

# Local Linear Convergence Analysis of Primal–Dual Splitting Methods

Jingwei Liang<sup>a</sup> and Jalal Fadili<sup>b</sup> and Gabriel Peyré<sup>c</sup>

<sup>a</sup>DAMTP, University of Cambridge, UK; <sup>b</sup>Normandie Université, ENSICAEN, CNRS, GREYC, France; <sup>c</sup>CNRS, DMA, ENS Paris, France

## ARTICLE HISTORY

Compiled January 8, 2018

## ABSTRACT

In this paper, we study the local linear convergence properties of a versatile class of Primal–Dual splitting methods for minimizing composite non-smooth convex optimization problems. Under the assumption that the non-smooth components of the problem are partly smooth relative to smooth manifolds, we present a unified local convergence analysis framework for these methods. More precisely, in our framework we first show that (i) the sequences generated by Primal–Dual splitting methods identify a pair of primal and dual smooth manifolds in a finite number of iterations, and then (ii) enter a local linear convergence regime, which is characterized based on the structure of the underlying active smooth manifolds. We also show how our results for Primal–Dual splitting can be specialized to cover existing ones on Forward–Backward splitting and Douglas–Rachford splitting/ADMM (alternating direction methods of multipliers). Moreover, based on these obtained local convergence analysis result, several practical acceleration techniques are discussed. To exemplify the usefulness of the obtained result, we consider several concrete numerical experiments arising from fields including signal/image processing, inverse problems and machine learning, etc. The demonstration not only verifies the local linear convergence behaviour of Primal–Dual splitting methods, but also the insights on how to accelerate them in practice.

## KEYWORDS

Primal–Dual splitting, Forward–Backward splitting, Douglas–Rachford/ADMM Partial Smoothness, Local Linear Convergence.

## AMS CLASSIFICATION

49J52, 65K05, 65K10.

## 1. Introduction

### 1.1. *Composed optimization problem*

In various fields such as inverse problems, signal and image processing, statistics and machine learning *etc.*, many problems are (eventually) formulated as structured optimization problems (see Section 6 for some specific examples). A typical example of

---

CONTACT Jingwei Liang. Email: jl993@cam.ac.uk. Most of this work was done while Jingwei Liang was at Normandie Université, ENSICAEN, France.

these optimization problems, given in its primal form, reads

$$\min_{x \in \mathbb{R}^n} R(x) + F(x) + (J \blacktriangledown G)(Lx), \quad (\mathcal{P}_P)$$

where  $(J \blacktriangledown G)(\cdot) \stackrel{\text{def}}{=} \inf_{v \in \mathbb{R}^m} J(\cdot) + G(\cdot - v)$  denotes the infimal convolution of  $J$  and  $G$ . Throughout, we assume the following:

- (A.1)  $R, F \in \Gamma_0(\mathbb{R}^n)$  with  $\Gamma_0(\mathbb{R}^n)$  being the class of proper convex and lower semi-continuous functions on  $\mathbb{R}^n$ , and  $\nabla F$  is  $(1/\beta_F)$ -Lipschitz continuous for some  $\beta_F > 0$ ,
- (A.2)  $J, G \in \Gamma_0(\mathbb{R}^m)$ , and  $G$  is  $\beta_G$ -strongly convex for some  $\beta_G > 0$ .
- (A.3)  $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is a linear mapping.
- (A.4)  $0 \in \text{ran}(\partial R + \nabla F + L^*(\partial J \square \partial G)L)$ , where  $\partial J \square \partial G \stackrel{\text{def}}{=} (\partial J^{-1} + \partial G^{-1})^{-1}$  is the parallel sum of the subdifferential  $\partial J$  and  $\partial G$ , and  $\text{ran}(\cdot)$  denotes the range of a set-valued operator. See Remark 3 for the reasoning of this condition.

The main difficulties of solving such a problem are that the objective function is non-smooth, the presence of the linear operator  $L$  and the infimal convolution. Consider also the Fenchel-Rockafellar dual problem [41] of  $(\mathcal{P}_P)$ ,

$$\min_{v \in \mathbb{R}^m} J^*(v) + G^*(v) + (R^* \blacktriangledown F^*)(-L^*v). \quad (\mathcal{P}_D)$$

