# Convergence Rate Analysis of the Majorize-Minimize Subspace Algorithm – Extended Version

Emilie Chouzenoux and Jean-Christophe Pesquet *

July 20, 2016

**Abstract**

State-of-the-art methods for solving smooth optimization problems are nonlinear conjugate gradient, low memory BFGS, and Majorize-Minimize (MM) subspace algorithms. The MM subspace algorithm which has been introduced more recently has shown good practical performance when compared with other methods on various optimization problems arising in signal and image processing. However, to the best of our knowledge, no general result exists concerning the theoretical convergence rate of the MM subspace algorithm. This paper aims at deriving such convergence rates both for batch and online versions of the algorithm and, in particular, discusses the influence of the choice of the subspace.

*Keywords*: convergence rate, optimization, subspace algorithms, memory gradient methods, descent methods, majorization-minimization, online optimization, learning.

## 1 Introduction

The Majorize-Minimize (MM) subspace algorithm [1] is based on the idea of constructing, at the current iteration, a quadratic majorizing approximation of the cost function of interest [2], and generating the next iterate by minimizing this surrogate function within a subspace spanned by few directions [3–5]. Note that the MM subspace algorithm can be viewed as a special instance of nonlinear conjugate gradient (NLCG) [6] with closed form formula for the stepsize and conjugacy parameter, or as a particular low memory BFGS (L-BFGS) algorithm [7] with a specific combination of memory directions. The MM subspace algorithm enjoys nice convergence properties [8], and shows good performance in practice, when compared with NLCG, L-BFGS, and also with graph-cut based discrete optimization methods, and proximal algorithms [1,9,10]. It has recently been extended to the online case when only a stochastic approximation of the criterion is employed at each iteration [11]. All these works illustrate the fact that the choice of the subspace has a major impact on the practical convergence speed of the algorithm (see, for instance [1, Section 5], [8, Section 5.1]). In particular, it seems that the best performance is obtained for the memory gradient subspace [12], spanned by the current gradient and the previous direction, leading to the so-called MM Memory Gradient (3MG) algorithm. However, only an analysis concerning the convergence rates of half-quadratic algorithms (corresponding to the case when the subspace spans the whole Euclidean space) is available [13, 14].

Section 2 describes the general form of the MM subspace algorithm and its main known properties. In Section 3, a convergence rate analysis is performed for both batch and online versions of the algorithm for minimizing a wide class of strongly convex cost functions.

## 2 MM subspace algorithm

### 2.1 Optimization problem

In this paper, we will be interested in the minimization of the penalized quadratic cost function:

$$F \colon \mathbb{R}^N \to \mathbb{R} \colon \boldsymbol{h} \mapsto \frac{1}{2} \boldsymbol{h}^\top \boldsymbol{R} \boldsymbol{h} - \boldsymbol{r}^\top \boldsymbol{h} + \Psi(\boldsymbol{h}), \tag{1}$$

where $\boldsymbol{r} \in \mathbb{R}^N$, $\boldsymbol{R} \in \mathbb{R}^{N \times N}$ is a symmetric positive definite matrix, and $\Psi$ is a lower-bounded twice-continuously differentiable convex function. In this paper, it will be assumed that $F$ is only accessible through a sequence $(F_n)_{n \geqslant 1}$ of approximations estimated in an online manner, such that, for every $n \in \mathbb{N}^*$,

$$F_n \colon \mathbb{R}^N \to \mathbb{R} \colon \boldsymbol{h} \mapsto \frac{1}{2} \boldsymbol{h}^\top \boldsymbol{R}_n \boldsymbol{h} - \boldsymbol{r}_n^\top \boldsymbol{h} + \Psi(\boldsymbol{h}), \tag{2}$$

where the vector $\boldsymbol{r}_n$ and the symmetric nonnegative definite matrix $\boldsymbol{R}_n$ are approximations of $\boldsymbol{r}$ and $\boldsymbol{R}$. For simplicity, we will suppose that

**Assumption 1.**

(i) $(\|\boldsymbol{r}_n - \boldsymbol{r}_{n+1}\|)_{n \geqslant 1}$ and $(\|\boldsymbol{R}_n - \boldsymbol{R}_{n+1}\|)_{n \geqslant 1}$ are summable sequences,

(ii) $(\boldsymbol{r}_n)_{n \geqslant 1}$, and $(\boldsymbol{R}_n)_{n \geqslant 1}$ converge to $\boldsymbol{r}$ and $\boldsymbol{R}$, respectively.

It is worth emphasizing that Assumption 1 encompasses the batch case when $F_n \equiv F$. Moreover, it should be pointed out that all the results presented subsequently can be easily extended to a stochastic framework where $\boldsymbol{r}_n$ and $\boldsymbol{R}_n$ are consistent statistical estimates of $\boldsymbol{r}$ and $\boldsymbol{R}$, and convergence arises almost surely.

### 2.2 Majorant function

At each iteration $n \in \mathbb{N}^*$ of the MM subspace algorithm, the available estimate $F_n$ of $F$ is replaced by a surrogate function $\Theta_n(\cdot, \boldsymbol{h}_n)$ based on the current point $\boldsymbol{h}_n$ (computed at the previous iteration). This surrogate function [15–17] must be such that

$$(\forall \boldsymbol{h} \in \mathbb{R}^N) \quad F_n(\boldsymbol{h}) - F_n(\boldsymbol{h}_n) \leqslant \Theta_n(\boldsymbol{h}, \boldsymbol{h}_n) - \Theta_n(\boldsymbol{h}_n, \boldsymbol{h}_n). \tag{3}$$

We assume that $\Theta_n(\cdot, \boldsymbol{h}_n)$ is a quadratic function of the form

$$(\forall \boldsymbol{h} \in \mathbb{R}^N) \quad \Theta_n(\boldsymbol{h}, \boldsymbol{h}_n) = F_n(\boldsymbol{h}_n) + \nabla F_n(\boldsymbol{h}_n)^\top (\boldsymbol{h} - \boldsymbol{h}_n)$$
$$+ \frac{1}{2} (\boldsymbol{h} - \boldsymbol{h}_n)^\top \boldsymbol{A}_n(\boldsymbol{h}_n)(\boldsymbol{h} - \boldsymbol{h}_n), \tag{4}$$

where $\boldsymbol{A}_n(\boldsymbol{h}_n) = \boldsymbol{R}_n + \boldsymbol{B}(\boldsymbol{h}_n)$ and $\boldsymbol{B}(\boldsymbol{h}_n) \in \mathbb{R}^{N \times N}$ is some symmetric nonnegative definite matrix (see [18–22] for examples).

