a buffer via the routine **YSTBLR**.
- **MLISOR**, **ECUBER**: subroutines realizing the line-search.
- [ **SIMULR** ]: dummy argument generic name of the simulator.
- **YSTBLR**: is used for very large problems.

\* **Additional remarks:** A more detailed documentation of **M1QN3R** can be found in (IDQN3). Preconditioning is a left one and can be done by a quasi-Newton preconditioner.

### 3.3 Minimizer N1CG1.

N1CG1 follows a Fletcher-Reeves conjugate gradient method. It can be used to solve positive definite quadratic functions or to solve linear systems involving a symmetric positive definite matrix.

#### \* Organigramme of N1CG1:

```
N1CG1 ->
* ALLOCATE_CTLVEC
* [SIMUL]
* N1CGA
  - DPSEUCLID -> DOT_PRODUCT
  - DBFGSL -> DPSEUCLID -> DOT_PRODUCT
  - [SIMUL]
  - DYSAVE -> DPSEUCLID -> DOT_PRODUCT
* DEALLOCATE_CTLVEC
```

#### \* Comments about routines involved:

- **N1CG1**: is the main routine; calls **N1CG1A**.
- **N1CG1A**: is the body of the minimizer.
- **ALLOCATE\_CTLVEC**  
(resp. **DEALLOCATE\_CTLVEC**) (in module-routine **CONTROL\_VECTORS**): allocation (resp. deallocation) of some “CONTROL\_VECTOR”-type arrays.
- **DPSEUCLID** does the “basic” scalar product and calls **DOT\_PRODUCT**.
- **DBFGSL**: computes the descent direction; manages multiplications by the inverse of the preconditioner. This routine is nearly identical to the routine **DD** called under **M1QN3**.
- **DYSAVE** is called when asking for preconditioning does the selection and saving of pairs (**Y**;**S**).
- [ **SIMUL** ]: dummy argument generic name of the simulator.

\* **Additional remarks:** This routine can be used with a left quasi-Newton preconditioner; see for example (Morales and Nocedal, 1999) for the description of such a preconditioner.

### 3.4 Minimizer CONGRAD.

**CONGRAD** follows a conjugate gradient method combined with a Lanczos algorithm. It can be used to solve positive definite quadratic functions or to solve linear systems involving a symmetric positive definite matrix. Preconditioning is a right one.

#### \* Organigramme of CONGRAD:

```
CONGRAD ->
* ALLOCATE_CTLVEC
* SIM4D -> (call tree not described)
* OPK -> (call tree not described)
* PRECOND
* PTSV -> [LAPACK]/SPTSV ou [LAPACK]/DPTSV
* STEQR -> [LAPACK]/SSTEQR ou [LAPACK]/DSTEQR
* WREVECS ->
  - ALLOCATE_CTLVEC
  - DEALLOCATE_CTLVEC
* DEALLOCATE_CTLVEC
```

The call to **XFORMEV** has been moved in the routines which call **CONGRAD**.

#### \* Organigramme of the setup part computing the preconditioner:

```
CVA1 ->
* PREPPCM
  - ALLOCATE_CTLVEC
  - DEALLOCATE_CTLVEC
  - ALLGATHER_CVSECTION
  - SET_CVSECTION
```

#### \* Comments about routines involved:

- **CONGRAD**: is the main routine and the body of the minimizer.
- **ALLGATHER\_CVSECTION** (in module-routine **CONTROL\_VECTORS\_COMM\_MOD**): gather part of the control variable on all processors.
- **ALLOCATE\_CTLVEC** (resp. **DEALLOCATE\_CTLVEC**) (in module-routine **CONTROL\_VECTORS**): allocation (resp. deallocation) of some “CONTROL\_VECTOR”-type arrays.
- **SIM4D**: is the simulator used in the 4D-VAR (no other simulator is currently allowed here).
- **OPK**: is the simulator used in the calculation of singular vectors.
- **PRECOND**: applies the preconditioner.
- **PREPPCM**: pre-computes the preconditioner.
- **PTSV**: resolution of a real linear system of equations, the matrix of which being symmetric positive definite and tridiagonal.
- **SET\_CVSECTION** (in module-routine **CONTROL\_VECTORS\_COMM\_MOD**): set part of the control variable on all processors.
- **STEQR**: computes the eigenvalues and, optionally, eigenvectors of a symmetric tridiagonal matrix.
- **CVA1**: controls variational assimilation job at level 1.
- **WREVECS**: computes and saves eigenvectors.
- **XFORMEV**: transforms eigenvectors towards model variable space.
- The “basic” scalar product is directly done by calls of **DOT\_PRODUCT**.