The classical Kuhn-Tucker theory asserts that a pair  $(x^*, v^*) \in \mathbb{R}^n \times \mathbb{R}^m$  solves  $(\mathcal{P}_P)$ - $(\mathcal{P}_D)$  if it satisfies the monotone inclusion

$$0 \in \begin{bmatrix} \partial R & L^* \\ -L & \partial J^* \end{bmatrix} \begin{pmatrix} x^* \\ v^* \end{pmatrix} + \begin{bmatrix} \nabla F & 0 \\ 0 & \nabla G^* \end{bmatrix} \begin{pmatrix} x^* \\ v^* \end{pmatrix}. \quad (1)$$

One can observe that in (1), the composition of the linear operator and the infimal convolution is decoupled, hence providing possibilities to achieve full splitting. This is a key property used by all Primal-Dual algorithms we are about to review. In turn, solving (1) provides a pair of points that are solutions to  $(\mathcal{P}_P)$  and  $(\mathcal{P}_D)$  respectively.

More complex forms of  $(\mathcal{P}_P)$  involving for instance a sum of infimal convolutions can be tackled in a similar way using a product space trick, as we will see in Section 5.

## 1.2. Primal-Dual splitting methods

Primal-Dual splitting methods to solve more or less complex variants of  $(\mathcal{P}_P)$ - $(\mathcal{P}_D)$  have witnessed a recent wave of interest in the literature [12,13,16,18,21,28,48]. All these methods achieve full splitting, they involve the resolvents of  $R$  and  $J^*$ , the gradients of  $F$  and  $G^*$  and the linear operator  $L$ , all separately at various points in the course of iteration. For instance, building on the work of [3], the now-popular scheme proposed in [13] solves  $(\mathcal{P}_P)$ - $(\mathcal{P}_D)$  with  $F = G^* = 0$ . The authors in [28] have shown that the Primal-Dual splitting method of [13] can be seen as a proximal point algorithm (PPA) in  $\mathbb{R}^n \times \mathbb{R}^m$  endowed with a suitable norm. Exploiting the same idea, the author in [21] considered  $(\mathcal{P}_P)$  with  $G^* = 0$ , and proposed an iterative scheme which can be interpreted as a Forward-Backward (FB) splitting again under an appropriately renormed space. This idea is further extended in [48] to solve more complex problems such as that in  $(\mathcal{P}_P)$ . A variable metric version was proposed in [19]. Motivated by the structure of (1), [12] and [18] proposed a Forward-Backward-Forward

scheme [46] to solve it.

In this paper, we focus the unrelaxed Primal–Dual splitting method summarized in Algorithm 1, which covers [13,16,21,28,48]. Though we omit the details here for brevity, our analysis and conclusions carry through to the method proposed in [12,18].

---

**Algorithm 1:** A Primal–Dual splitting method

---

**Initial:** Choose  $\gamma_R, \gamma_J > 0$  and  $\theta \in [-1, +\infty[$ . For  $k = 0$ ,  $x_0 \in \mathbb{R}^n$ ,  $v_0 \in \mathbb{R}^m$ ;

**repeat**

$$\begin{cases} x_{k+1} = \text{prox}_{\gamma_R R}(x_k - \gamma_R \nabla F(x_k) - \gamma_R L^* v_k), \\ \bar{x}_{k+1} = x_{k+1} + \theta(x_{k+1} - x_k), \\ v_{k+1} = \text{prox}_{\gamma_J J^*}(v_k - \gamma_J \nabla G^*(v_k) + \gamma_J L \bar{x}_{k+1}), \end{cases} \quad (2)$$

$k = k + 1$ ;

**until** convergence;

---

**Remark 1.**

- (i) Algorithm 1 is somehow an interesting extension to the literature given the choice of  $\theta$  that we advocate. Indeed, the range  $[-1, +\infty[$  is larger than the one proposed in [28], which is  $[-1, 1]$ . It encompasses the Primal–Dual splitting method in [48] when  $\theta = 1$ , the one in [21] when moreover  $G^* = 0$ . When both  $F = G^* = 0$ , it reduces to the Primal–Dual splitting method proposed in [13,28].
- (ii) It can also be verified that Algorithm 1 covers the Forward–Backward (FB) splitting [37] ( $J^* = G^* = 0$ ), Douglas–Rachford (DR) splitting [24] (if  $F = G^* = 0$ ,  $L = \text{Id}$  and  $\gamma_R = 1/\gamma_J$ ) as special cases; see Section 4 for a discussion or [13] and references therein for more details. Exploiting this relation, in Section 4, we build connections with the results provided in [34,35] for FB-type methods and DR/ADMM. It also should be noted that, DR splitting is the limiting case of the Primal–Dual splitting [13], and the global convergence result of Primal–Dual splitting does not apply to DR.