## 2.3 MM subspace algorithm

The MM subspace algorithm consists of defining the following sequence of vectors $(\boldsymbol{h}_n)_{n \geqslant 1}$:

$$(\forall n \in \mathbb{N}^*) \qquad \boldsymbol{h}_{n+1} \in \operatorname*{Argmin}_{\boldsymbol{h} \in \operatorname{ran} \boldsymbol{D}_n} \ \Theta_n(\boldsymbol{h}, \boldsymbol{h}_n), \tag{5}$$

where $\boldsymbol{h}_1$ is set to an initial value, and $\operatorname{ran} \boldsymbol{D}_n$ is the range of matrix $\boldsymbol{D}_n \in \mathbb{R}^{N \times M_n}$ with $M_n \geqslant 1$, constructed in such a way that the steepest descent direction $-\nabla F_n(\boldsymbol{h}_n)$ belongs to $\operatorname{ran} \boldsymbol{D}_n$. Several choices have been proposed in the literature for matrices $(\boldsymbol{D}_n)_{n \in \mathbb{N}^*}$. On the one hand, if, for every $n \in \mathbb{N}^*$, $\operatorname{rank}(\boldsymbol{D}_n) = N$, Algorithm (5) becomes equivalent to a half-quadratic method with unit stepsize [13, 23, 24]. Half-quadratic algorithms are known to be effective optimization methods, but the resolution of the minimization subproblem involved in (5) requires the inversion of matrix $\boldsymbol{A}_n(\boldsymbol{h}_n)$ which may have a high computational cost. On the other hand, if for every $n \in \mathbb{N}^*$, $\boldsymbol{D}_n$ reduces to $[-\nabla F_n(\boldsymbol{h}_n), \boldsymbol{h}_n]$, then (5) reads: for every $n \in \mathbb{N}^*$ $\boldsymbol{h}_{n+1} = u_{n,2}\boldsymbol{h}_n - u_{n,1}\nabla F_n(\boldsymbol{h}_n)$, where $(u_{n,1}, u_{n,2}) \in \mathbb{R}^2$. In the special case when $u_{n,2} = 1$, we recover the form of a gradient-like algorithm with step-size $u_{n,1}$ [25, 26]. An intermediate size subspace matrix is obtained by choosing, for every $n > 1$, $\boldsymbol{D}_n = [-\nabla F_n(\boldsymbol{h}_n), \boldsymbol{h}_n, \boldsymbol{h}_n - \boldsymbol{h}_{n-1}]$. This particular choice for the subspace yields the 3MG algorithm [8, 11].

## 2.4 Convergence result

The convergence of the MM subspace Algorithm (5) has been studied in [1, 8, 11] under various assumptions. We now provide a convergence result which is a deterministic version of the one in [11, Section IV]. This result requires the following additional assumption:

**Assumption 2.**

(i) *For every $n \in \mathbb{N}^*$, $\{\nabla F_n(\boldsymbol{h}_n), \boldsymbol{h}_n\} \subset \operatorname{ran} \boldsymbol{D}_n$,*

(ii) *There exists a positive definite matrix $\boldsymbol{V}$ such that, for every $n \in \mathbb{N}^*$, $\nabla^2 \Psi(\boldsymbol{h}_n) \preceq \boldsymbol{B}(\boldsymbol{h}_n) \preceq \boldsymbol{V}$, where $\nabla^2 \Psi$ denotes the Hessian of $\Psi$,* [1]

(iii) *At least one of the following statements holds:*

    (a) *$\boldsymbol{r}_n \equiv \boldsymbol{r}$ and $\boldsymbol{R}_n \equiv \boldsymbol{R}$,*

    (b) *$\boldsymbol{h} \mapsto \boldsymbol{B}(\boldsymbol{h})\boldsymbol{h} - \nabla \Psi(\boldsymbol{h})$ is a bounded function.*

**Remark 1.** *Note that the convexity of $\Psi$ and Assumption 2(ii) implies that $\Psi$ is Lipschitz differentiable on $\mathbb{R}^N$, with Lipschitz constant $|||\boldsymbol{V}|||$. Conversely, if $\Psi$ is $\beta$-Lipschitz differentiable with $\beta \in ]0, +\infty[$, Assumption 2(ii) is satisfied with $\boldsymbol{V} = \boldsymbol{B}(\boldsymbol{h}_n) = \beta \boldsymbol{I}_N$ [27]. However, better choices for the curvature matrix are often possible [20, 22]. In particular, Assumption 2(iii)(b), required in the online case, is satisfied for a wide class of functions and majorants [1, 11].*

**Proposition 1.** *Assume that Assumptions 1 and 2 are fulfilled. Then, the following hold:*

(i) *$(\|\nabla F_n(\boldsymbol{h}_n)\|)_{n \geqslant 1}$ is square-summable.*

(ii) *$(\boldsymbol{h}_n)_{n \geqslant 1}$ converges to the unique (global) minimizer $\widehat{\boldsymbol{h}}$ of $F$.*

*Proof.* See Appendix A. □

---

[1] $\preceq$ and $\prec$ denote the weak and strict Loewner orders, respectively,

# 3 Convergence rate analysis

## 3.1 Convergence rate results

We will first give a technical lemma the proof of which is in the spirit of classical approximation techniques for the study of first-order optimization methods (see [28, Section 1]):

**Lemma 1.** *Suppose that Assumptions 1 and 2 hold. Let $\epsilon \in ]0, +\infty[$ be such that $\epsilon \boldsymbol{I}_N \prec \boldsymbol{R}$. Then, there exists $n_\epsilon \in \mathbb{N}^*$ such that, for every $n \geqslant n_\epsilon$, $\nabla^2 F_n(\boldsymbol{h}_n) \succeq \boldsymbol{R} - \epsilon \boldsymbol{I}_N$ and*

$$F_n(\boldsymbol{h}_n) - \inf F_n \leqslant \frac{1}{2}(1 + \epsilon)\big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \big(\nabla^2 F_n(\boldsymbol{h}_n)\big)^{-1} \nabla F_n(\boldsymbol{h}_n). \tag{6}$$

*Proof.* See Appendix B. $\qquad \square$

We now state our main result which basically allows us to quantify how fast the proposed iterative approach is able to decrease asymptotically the cost function:

**Proposition 2.** *Suppose that Assumptions 1 and 2 hold. Let $\epsilon \in ]0, +\infty[$ be such that $\epsilon \boldsymbol{I}_N \prec \boldsymbol{R}$. Then, there exists $n_\epsilon \in \mathbb{N}^*$ such that, for every $n \geqslant n_\epsilon$, $\nabla^2 F_n(\boldsymbol{h}_n) \succeq \boldsymbol{R} - \epsilon \boldsymbol{I}_N$ and*

$$F_n(\boldsymbol{h}_{n+1}) - \inf F_n \leqslant \theta_n \big(F_n(\boldsymbol{h}_n) - \inf F_n\big) \tag{7}$$

*where $\theta_n = 1 - (1 + \epsilon)^{-1}\widetilde{\theta}_n$,*

$$\widetilde{\theta}_n = \frac{\big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \boldsymbol{C}_n(\boldsymbol{h}_n)\nabla F_n(\boldsymbol{h}_n)}{\big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \big(\nabla^2 F_n(\boldsymbol{h}_n)\big)^{-1}\nabla F_n(\boldsymbol{h}_n)}, \tag{8}$$

*$\boldsymbol{C}_n(\boldsymbol{h}_n) = \boldsymbol{D}_n(\boldsymbol{D}_n^\top \boldsymbol{A}_n(\boldsymbol{h}_n)\boldsymbol{D}_n)^\dagger \boldsymbol{D}_n^\top$, and $(\cdot)^\dagger$ denotes the pseudo-inverse operation. Furthermore, some lower and upper bounds on $\theta_n$ are given by*

$$\underline{\theta}_n = 1 - (1 + \epsilon)^{-1}\underline{\kappa}_n^{-1} > 0, \tag{9}$$

$$\overline{\theta}_n = 1 - (1 + \epsilon)^{-1}\overline{\kappa}_n^{-1}\left(1 - \left(\frac{\overline{\sigma}_n - \underline{\sigma}_n}{\overline{\sigma}_n + \underline{\sigma}_n}\right)^2\right) < 1, \tag{10}$$

*where $\underline{\kappa}_n \geqslant 1$ (resp. $\overline{\kappa}_n$) is the minimum (resp. maximum) eigenvalue of $\big(\boldsymbol{A}_n(\boldsymbol{h}_n)\big)^{\frac{1}{2}}\big(\nabla^2 F_n(\boldsymbol{h}_n)\big)^{-1}\big(\boldsymbol{A}_n(\boldsymbol{h}_n)\big)^{\frac{1}{2}}$, and $\underline{\sigma}_n$ (resp. $\overline{\sigma}_n$) is the minimum (resp. maximum) eigenvalue of $\nabla^2 F_n(\boldsymbol{h}_n)$.*