\* **Additional remarks:** A more detailed description of the algorithm used by **CONGRAD** can be found in (Fisher, 1998) and in the part 6 of ECMWF documentation (TDECDCAS). The preconditioning is done by a change of variable inside routine **PRECOND** (solve  $AH^{-1}[HX] = B$  instead of  $AX = B$ ); calculation of the preconditioner is described in the part 6.5 of ECMWF documentation (TDECDCAS). The preconditioning matrix is computed in routine **PREPPCM**; the change of variable is done inside routine **PRECOND**.

## 4 Configurations of use of minimizations, and simulators.

### 4.1 Minimization for orography in climatological files (configuration 923).

It uses a non-quadratic function, and a minimization by the quasi-Newton algorithm:

- Minimizer: **M1QN3R**.
- Simulator: **SIMREL**.
- Size of the problem: a local variable **IRZ** computed in **RELSPE**, which is a function of **NSEFRE**.
- Caller: **APLM1G**.

### 4.2 Minimization in the 4DVAR assimilation.

It uses a quadratic function.

- Minimizer: the minimizers which can be used are **M1QN3**, **N1CG1** or **CONGRAD**.
- Simulator: **SIM4D**. A more detailed description of **SIM4D** can be found in part 4.2.3 of (IDVAR) and in the IFS technical documentation part II (section 2.3).
- Size of the problem: **NSEFRE**.
- Callers: **CFCSSENS2OBS**, **CVA2**, **FORECAST\_ERROR**.

### 4.3 Other applications using minimizers.

- **NALAN1** calls **CONGRAD** (simulator **SIM4D** or **OPK**).



## 5 Modules and namelists.

These modules are auto-documented so description of each variable is provided in the code source. We can recall here the most important variables to know for each module:

- **PTRSPOR**: used only in the minimization of orography done in the part 1 of configuration 923.
- **YOMCVA** (control variable).
- **YOMDIM** (dimensioning), in particular NSEFRE.
- **YOMGETMINI**: contains variables which are at least used when getting the pre-conditioner for CONGRAD.
- **YOMIOMI** (warm restart I/O mechanism for minimization). Some of these variables are in namelist **NAMIOMI**.
- **YOMSENS** (variables for sensitivity). Some of these variables are in namelist **NAMSENS**.
- **YOMVAR**: a subset of variables are used in minimizers or in applications calling minimizers.
  - NITER, NITER\_MIN, NSIMU, NINFRA.
  - NMIMP, RXMIN (used for M1QN3).
  - LN1CG1, ZEPSNEG, N1IMP, NSELECT, NPRECO, NBFGB (used for N1CG1).
  - LCONGRAD, NPCVECS, LWRIEVEC, NWRIEVEC, N\_DIAGS\_CONVERGENCE, N\_DIAGS\_EIGENVECS, LEVECCNTL, LEVECGRIB, RTOL\_CHECK\_GRADIENT, LMPCGL, EVBCGL, MCGLVEC, LWCGL (used for CONGRAD).
  - RCVGE, R\_NORM\_REDUCTION\_ABORT\_LEVEL, L\_CHECK\_CONVERGENCE, L\_ABS\_CONVERGENCE, L\_INFO\_CONVERGENCE.
  - L\_INFO\_CONTENT, N\_INFO\_CONTENT\_METHOD, N\_INFO\_CONTENT\_SEED.
  - and maybe also LSKIPMIN, NUPTRA, NHEVECS.