### 1.3. Contributions

In the literature, most studies on the convergence rate of Primal–Dual splitting methods mainly focus on the global behaviour [9,10,13,14,22,33]. For instance, it is now known that the (partial) duality gap decreases sublinearly (pointwise or in ergodic sense) at the rate  $O(1/k)$  [9,13]. This can be accelerated to  $O(1/k^2)$  on the iterates sequence under strong convexity of either the primal or the dual problem [10,13]. Linear convergence is achieved if both  $R$  and  $J^*$  are strongly convex [10,13]. However, in practice, local linear convergence of the sequence generated by Algorithm 1 has been observed for many problems in the absence of strong convexity (as confirmed by our numerical experiments in Section 6). None of the existing theoretical analysis was able to explain this behaviour so far. Providing the theoretical underpinnings of this local behaviour is the main goal pursued in this paper. Our main findings can be summarized as follows.

**Finite time activity identification** For Algorithm 1, let  $(x^*, v^*)$  be a Kuhn–Tucker pair, *i.e.* a solution of (1). Under a non-degeneracy condition, and provided that

both  $R$  and  $J^*$  are partly smooth relative to  $C^2$ -smooth manifolds, respectively  $\mathcal{M}_{x^*}^R$  and  $\mathcal{M}_{v^*}^{J^*}$  near  $x^*$  and  $v^*$  (see Definition 2.8), we show that the generated primal-dual sequence  $\{(x_k, v_k)\}_{k \in \mathbb{N}}$  which converges to  $(x^*, v^*)$  will identify in finite time the manifold  $\mathcal{M}_{x^*}^R \times \mathcal{M}_{v^*}^{J^*}$  (see Theorem 3.2). In plain words, this means that after a finite number of iterations, say  $K$ , we have  $x_k \in \mathcal{M}_{x^*}^R$  and  $v_k \in \mathcal{M}_{v^*}^{J^*}$  for all  $k \geq K$ .

**Local linear convergence** Capitalizing on this finite identification result, we first show in Proposition 3.4 that the globally non-linear iteration (2) locally linearizes along the identified smooth manifolds, then we deduce that the convergence of the sequence becomes locally linear (see Theorem 3.7). The rate of linear convergence is characterized precisely based on the properties of the identified partly smooth manifolds and the involved linear operator  $L$ .

Moreover, when  $F = G^* = 0, L = \text{Id}$  and  $R, J^*$  are locally polyhedral around  $(x^*, v^*)$ , we show that the convergence rate is parameterized by the *cosine* of the largest principal angle (see Definition 2.5) between the tangent spaces of the two manifolds at  $(x^*, v^*)$  (see Lemma 3.5). This builds a clear connection between the results in this paper and those we drew in our previous works on DR and ADMM [35,36].

#### 1.4. Related work

For the past few years, an increasing attention has been paid to investigate the local linear convergence of first-order proximal splitting methods in the absence of strong convexity. This has been done for instance for FB-type splitting [2,11,26,29,45], and DR/ADMM [4,8,23] for special objective functions. In our previous work [32,34–36], based on the powerful framework provided by partial smoothness, we unified all the above-mentioned work and provide even stronger claims.

To the best of our knowledge, we are aware of only one recent paper [44] which investigated finite identification and local linear convergence of a Primal–Dual splitting method to solve a very special instance of  $(\mathcal{P}_P)$ . More precisely, they assumed  $R$  to be gauge,  $F = \frac{1}{2} \|\cdot\|^2$  (hence strong convexity of the primal problem),  $G^* = 0$  and  $J$  the indicator function of a point. Our work goes much beyond this limited case. It also deepens our current understanding of local behaviour of proximal splitting algorithms by complementing the picture we started in [34,35] for FB and DR methods.

**Paper organization** The rest of the paper is organized as follows. Some useful pre-requisites, including partial smoothness, are collected in Section 2. The main contributions of this paper, *i.e.* finite time activity identification and local linear convergence of Primal–Dual splitting under partial smoothness are the core of Section 3. Several discussions on the obtained result are delivered in Section 4. Section 5 extends the results to the case of more than one infimal convolution. In Section 6, we report various numerical experiments to support our theoretical findings.

## 2. Preliminaries

Throughout the paper,  $\mathbb{N}$  is the set of non-negative integers,  $\mathbb{R}^n$  is a  $n$ -dimensional real Euclidean space equipped with scalar product  $\langle \cdot, \cdot \rangle$  and associated norm  $\|\cdot\|$ .  $\text{Id}_n$  denotes the identity operator on  $\mathbb{R}^n$ , where  $n$  will be dropped if the dimension is clear from the context.