*Proof.* See Appendix C. $\qquad \square$

## 3.2 Discussion on the choice of the subspace

Let us make some comments about the above results. First, as enlightened by our proof, at iteration $n \geqslant n_\epsilon$, the upper value of $\theta_n$ (i.e. the slowest convergence) is obtained in the case of a gradient-like algorithm. As expected, $\overline{\theta}_n$ has a larger value when the eigenvalues of the Hessian of $F_n$ are dispersed. Note that, according to (50),

$$\frac{\overline{\sigma}_n - \underline{\sigma}_n}{\overline{\sigma}_n + \underline{\sigma}_n} \leqslant \frac{\overline{\eta} - \underline{\eta} + 2\epsilon}{\overline{\eta} + \underline{\eta}}, \tag{11}$$

where $\underline{\eta} > 0$ is the minimum eigenvalue of $\boldsymbol{R}$ and $\overline{\eta}$ is the maximum eigenvalue of $\boldsymbol{R} + \boldsymbol{V}$. Since $\big(\big(\boldsymbol{A}_n(\boldsymbol{h}_n)\big)^{\frac{1}{2}}\big(\nabla^2 F_n(\boldsymbol{h}_n)\big)^{-1}\big(\boldsymbol{A}_n(\boldsymbol{h}_n)\big)^{\frac{1}{2}}\big)_{n \geqslant n_\epsilon}$ is bounded, there exists $\overline{\kappa}_{\max} \in [1, +\infty[$

such that $(\forall n \geqslant n_\epsilon) \; \overline{\kappa}_n \leqslant \overline{\kappa}_{\max}$. All these show that the decay rate is uniformly strictly lower than 1.

In contrast, when the search subspace is the full space, the lower value of $\theta_n$ (i.e. the fastest convergence) is obtained. The expression $\underline{\theta}_n$ in (9) shows that the decay is then faster when the quadratic majorant constitutes a tight approximation of function $F_n$ at $\boldsymbol{h}_n$. Ideally, if $\boldsymbol{A}_n(\boldsymbol{h}_n)$ can be chosen equal to $\nabla^2 F_n(\boldsymbol{h}_n)$ and $\boldsymbol{D}_n$ is full rank, then $\theta_n = O(\epsilon)$. Such a behavior similar to Newton's method behavior leads to the best performance one can reasonably expect from the available data at iteration $n$.

Finally, when a mid-size subspace is chosen (as in the 3MG algorithm), an intermediate decay rate is obtained. Provided that $\boldsymbol{D}_n$ captures the main eigendirections in $\boldsymbol{A}_n(\boldsymbol{h}_n)$, a behavior close to the one previously mentioned can be expected in practice with the potential advantage of a reduced computational complexity per iteration.

## 3.3 Batch case

The case when $F \equiv F_n$ is of main interest since it is addressed in most of the existing works. Then, Proposition 2 and (11) lead to

$$(\forall n \geqslant n_\epsilon) \quad F(\boldsymbol{h}_n) - \inf F \leqslant \mu \vartheta^n, \tag{12}$$

where $\mu = \big(F(\boldsymbol{h}_{n_\epsilon}) - \inf F\big)/\vartheta^{n_\epsilon}$ and the worst-case geometrical decay rate $\vartheta \in ]0,1[$ is given by

$$\vartheta = 1 - \frac{1}{(1+\epsilon)\overline{\kappa}_{\max}} \left( 1 - \Big(\frac{\overline{\eta} - \underline{\eta} + 2\epsilon}{\overline{\eta} + \underline{\eta}}\Big)^2 \right). \tag{13}$$

Since $F$ is an $\underline{\eta}$-strongly convex function, the following inequality is satisfied [27, Definition 10.5], for every $\alpha \in ]0,1[$,

$$F\big(\alpha \boldsymbol{h}_n + (1-\alpha)\widehat{\boldsymbol{h}}\big) + \frac{1}{2}\alpha(1-\alpha)\underline{\eta}\|\boldsymbol{h}_n - \widehat{\boldsymbol{h}}\|^2 \leqslant \alpha F(\boldsymbol{h}_n) + (1-\alpha)F(\widehat{\boldsymbol{h}}), \tag{14}$$

or, equivalently,

$$\frac{1}{2}\alpha(1-\alpha)\underline{\eta}\|\boldsymbol{h}_n - \widehat{\boldsymbol{h}}\|^2 \leqslant \alpha\big(F(\boldsymbol{h}_n) - F(\widehat{\boldsymbol{h}})\big) + F(\widehat{\boldsymbol{h}}) - F\big(\alpha \boldsymbol{h}_n + (1-\alpha)\widehat{\boldsymbol{h}}\big). \tag{15}$$

Thus,

$$\frac{1}{2}(1-\alpha)\underline{\eta}\|\boldsymbol{h}_n - \widehat{\boldsymbol{h}}\|^2 \leqslant F(\boldsymbol{h}_n) - F(\widehat{\boldsymbol{h}}). \tag{16}$$

Letting $\alpha$ tend to 0 in the latter inequality implies that

$$\frac{1}{2}\underline{\eta}\|\boldsymbol{h}_n - \widehat{\boldsymbol{h}}\|^2 \leqslant F(\boldsymbol{h}_n) - F(\widehat{\boldsymbol{h}}) \leqslant \mu \vartheta^n. \tag{17}$$

This shows that the MM subspace algorithm converges linearly with rate $\sqrt{\vartheta}$.

## 4 Conclusion

In this paper, we have established expressions of the convergence rate of an online version of the MM subspace algorithm. These results help in better understanding the good numerical behaviour of this algorithm in signal/image processing applications and the role played by the subspace choice. Even in the batch case, the provided linear convergence result appears to be new. In future work, it could be interesting to investigate extensions of these properties to more general cost functions than (1).

# A  Proof of Proposition 1

## A.1  Boundedness of $(\boldsymbol{h}_n)_{n \geqslant 1}$ (online case)

Assume that Assumption 2(iii)(b) holds. For every $n \in \mathbb{N}^*$, minimizing $\Theta_n(\cdot, \boldsymbol{h}_n)$ is equivalent to minimizing the function

$$(\forall \boldsymbol{h} \in \mathbb{R}^N) \quad \widetilde{\Theta}_n(\boldsymbol{h}, \boldsymbol{h}_n) = \frac{1}{2}\boldsymbol{h}^\top \boldsymbol{A}_n(\boldsymbol{h}_n)\boldsymbol{h} - \boldsymbol{c}_n(\boldsymbol{h}_n)^\top \boldsymbol{h}, \tag{18}$$

with

$$\begin{aligned}
\boldsymbol{c}_n(\boldsymbol{h}_n) &= \boldsymbol{A}_n(\boldsymbol{h}_n)\boldsymbol{h}_n - \nabla F_n(\boldsymbol{h}_n) \\
&= \boldsymbol{r}_n + \boldsymbol{B}(\boldsymbol{h}_n)\boldsymbol{h}_n - \nabla \Psi(\boldsymbol{h}_n)
\end{aligned} \tag{19}$$

According to Assumption 2(iii)(b), these exists $\eta \in ]0, +\infty[$ such that

$$(\forall n \geqslant 1) \qquad \|\boldsymbol{c}_n(\boldsymbol{h}_n)\| \leqslant \eta, \tag{20}$$