Some of them are in namelist **NAMVAR**.

- **YOMVCGL**: global arrays for intermediate results of the forecast error calculation or used in the CONGRAD minimization (preconditioner).
- **CONTROL\_VECTORS\_COMM\_MOD**: contains routines which do operations on “derivated type” variables used in the minimizations (routines in this module seem linked with distributed memory).
- **CONTROL\_VECTORS** and **CONTROL\_VECTORS\_...** (in library “algor”): contains routines which do operations on “derivated type” variables used in the minimizations.

## 6 Bibliography.

### 6.1 Publications.

- Barrett R., M. Berry, T. F. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine and H. Van der Vorst, 1994: Templates for the solution of linear systems: building blocks for iterative methods, 2nd edition, Pub. SIAM, Philadelphia. Available on <http://www.netlib.org/templates/Templates.html>
- Bonnans, J.F., J.C. Gilbert, C. Lemaréchal and C. Sagastizabal, 1997: Optimisation numérique, aspect théorique et pratique. (Collection Mathématiques et Applications, édité chez Springer).
- Ciarlet, P.G., 1982: Introduction à l'analyse numérique matricielle et à l'optimisation. (Collection Mathématiques appliquées pour la maîtrise, édité chez Masson, 279pp).
- Courtier, Ph., C. Freyrier, J.F. Geleyn, F. Rabier and M. Rochas, 1991: The ARPEGE project at METEO-FRANCE. ECMWF Seminar Proceedings 9-13 September 1991, Volume II, 193-231.
- Fisher M., 1998: Minimization algorithms for variational data assimilation. Proceedings of ECMWF seminar on recent developments in numerical methods for atmospheric modelling, 7-11 Sept 1998, pp 364-385.
- Fletcher, R., 1993: An overview of unconstrained optimization. Pre-print, 31pp. Postscript file available as <http://www.cerfacs.fr/~wlay/LAY/Publications/Optifletcher.ps.gz>
- Morales, J.L., and J. Nocedal, 1999: Automatic preconditioning by limited memory quasi-Newton updating. Pre-print, 22pp.
- Shewchuk, J.R., 1994: An introduction to the conjugate gradient method without the agonizing pain. Edition 1 $\frac{1}{4}$ . Note, 64pp, available on <http://www.cs.cmu.edu/~jrs/>

### 6.2 Some internal notes and other ARPEGE notes.

- (TDECDAS) 2016: IFS technical documentation (CY41R2). Part II: data assimilation. Available at "<https://software.ecmwf.int/wiki/display/IFS/Official+IFS+Documentation>".
- (TDECTEC) 2016: IFS technical documentation (CY41R2). Part VI: technical and computational procedures. Available at "<https://software.ecmwf.int/wiki/display/IFS/Official+IFS+Documentation>".
- (IDVAR) Fischer, C., and L. Berre, 2007: The variational computations inside ARPEGE/ALADIN: cycle CY32. Internal note, 77pp. Available on the intranet server "<http://www.umr-cnrm.fr/gmapdoc/>".
- (IDQN3) Gilbert, J.C., and C. Lemaréchal, 1993: The modules M1QN3 and N1QN3. Version 2.0b. Internal note, 12pp.
- (IDNEWT) Veersé, F., 1996. Quelques algorithmes de minimisation pour l'assimilation variationnelle de données météorologiques. Internal note (slides-sized), 26pp.
- (IDBAS) Yessad, K., 2017: Basics about ARPEGE/IFS, ALADIN and AROME in the cycle 44 of ARPEGE/IFS (internal note).
- (IDEUL) Yessad, K., 2017: Integration of the model equations, and Eulerian dynamics, in the cycle 44 of ARPEGE/IFS (internal note).
- (IDDM) Yessad, K., 2017: Distributed memory features in the cycle 44 of ARPEGE/IFS (internal note).