**Sets** For a nonempty convex set  $C \subset \mathbb{R}^n$ , denote  $\text{aff}(C)$  its affine hull, and  $\text{par}(C)$  the smallest subspace parallel to  $\text{aff}(C)$ . Denote  $\iota_C$  the indicator function of  $C$ , and  $P_C$  the orthogonal projection operator onto the set.

**Functions** The sub-differential of a function  $R \in \Gamma_0(\mathbb{R}^n)$  is a set-valued operator,

$$\partial R : \mathbb{R}^n \rightrightarrows \mathbb{R}^n, x \mapsto \{g \in \mathbb{R}^n \mid R(x') \geq R(x) + \langle g, x' - x \rangle, \forall x' \in \mathbb{R}^n\}, \quad (3)$$

which is maximal monotone (see Definition 2.2). For  $R \in \Gamma_0(\mathbb{R}^n)$ , the proximity operator of  $R$  is

$$\text{prox}_R(x) \stackrel{\text{def}}{=} \underset{z \in \mathbb{R}^n}{\text{argmin}} \frac{1}{2} \|z - x\|^2 + R(z). \quad (4)$$

Given a function  $J \in \Gamma_0(\mathbb{R}^m)$ , its Legendre-Fenchel conjugate is defined as

$$J^*(y) = \sup_{v \in \mathbb{R}^m} \langle v, y \rangle - J(v). \quad (5)$$

**Lemma 2.1** (Moreau Identity [39]). *Let  $J \in \Gamma_0(\mathbb{R}^m)$ , then for any  $v \in \mathbb{R}^m$  and  $\gamma > 0$ ,*

$$v = \text{prox}_{\gamma J}(v) + \gamma \text{prox}_{J^*/\gamma}(v/\gamma). \quad (6)$$

Using the Moreau identity, it is straightforward to see that the update of  $v_k$  in Algorithm 1 can be obtained also from  $\text{prox}_{J/\gamma_J}$ .

**Operators** Given a set-valued mapping  $A : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$ , its range is  $\text{ran}(A) = \{y \in \mathbb{R}^n : \exists x \in \mathbb{R}^n \text{ s.t. } y \in A(x)\}$ , and graph is  $\text{gph}(A) \stackrel{\text{def}}{=} \{(x, u) \in \mathbb{R}^n \times \mathbb{R}^n \mid u \in A(x)\}$ .

**Definition 2.2** (Monotone operator). A set-valued mapping  $A : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$  is called monotone if,

$$\langle x_1 - x_2, v_1 - v_2 \rangle \geq 0, \forall (x_1, v_1) \in \text{gph}(A) \text{ and } (x_2, v_2) \in \text{gph}(A). \quad (7)$$

It is moreover maximal monotone if  $\text{gph}(A)$  can not be contained in the graph of any other monotone operator.

Let  $\beta \in ]0, +\infty[$ ,  $B : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , then  $B$  is  $\beta$ -cocoercive if

$$\langle B(x_1) - B(x_2), x_1 - x_2 \rangle \geq \beta \|B(x_1) - B(x_2)\|^2, \forall x_1, x_2 \in \mathbb{R}^n, \quad (8)$$

which implies that  $B$  is  $\beta^{-1}$ -Lipschitz continuous.

For a maximal monotone operator  $A$ ,  $(\text{Id} + A)^{-1}$  is called its resolvent. It is known that for  $R \in \Gamma_0(\mathbb{R}^n)$ ,  $\text{prox}_R = (\text{Id} + \partial R)^{-1}$  [6, Example 23.3].

**Definition 2.3** (Non-expansive operator). An operator  $\mathcal{F} : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is non-expansive if

$$\|\mathcal{F}(x_1) - \mathcal{F}(x_2)\| \leq \|x_1 - x_2\|, \forall x_1, x_2 \in \mathbb{R}^n.$$

For any  $\alpha \in ]0, 1[$ ,  $\mathcal{F}$  is called  $\alpha$ -averaged if there exists a non-expansive operator  $\mathcal{F}'$  such that  $\mathcal{F} = \alpha\mathcal{F}' + (1 - \alpha)\text{Id}$ .

The class of  $\alpha$ -averaged operators is closed under relaxation, convex combination and composition [6,20]. In particular when  $\alpha = \frac{1}{2}$ ,  $\mathcal{F}$  is called firmly non-expansive.