In addition, because of Assumption 1(ii), there exists $\epsilon \in ]0, +\infty[$ and $n_0 \in \mathbb{N}^*$ such that

$$(\forall n \geqslant n_0) \qquad \boldsymbol{A}_n(\boldsymbol{h}_n) \succeq \boldsymbol{R} - \epsilon \boldsymbol{I}_N \succ \boldsymbol{O}_N, \tag{21}$$

Using now the Cauchy-Schwarz inequality, we have

$$(\forall n \geqslant n_0)(\forall \boldsymbol{h} \in \mathbb{R}^N) \quad \frac{1}{2}\boldsymbol{h}^\top(\boldsymbol{R} - \epsilon \boldsymbol{I}_N)\boldsymbol{h} - \|\boldsymbol{h}\|\eta \leqslant \widetilde{\Theta}_n(\boldsymbol{h}, \boldsymbol{h}_n). \tag{22}$$

Since $\boldsymbol{R} - \epsilon \boldsymbol{I}_N$ is a positive definite matrix, the lower bound corresponds to a coercive function with respect to $\boldsymbol{h}$. There thus exists $\zeta \in ]0, +\infty[$ such that, for every $\boldsymbol{h} \in \mathbb{R}^N$,

$$\|\boldsymbol{h}\| > \zeta \quad \Rightarrow \quad (\forall n \geqslant n_0) \ \widetilde{\Theta}_n(\boldsymbol{h}, \boldsymbol{h}_n) > 0. \tag{23}$$

On the other hand, since $\boldsymbol{0} \in \operatorname{ran} \boldsymbol{D}_n$, we have

$$\widetilde{\Theta}_n(\boldsymbol{h}_{n+1}, \boldsymbol{h}_n) \leqslant \widetilde{\Theta}_n(\boldsymbol{0}, \boldsymbol{h}_n) = 0. \tag{24}$$

The last two inequalities allow us to conclude that

$$(\forall n \geqslant n_0) \qquad \|\boldsymbol{h}_{n+1}\| \leqslant \zeta. \tag{25}$$

## A.2  Convergence of $(F_n(\boldsymbol{h}_n))_{n \geqslant 1}$

According to Assumption 2(i), the proposed algorithm is actually equivalent to

$$(\forall n \in \mathbb{N}^*) \qquad \boldsymbol{h}_{n+1} = \boldsymbol{h}_n + \boldsymbol{D}_n \widetilde{\boldsymbol{u}}_n \tag{26}$$

$$\widetilde{\boldsymbol{u}}_n = \underset{\widetilde{\boldsymbol{u}} \in \mathbb{R}^{M_n}}{\arg\min} \Theta_n(\boldsymbol{h}_n + \boldsymbol{D}_n \widetilde{\boldsymbol{u}}, \boldsymbol{h}_n). \tag{27}$$

By using (4) and cancelling the derivative of the function $\widetilde{\boldsymbol{u}} \mapsto \Theta_n(\boldsymbol{h}_n + \boldsymbol{D}_n \widetilde{\boldsymbol{u}}, \boldsymbol{h}_n)$,

$$\boldsymbol{D}_n^\top \nabla F_n(\boldsymbol{h}_n) + \boldsymbol{D}_n^\top \boldsymbol{A}_n(\boldsymbol{h}_n)\boldsymbol{D}_n \widetilde{\boldsymbol{u}}_n = \boldsymbol{0}. \tag{28}$$

Hence,

$$\begin{aligned}
\Theta(\boldsymbol{h}_{n+1}, \boldsymbol{h}_n) &= F_n(\boldsymbol{h}_n) - \frac{1}{2}\widetilde{\boldsymbol{u}}_n^\top \boldsymbol{D}_n^\top \boldsymbol{A}_n(\boldsymbol{h}_n)\boldsymbol{D}_n \widetilde{\boldsymbol{u}}_n \\
&= F_n(\boldsymbol{h}_n) - \frac{1}{2}(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)^\top \boldsymbol{A}_n(\boldsymbol{h}_n)(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n).
\end{aligned} \tag{29}$$

In view of (3) and (4), this yields

$$(\forall n \in \mathbb{N}^*) \quad F_n(\boldsymbol{h}_{n+1}) + \frac{1}{2}(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)^\top \boldsymbol{A}_n(\boldsymbol{h}_n)(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n) \leqslant F_n(\boldsymbol{h}_n). \tag{30}$$

In addition, the following recursive relation holds

$$(\forall \boldsymbol{h} \in \mathbb{R}^N) \quad F_{n+1}(\boldsymbol{h}) = F_n(\boldsymbol{h}) - (\boldsymbol{r}_{n+1} - \boldsymbol{r}_n)^\top \boldsymbol{h} + \frac{1}{2}\boldsymbol{h}^\top(\boldsymbol{R}_{n+1} - \boldsymbol{R}_n)\boldsymbol{h}. \tag{31}$$

It can thus be deduced that

$$F_{n+1}(\boldsymbol{h}_{n+1}) + \frac{1}{2}(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)^\top \boldsymbol{A}_n(\boldsymbol{h}_n)(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n) \leqslant F_n(\boldsymbol{h}_n) + \chi_n \tag{32}$$

where

$$\chi_n = -(\boldsymbol{r}_n - \boldsymbol{r}_{n+1})^\top \boldsymbol{h}_{n+1} + \frac{1}{2}\boldsymbol{h}_{n+1}^\top(\boldsymbol{R}_n - \boldsymbol{R}_{n+1})\boldsymbol{h}_{n+1}. \tag{33}$$

We have

$$|\chi_n| \leqslant \|\boldsymbol{r}_n - \boldsymbol{r}_{n+1}\| \, \|\boldsymbol{h}_{n+1}\| + \frac{1}{2}|||\boldsymbol{R}_n - \boldsymbol{R}_{n+1}||| \, \|\boldsymbol{h}_{n+1}\|^2. \tag{34}$$

If Assumption 2(iii)(b) holds, then, according to (25), $(\boldsymbol{h}_n)_{n \geqslant 1}$ is bounded, so that Assumption 1(i) guarantees that

$$\sum_{n=1}^{+\infty} |\chi_n| < +\infty. \tag{35}$$

Otherwise, if Assumption 2(iii)(a) holds, then $\chi_n \equiv 0$ and (35) is obviously fulfilled. The lower-boundedness property of $\Psi$ entails that, for every $n \in \mathbb{N}^*$, $F_n$ is lower bounded by $\inf \Psi > -\infty$. Furthermore, (32) leads to

$$F_{n+1}(\boldsymbol{h}_{n+1}) - \inf \Psi + \frac{1}{2}(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)^\top \boldsymbol{A}_n(\boldsymbol{h}_n)(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n) \leqslant F_n(\boldsymbol{h}_n) - \inf \Psi + |\chi_n|. \tag{36}$$

Since, for every $n \in \mathbb{N}^*$, $F_n(\boldsymbol{h}_n) - \inf \Psi$ and $(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)^\top \boldsymbol{A}_n(\boldsymbol{h}_n)(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)$ are nonnegative, $\big((\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)^\top \boldsymbol{A}_n(\boldsymbol{h}_n)(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)\big)_{n \geqslant 1}$ is a summable sequence, and $(F_n(\boldsymbol{h}_n))_{n \geqslant 1}$ is convergent.

## A.3   Convergence of $(\nabla \mathrm{F}_n(\boldsymbol{h}_n))_{n \geqslant 1}$

According to (4), we have, for every $\phi \in \mathbb{R}$ and $n \in \mathbb{N}^*$,

$$\Theta_n\big(\boldsymbol{h}_n - \phi \nabla F_n(\boldsymbol{h}_n), \boldsymbol{h}_n\big) = F_n(\boldsymbol{h}_n) - \phi\|\nabla F_n(\boldsymbol{h}_n)\|^2 + \frac{\phi^2}{2}\big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \boldsymbol{A}_n(\boldsymbol{h}_n)\nabla F_n(\boldsymbol{h}_n). \tag{37}$$

Let

$$\Phi_n \in \underset{\phi \in \mathbb{R}}{\mathrm{Argmin}} \ \Theta_n\big(\boldsymbol{h}_n - \phi \nabla F_n(\boldsymbol{h}_n), \boldsymbol{h}_n\big). \tag{38}$$

The following optimality condition holds:

$$\big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \boldsymbol{A}_n(\boldsymbol{h}_n)\nabla F_n(\boldsymbol{h}_n)\, \Phi_n = \|\nabla F_n(\boldsymbol{h}_n)\|^2. \tag{39}$$

As a consequence of Assumption 2(i), $(\forall \phi \in \mathbb{R}) \ \boldsymbol{h}_n - \phi \nabla F_n(\boldsymbol{h}_n) \in \mathrm{ran}\, \boldsymbol{D}_n$. It then follows from (5) and (39) that

$$\Theta_n\big(\boldsymbol{h}_{n+1}, \boldsymbol{h}_n\big) \leqslant \Theta_n\big(\boldsymbol{h}_n - \Phi_n \nabla F_n(\boldsymbol{h}_n), \boldsymbol{h}_n\big) = F_n(\boldsymbol{h}_n) - \frac{\Phi_n}{2}\|\nabla F_n(\boldsymbol{h}_n)\|^2, \tag{40}$$

7

which, by using (29), leads to

$$\Phi_n \|\nabla F_n(\boldsymbol{h}_n)\|^2 \leqslant (\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)^\top \boldsymbol{A}_n(\boldsymbol{h}_n)(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n). \tag{41}$$

Let $\epsilon > 0$. Assumption 2(ii) yields, for every $n \in \mathbb{N}^*$,

$$\boldsymbol{A}_n(\boldsymbol{h}_n) \preceq (|||\boldsymbol{R}_n||| + |||\boldsymbol{V}|||)\boldsymbol{I}_N. \tag{42}$$

Therefore, according to Assumption 1(ii),

$$(\exists n_0 \in \mathbb{N}^*)(\forall n \geqslant n_0) \quad \boldsymbol{O}_N \prec \boldsymbol{A}_n(\boldsymbol{h}_n) \preceq \alpha_\epsilon^{-1} \boldsymbol{I}_N \tag{43}$$

where

$$\alpha_\epsilon = (|||\boldsymbol{R}||| + |||\boldsymbol{V}||| + \epsilon)^{-1} > 0. \tag{44}$$

By using now (39), it can be deduced from (43) that, if $n \geqslant n_0$ and $\nabla F_n(\boldsymbol{h}_n) \neq \boldsymbol{0}$, then $\Phi_n \geqslant \alpha_\epsilon$. Then, it follows from (41) that

$$\alpha_\epsilon \sum_{n=n_0}^{+\infty} \|\nabla F_n(\boldsymbol{h}_n)\|^2 \leqslant \sum_{n=n_0}^{+\infty} (\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)^\top \boldsymbol{A}_n(\boldsymbol{h}_n)(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n). \tag{45}$$

By invoking the summability property of $\big((\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)^\top \boldsymbol{A}_n(\boldsymbol{h}_n)(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)\big)_{n \geqslant 1}$, we can conclude that $(\|\nabla F_n(\boldsymbol{h}_n)\|^2)_{n \geqslant 1}$ is itself summable.

## A.4   Convergence of $(\boldsymbol{h}_n)_{n \geqslant 1}$

We have shown that $\big((\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)^\top \boldsymbol{A}_n(\boldsymbol{h}_n)(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)\big)_{n \geqslant 1}$ converges to 0. In addition, we have seen that (21) holds for a given $\epsilon \in ]0, +\infty[$ and $n_0 \in \mathbb{N}^*$. This implies that, for every $n \geqslant n_0$,

$$|||\boldsymbol{R} - \epsilon \boldsymbol{I}_N||| \, \|\boldsymbol{h}_{n+1} - \boldsymbol{h}_n\|^2 \leqslant (\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)^\top \boldsymbol{A}_n(\boldsymbol{h}_n)(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n) \tag{46}$$

where $|||\boldsymbol{R} - \epsilon \boldsymbol{I}_N||| > 0$. Consequently, $(\boldsymbol{h}_{n+1} - \boldsymbol{h}_n)_{n \geqslant 1}$ converges to $\boldsymbol{0}$. In addition, $(\boldsymbol{h}_n)_{n \geqslant 1}$ belongs to a compact set. Thus, invoking Ostrowski's theorem [29, Theorem 26.1] implies that the set of cluster points of $(\boldsymbol{h}_n)_{n \geqslant 1}$ is a nonempty compact connected set. By using (1)-(2), we have

$$(\forall n \in \mathbb{N}^*) \quad \nabla F_n(\boldsymbol{h}_n) - \nabla F(\boldsymbol{h}_n) = (\boldsymbol{R}_n - \boldsymbol{R})\boldsymbol{h}_n - \boldsymbol{r}_n + \boldsymbol{r}. \tag{47}$$

Since $(\boldsymbol{h}_n)_{n \geqslant 1}$ is bounded, it follows from that $\big(\nabla F_n(\boldsymbol{h}_n) - \nabla F(\boldsymbol{h}_n)\big)_{n \geqslant 1}$ converges to $\boldsymbol{0}$. Since $\big(\nabla F_n(\boldsymbol{h}_n)\big)_{n \geqslant 1}$ converges to $\boldsymbol{0}$, this implies that $\big(\nabla F(\boldsymbol{h}_n)\big)_{n \geqslant 1}$ also converges to $\boldsymbol{0}$. Let $\widehat{\boldsymbol{h}}$ be a cluster point of $(\boldsymbol{h}_n)_{n \geqslant 1}$. There exists a subsequence $(\boldsymbol{h}_{k_n})_{n \geqslant 1}$ such that $\boldsymbol{h}_{k_n} \to \widehat{\boldsymbol{h}}$. As $F$ is continuously differentiable, we have

$$\nabla F(\widehat{\boldsymbol{h}}) = \lim_{n \to +\infty} \nabla F(\boldsymbol{h}_{k_n}) = \boldsymbol{0}. \tag{48}$$

This means that $\widehat{\boldsymbol{h}}$ is a critical point of $F$. Since $F$ is a strongly convex function, it possesses a unique critical point $\widehat{\boldsymbol{h}}$, which is the global minimizer of $F$ [27, Prop.11.7]. Since the unique cluster point of $(\boldsymbol{h}_n)_{n \geqslant 1}$ is $\widehat{\boldsymbol{h}}$, this shows that $\boldsymbol{h}_n \to \widehat{\boldsymbol{h}}$.

# B Proof of Lemma 1

Because $\boldsymbol{R}$ is positive definite, according to Assumption 1(ii), there exists $n_0 \in \mathbb{N}^*$ such that, for every $n \geqslant n_0$,

$$\boldsymbol{O}_N \prec \boldsymbol{R} - \epsilon\boldsymbol{I}_N \preceq \boldsymbol{R}_n \preceq \boldsymbol{R} + \epsilon\boldsymbol{I}_N. \tag{49}$$

Let $n \geqslant n_0$. Then, $F_n$ is a strongly convex continuous function. From standard results, this function possesses a unique global minimizer $\widehat{\boldsymbol{h}}_n$. According to Assumption 2(ii), and (49), $\nabla^2 F_n$ is such that

$$\begin{aligned}(\forall \boldsymbol{h} \in \mathbb{R}^N) \quad \boldsymbol{O}_N &\prec \boldsymbol{R} - \epsilon\boldsymbol{I}_N \\ &\preceq \boldsymbol{R}_n + \nabla^{(2)}\Psi(\boldsymbol{h}) = \nabla^2 F_n(\boldsymbol{h}) \\ &\preceq \boldsymbol{R} + \epsilon\boldsymbol{I}_N + \boldsymbol{V}.\end{aligned} \tag{50}$$

By using now the second-order Taylor formula with integral remainder, we get

$$F_n(\widehat{\boldsymbol{h}}_n) = F_n(\boldsymbol{h}_n) + \big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \big(\widehat{\boldsymbol{h}}_n - \boldsymbol{h}_n\big) + \frac{1}{2}\big(\widehat{\boldsymbol{h}}_n - \boldsymbol{h}_n\big)^\top \boldsymbol{H}_n^{(2)}(\boldsymbol{h}_n)\big(\widehat{\boldsymbol{h}}_n - \boldsymbol{h}_n\big), \tag{51}$$

where

$$\begin{aligned}\nabla F_n(\boldsymbol{h}_n) &= \nabla F_n(\widehat{\boldsymbol{h}}_n) + \boldsymbol{H}_n^{(1)}(\boldsymbol{h}_n)(\boldsymbol{h}_n - \widehat{\boldsymbol{h}}_n) \\ &= \boldsymbol{H}_n^{(1)}(\boldsymbol{h}_n)(\boldsymbol{h}_n - \widehat{\boldsymbol{h}}_n)\end{aligned} \tag{52}$$

and, for every $\boldsymbol{h} \in \mathbb{R}^N$,

$$\begin{aligned}\boldsymbol{H}_n^{(1)}(\boldsymbol{h}) &= \int_0^1 \nabla^2 F_n\big(\widehat{\boldsymbol{h}}_n + t(\boldsymbol{h} - \widehat{\boldsymbol{h}}_n)\big)dt \\ &= \boldsymbol{R}_n + \int_0^1 \nabla^2\Psi\big(\widehat{\boldsymbol{h}}_n + t(\boldsymbol{h} - \widehat{\boldsymbol{h}}_n)\big)dt \end{aligned} \tag{53}$$

$$\begin{aligned}\boldsymbol{H}_n^{(2)}(\boldsymbol{h}) &= 2\int_0^1 (1-t)\nabla^2 F_n\big(\widehat{\boldsymbol{h}}_n + t(\boldsymbol{h} - \widehat{\boldsymbol{h}}_n)\big)dt \\ &= \boldsymbol{R}_n + 2\int_0^1 (1-t)\nabla^2\Psi\big(\widehat{\boldsymbol{h}}_n + t(\boldsymbol{h} - \widehat{\boldsymbol{h}}_n)\big)dt. \end{aligned} \tag{54}$$

Because of the lower bound in (50),

$$(\forall \boldsymbol{h} \in \mathbb{R}^N) \quad \boldsymbol{O}_N \prec \boldsymbol{R} - \epsilon\boldsymbol{I}_N \preceq \boldsymbol{H}_n^{(1)}(\boldsymbol{h}) \tag{55}$$

and $\boldsymbol{H}_n^{(1)}(\boldsymbol{h})$ is thus invertible. Therefore, combining (51) and (52) yields

$$\begin{aligned}F_n(\widehat{\boldsymbol{h}}_n) = F_n(\boldsymbol{h}_n) &- \big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \big(\boldsymbol{H}_n^{(1)}(\boldsymbol{h}_n)\big)^{-1}\nabla F_n(\boldsymbol{h}_n) \\ &+ \frac{1}{2}\big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \big(\boldsymbol{H}_n^{(1)}(\boldsymbol{h}_n)\big)^{-1}\boldsymbol{H}_n^{(2)}(\boldsymbol{h}_n)\big(\boldsymbol{H}_n^{(1)}(\boldsymbol{h}_n)\big)^{-1}\nabla F_n(\boldsymbol{h}_n). \end{aligned} \tag{56}$$

According to Assumption 2(ii), for every $t \in [0,1]$,

$$|||\nabla^2\Psi\big(\widehat{\boldsymbol{h}}_n + t(\boldsymbol{h}_n - \widehat{\boldsymbol{h}}_n)\big)||| \leqslant |||\boldsymbol{V}|||, \tag{57}$$

where $|||\cdot|||$ denotes the matrix spectral norm. As Proposition 1(ii) guarantees that $(\boldsymbol{h}_n)_{n\geqslant 1}$ converges to the unique minimizer $\widehat{\boldsymbol{h}}$ of $F$, it follows from Proposition 1(i), (52), and (55) that $(\widehat{\boldsymbol{h}}_n)_{n\geqslant 1}$ also converges to $\widehat{\boldsymbol{h}}$. By using the continuity of $\nabla^2\Psi$, $\big(\nabla^2\Psi\big(\widehat{\boldsymbol{h}}_n + t(\boldsymbol{h}_n -$

$\widehat{h}_n)))_{n \geqslant 1}$ converges to $\nabla^2 \Psi(\widehat{h})$ and, by invoking the dominated convergence theorem, it can be deduced that

$$\int_0^1 \nabla^2 \Psi\big(\widehat{h}_n + t(h_n - \widehat{h}_n)\big) dt \to \nabla^2 \Psi(\widehat{h}). \tag{58}$$

Since $(R_n)_{n \geqslant 1}$ converges to $R$, this allows us to conclude that $\big(H_n^{(1)}(h_n)\big)_{n \geqslant 1}$ converges to $\nabla^2 F(\widehat{h})$. Proceeding similarly, it can be proved that $\big(H_n^{(2)}(h_n)\big)_{n \geqslant 1}$ also converges to $\nabla^2 F(\widehat{h})$. This entails that

$$\big(H_n^{(1)}(h_n)\big)^{-1} - \frac{1}{2}\big(H_n^{(1)}(h_n)\big)^{-1} H_n^{(2)}(h_n)\big(H_n^{(1)}(h_n)\big)^{-1} \to \frac{1}{2}\big(\nabla^2 F(\widehat{h})\big)^{-1}. \tag{59}$$

Besides, since $\big(\nabla^2 F_n(h_n)\big)_{n \geqslant 1} = \big(R_n + \nabla^2 \Psi(h_n)\big)_{n \geqslant 1}$ converges to $\nabla^2 F(\widehat{h})$, there exists $n_\epsilon \geqslant n_0$ such that, for every $n \geqslant n_\epsilon$,

$$\big(H_n^{(1)}(h_n)\big)^{-1} - \frac{1}{2}\big(H_n^{(1)}(h_n)\big)^{-1} H_n^{(2)}(h_n)\big(H_n^{(1)}(h_n)\big)^{-1} - \frac{1}{2}\big(\nabla^2 F_n(h_n)\big)^{-1} \tag{60}$$

$$\preceq \frac{1}{2}\epsilon(R + \epsilon I_N + V)^{-1}$$

$$\preceq \frac{1}{2}\epsilon\big(\nabla^2 F_n(h_n)\big)^{-1}, \tag{61}$$

where the last inequality follows from (50). This implies that

$$\big(H_n^{(1)}(h_n)\big)^{-1} - \frac{1}{2}\big(H_n^{(1)}(h_n)\big)^{-1} H_n^{(2)}(h_n)\big(H_n^{(1)}(h_n)\big)^{-1} \preceq \frac{1}{2}(1 + \epsilon)\big(\nabla^2 F_n(h_n)\big)^{-1}. \tag{62}$$

By coming back to (56), we deduce that, for every $n \geqslant n_\epsilon$, (6) holds.

## C  Proof of Proposition 2

Let $n \in \mathbb{N}^*$. If $\nabla F_n(h_n)$ is zero, then $h_n$ is a global minimizer of $F_n$ and, according to (3)-(5), $F(h_{n+1}) \leqslant \Theta_n(h_{n+1}, h_n) - \Theta_n(h_n, h_n) + F(h_n) \leqslant F(h_n)$ so that $h_{n+1}$ is also a global minimizer of $F_n$, and (7) is obviously satisfied. So, without loss of generality, it will be assumed in the rest of the proof that $\nabla F_n(h_n)$ is nonzero. Because of Assumption 2(ii) and (49), there exists $n_0 \in \mathbb{N}^*$ such that, for every $n \geqslant n_0$,

$$O_N \prec R - \epsilon I_N \preceq R_n \preceq A_n(h_n). \tag{63}$$

Using (30) and the definition of $C_n(h_n)$,

$$F_n(h_{n+1}) \leqslant F_n(h_n) - \frac{1}{2}(h_{n+1} - h_n)^\top A_n(h_n)(h_{n+1} - h_n)$$

$$= F_n(h_n) - \frac{1}{2}\big(\nabla F_n(h_n)\big)^\top C_n(h_n)\nabla F_n(h_n). \tag{64}$$

Combining (63), (64) and (40) yields

$$\frac{\|\nabla F_n(h_n)\|^4}{\big(\nabla F_n(h_n)\big)^\top A_n(h_n)\nabla F_n(h_n)} \leqslant \big(\nabla F_n(h_n)\big)^\top C_n(h_n)\nabla F_n(h_n). \tag{65}$$

In turn, we have

$$\Theta_n\big(\widetilde{h}_n, h_n\big) \leqslant \Theta_n\big(h_{n+1}, h_n\big), \tag{66}$$

10

where $\widetilde{\boldsymbol{h}}_n$ is a global minimizer of $\Theta_n(\cdot, \boldsymbol{h}_n)$. If $n \geqslant n_0$, then (63) shows that $\boldsymbol{A}_n(\boldsymbol{h}_n)$ is invertible, and

$$\widetilde{\boldsymbol{h}}_n = \boldsymbol{h}_n - \big(\boldsymbol{A}_n(\boldsymbol{h}_n)\big)^{-1} \nabla F_n(\boldsymbol{h}_n) \tag{67}$$

which, by using (64) and (66), yields

$$\big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \boldsymbol{C}_n(\boldsymbol{h}_n) \nabla F_n(\boldsymbol{h}_n) \leqslant \big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \big(\boldsymbol{A}_n(\boldsymbol{h}_n)\big)^{-1} \nabla F_n(\boldsymbol{h}_n). \tag{68}$$

It can be noticed that the lower bound in (65) is obtained when $\boldsymbol{D}_n = \nabla F_n(\boldsymbol{h}_n)$, while the upper bound in (68) is attained when $M_n = N$ and $\boldsymbol{D}_n$ is full rank.

Let us now apply Lemma 1. According to this lemma, there exists $n_\epsilon \geqslant n_0$ such that, for every $n \geqslant n_\epsilon$, (6) holds with $\nabla^2 F_n(\boldsymbol{h}_n) \succ \boldsymbol{O}_N$. Let us assume that $n \geqslant n_\epsilon$. By combining (6) and (64), we obtain

$$F_n(\boldsymbol{h}_n) - F_n(\boldsymbol{h}_{n+1}) \geqslant \frac{\widetilde{\theta}_n}{1+\epsilon} \big(F_n(\boldsymbol{h}_n) - \inf F_n\big)$$

$$\Leftrightarrow\; F_n(\boldsymbol{h}_{n+1}) - \inf F_n \leqslant \Big(1 - \frac{\widetilde{\theta}_n}{1+\epsilon}\Big)\big(F_n(\boldsymbol{h}_n) - \inf F_n\big), \tag{69}$$

which itself is equivalent to (7). The following lower bound is then be deduced from (65):

$$\widetilde{\theta}_n \geqslant \frac{\|\nabla F_n(\boldsymbol{h}_n)\|^4}{\beta_n \big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \boldsymbol{A}_n(\boldsymbol{h}_n) \nabla F_n(\boldsymbol{h}_n)}, \tag{70}$$

by setting $\beta_n = \big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \big(\nabla^2 F_n(\boldsymbol{h}_n)\big)^{-1} \nabla F_n(\boldsymbol{h}_n)$. Hence, we have

$$\widetilde{\theta}_n \geqslant \frac{\|\nabla F_n(\boldsymbol{h}_n)\|^4}{\beta_n \beta_n'} \frac{\big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \nabla^2 F_n(\boldsymbol{h}_n) \nabla F_n(\boldsymbol{h}_n)}{\nabla F_n(\boldsymbol{h}_n)^\top \boldsymbol{A}_n(\boldsymbol{h}_n) \nabla F_n(\boldsymbol{h}_n)},$$

$$\geqslant \frac{\|\nabla F_n(\boldsymbol{h}_n)\|^4}{\beta_n \beta_n'} \Bigg( \sup_{\substack{\boldsymbol{g} \in \mathbb{R}^N \\ \boldsymbol{g} \neq \boldsymbol{0}}} \frac{\boldsymbol{g}^\top \boldsymbol{A}_n(\boldsymbol{h}_n) \boldsymbol{g}}{\boldsymbol{g}^\top \nabla^2 F_n(\boldsymbol{h}_n) \boldsymbol{g}} \Bigg)^{-1}, \tag{71}$$

where $\beta_n' = \big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \nabla^2 F_n(\boldsymbol{h}_n) \nabla F_n(\boldsymbol{h}_n)$. The sup term in (71) corresponds to the generalized Rayleigh quotient of $\boldsymbol{A}_n(\boldsymbol{h}_n)$ and $\nabla^2 F_n(\boldsymbol{h}_n)$, which is equal to $\overline{\kappa}_n$. By invoking now Kantorovich inequality [28, Section 1.3.2], we get

$$\widetilde{\theta}_n \geqslant \frac{4 \underline{\sigma}_n \overline{\sigma}_n}{\overline{\kappa}_n (\underline{\sigma}_n + \overline{\sigma}_n)^2}, \tag{72}$$

which leads to

$$1 - \frac{\widetilde{\theta}_n}{1+\epsilon} \leqslant \overline{\theta}_n < 1 \tag{73}$$

since $\overline{\sigma}_n \geqslant \underline{\sigma}_n > 0$. An upper bound on $\widetilde{\theta}_n$ is derived from (68) and (8):

$$\widetilde{\theta}_n \leqslant \frac{\big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \big(\boldsymbol{A}_n(\boldsymbol{h}_n)\big)^{-1} \nabla F_n(\boldsymbol{h}_n)}{\big(\nabla F_n(\boldsymbol{h}_n)\big)^\top \big(\nabla^2 F_n(\boldsymbol{h}_n)\big)^{-1} \nabla F_n(\boldsymbol{h}_n)}$$

$$\leqslant \sup_{\substack{\boldsymbol{g} \in \mathbb{R}^N \\ \boldsymbol{g} \neq \boldsymbol{0}}} \frac{\boldsymbol{g}^\top \big(\boldsymbol{A}_n(\boldsymbol{h}_n)\big)^{-1} \boldsymbol{g}}{\boldsymbol{g}^\top \big(\nabla^2 F_n(\boldsymbol{h}_n)\big)^{-1} \boldsymbol{g}}. \tag{74}$$

The sup term in (74) is equal to $\underline{\kappa}_n^{-1}$. Altogether (69), (73), and (74) yield (7)-(10), by setting $\theta_n = 1 - (1 + \epsilon)^{-1}\widetilde{\theta}_n$. In view of Assumption 2(ii) and the equality in (50), the Hessian of $F_n$ is such that

$$(\forall \boldsymbol{h} \in \mathbb{R}^N) \quad \nabla^2 F_n(\boldsymbol{h}) \preceq \boldsymbol{A}_n(\boldsymbol{h}), \tag{75}$$

and therefore $\underline{\kappa}_n \geqslant 1$.

# References

[1] E. Chouzenoux, J. Idier, and S. Moussaoui, "A majorize-minimize subspace strategy for subspace optimization applied to image restoration," *IEEE Trans. Image Process.*, vol. 20, no. 18, pp. 1517–1528, Jun. 2011.

[2] D. R. Hunter and K. Lange, "A tutorial on MM algorithms," *Amer. Stat.*, vol. 58, no. 1, pp. 30–37, Feb. 2004.

[3] M. Elad, B. Matalon, and M. Zibulevsky, "Coordinate and subspace optimization methods for linear least squares with non-quadratic regularization," *Appl. Comput. Harmon. Anal.*, vol. 23, pp. 346–367, Nov. 2007.

[4] A. R. Conn, N. Gould, A. Sartenaer, and Ph. L. Toint, "On iterated-subspace minimization methods for nonlinear optimization," Tech. Rep. 94-069, Rutherford Appleton Laboratory, Oxfordshire, UK, May 1994, ftp://130.246.8.32/pub/reports/cgstRAL94069.ps.Z.

[5] Y. Yuan, "Subspace techniques for nonlinear optimization," in *Some Topics in Industrial and Applied Mathematics*, R. Jeltsh, T.-T. Li, and H I. Sloan, Eds., vol. 8, pp. 206–218. Series on Concrete and Applicable Mathematics, 2007.

[6] W. W. Hager and H. Zhang, "A survey of nonlinear conjugate gradient methods," *Pac. J. Optim.*, vol. 2, no. 1, pp. 35–58, Jan. 2006.

[7] D. C. Liu and J. Nocedal, "On the limited memory BFGS method for large scale optimization," *Math. Program.*, vol. 45, no. 3, pp. 503–528, Aug. 1989.

[8] E. Chouzenoux, A. Jezierska, J.-C. Pesquet, and H. Talbot, "A majorize-minimize subspace approach for $\ell_2$-$\ell_0$ image regularization," *SIAM J. Imag. Sci.*, vol. 6, no. 1, pp. 563–591, 2013.

[9] E. Chouzenoux, J.-C. Pesquet, H. Talbot, and A. Jezierska, "A memory gradient algorithm for $\ell_2$-$\ell_0$ regularization with applications to image restoration," in *18th IEEE Int. Conf. Image Process. (ICIP 2011)*, Brussels, Belgium, 11-14 Sep. 2011, pp. 2717–2720.

[10] A. Florescu, E. Chouzenoux, J.-C. Pesquet, P. Ciuciu, and S. Ciochina, "A majorize-minimize memory gradient method for complex-valued inverse problem," *Signal Process.*, vol. 103, pp. 285–295, Oct. 2014, Special issue on Image Restoration and Enhancement: Recent Advances and Applications.

[11] E. Chouzenoux and J.-C. Pesquet, "A stochastic majorize-minimize subspace algorithm for online penalized least squares estimation," Tech. Rep., 2015, http://arxiv.org/abs/1512.08722.

[12] A. Miele and J. W. Cantrell, "Study on a memory gradient method for the minimization of functions," *J. Optim. Theory Appl.*, vol. 3, no. 6, pp. 459–470, Nov. 1969.

[13] M. Allain, J. Idier, and Y. Goussard, "On global and local convergence of half-quadratic algorithms," *IEEE Trans. Image Process.*, vol. 15, no. 5, pp. 1130–1142, May 2006.

[14] M. Nikolova and M. Ng, "Analysis of half-quadratic minimization methods for signal and image recovery," *SIAM J. Sci. Comput.*, vol. 27, no. 3, pp. 937–966, 2005.

[15] M. W. Jacobson and J. A. Fessler, "An expanded theoretical treatment of iteration-dependent Majorize-Minimize algorithms," *IEEE Trans. Image Process.*, vol. 16, no. 10, pp. 2411–2422, Oct. 2007.

[16] Z. Zhang, J. T. Kwok, and D.-Y. Yeung, "Surrogate maximization/minimization algorithms and extensions," *Mach. Learn.*, vol. 69, pp. 1–33, Oct. 2007.

[17] M. Hong, M. Razaviyayn, Z. Q. Luo, and J. S. Pang, "A unified algorithmic framework for block-structured optimization involving big data: With applications in machine learning and signal processing," *IEEE Signal Process. Mag.*, vol. 33, no. 1, pp. 57–77, Jan. 2016.

[18] M. Figueiredo, J. Bioucas-Dias, and R. Nowak, "Majorization-minimization algorithms for wavelet-based image restoration," *IEEE Trans. Image Process.*, vol. 16, no. 12, pp. 2980–2991, Dec. 2007.

[19] J.A. Fessler and H. Erdogan, "A paraboloidal surrogates algorithm for convergent penalized-likelihood emission image reconstruction," Toronto, Canada, 8-14 Nov. 1998, vol. 2, pp. 1132–1135.

[20] A. Repetti, M. Q. Pham, L. Duval, E. Chouzenoux, and J.-C. Pesquet, "Euclid in a taxicab: Sparse blind deconvolution with smoothed l1/l2 regularization," *IEEE Signal Process. Letters*, vol. 22, no. 5, pp. 539–543, May 2015.

[21] X. Ning, I. W. Selesnick, and L. Duval, "Chromatogram baseline estimation and denoising using sparsity (beads)," *Chemometr. Intell. Lab. Syst.*, vol. 139, pp. 156–167, 2014.

[22] J. Song, P. Babu, and D. P. Palomar, "Sparse generalized eigenvalue problem via smooth optimization," *IEEE Trans. Signal Process.*, vol. 63, no. 7, pp. 1627–1642, Apr. 2015.

[23] J. Idier, "Convex half-quadratic criteria and interacting auxiliary variables for image restoration," *IEEE Trans. Image Process.*, vol. 10, no. 7, pp. 1001–1009, Jul. 2001.

[24] P. Charbonnier, L. Blanc-Féraud, G. Aubert, and M. Barlaud, "Deterministic edge-preserving regularization in computed imaging," *IEEE Trans. Image Process.*, vol. 6, no. 2, pp. 298–311, Feb. 1997.

[25] C. Labat and J. Idier, "Convergence of conjugate gradient methods with a closed-form stepsize formula," *J. Optim. Theory Appl.*, vol. 136, no. 1, pp. 43–60, Jan. 2008.

[26] K. Lange, "A gradient algorithm locally equivalent to the EM algorithm," *J. R. Stat. Soc. Series B Stat. Methodol.*, vol. 57, no. 2, pp. 425–437, 1995.

[27] H. H. Bauschke and P. L. Combettes, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, Springer, New York, 2011.

[28] E. Polak, *Optimization. Algorithms and Consistent Approximations*, Springer-Verlag, New York, 1997.

[29] A. M. Ostrowski, *Solution of Equations in Euclidean and Banach Spaces*, Academic Press, London, 1973.