Let  $\mathcal{S}(\mathbb{R}^n) = \{\mathcal{V} \in \mathbb{R}^{n \times n} | \mathcal{V}^T = \mathcal{V}\}$  the set of symmetric positive definite matrices acting on  $\mathbb{R}^n$ . The Loewner partial ordering on  $\mathcal{S}(\mathbb{R}^n)$  is defined as

$$\forall \mathcal{V}_1, \mathcal{V}_2 \in \mathcal{S}(\mathbb{R}^n), \quad \mathcal{V}_1 \succcurlyeq \mathcal{V}_2 \iff \forall x \in \mathbb{R}^n, \langle (\mathcal{V}_1 - \mathcal{V}_2)x, x \rangle \geq 0,$$

that is,  $\mathcal{V}_1 - \mathcal{V}_2$  is positive definite. Given any positive constant  $\nu > 0$ , define  $\mathcal{S}_\nu$  as

$$\mathcal{S}_\nu \stackrel{\text{def}}{=} \{\mathcal{V} \in \mathcal{S}(\mathbb{R}^n) : \mathcal{V} \succcurlyeq \nu \text{Id}\}, \quad (9)$$

*i.e.* the set of symmetric positive definite matrices whose eigenvalues are bounded below by  $\nu$ . For any  $\mathcal{V} \in \mathcal{S}_\nu$ , define the following induced scalar product and norm,

$$\langle x, x' \rangle_{\mathcal{V}} = \langle x, \mathcal{V}x' \rangle, \quad \|x\|_{\mathcal{V}} = \sqrt{\langle x, \mathcal{V}x \rangle}, \quad \forall x, x' \in \mathbb{R}^n.$$

By endowing the Euclidean space  $\mathbb{R}^n$  with the above scalar product and norm, we obtain the Hilbert space which is denoted by  $\mathbb{R}_\nu^n$ .

**Lemma 2.4.** *Let the operators  $A : \mathbb{R}^n \rightrightarrows \mathbb{R}^n$  be maximal monotone,  $B : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be  $\beta$ -cocoercive, and  $\mathcal{V} \in \mathcal{S}_\nu$ . Then for  $\gamma \in ]0, 2\beta\nu[$ ,*

- (i)  $(\text{Id} + \gamma\mathcal{V}^{-1}A)^{-1} : \mathbb{R}_\nu^n \rightarrow \mathbb{R}_\nu^n$  is firmly non-expansive;
- (ii)  $(\text{Id} - \gamma\mathcal{V}^{-1}B) : \mathbb{R}_\nu^n \rightarrow \mathbb{R}_\nu^n$  is  $\frac{\gamma}{2\beta\nu}$ -averaged non-expansive;
- (iii)  $(\text{Id} + \gamma\mathcal{V}^{-1}A)^{-1}(\text{Id} - \gamma\mathcal{V}^{-1}B) : \mathbb{R}_\nu^n \rightarrow \mathbb{R}_\nu^n$  is  $\frac{2\beta\nu}{4\beta\nu - \gamma}$ -averaged non-expansive.

**Proof.**

- (i) See [19, Lemma 3.7(ii)];
- (ii) Since  $B : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is  $\beta$ -cocoercive, given any  $x, x' \in \mathbb{R}^n$ , we have

$$\begin{aligned} \langle x - x', \mathcal{V}^{-1}B(x) - \mathcal{V}^{-1}B(x') \rangle_{\mathcal{V}} &\geq \beta \|B(x) - B(x')\|^2 \\ &= \beta \|\mathcal{V}(\mathcal{V}^{-1}B(x) - \mathcal{V}^{-1}B(x'))\|^2 \\ &= \beta \|\mathcal{V}^{1/2}(\mathcal{V}^{-1}B(x) - \mathcal{V}^{-1}B(x'))\|_{\mathcal{V}}^2 \\ &\geq \beta\nu \|\mathcal{V}^{-1}B(x) - \mathcal{V}^{-1}B(x')\|_{\mathcal{V}}^2, \end{aligned}$$

which means that  $\mathcal{V}^{-1}B : \mathbb{R}_\nu^n \rightarrow \mathbb{R}_\nu^n$  is  $(\beta\nu)$ -cocoercive. The rest of the proof follows [6, Proposition 4.33].

- (iii) See [40, Theorem 3]. □

## 2.1. Angles between subspaces

Let  $T_1$  and  $T_2$  be two subspaces of  $\mathbb{R}^n$ . Without the loss of generality, we assume  $1 \leq p \stackrel{\text{def}}{=} \dim(T_1) \leq q \stackrel{\text{def}}{=} \dim(T_2) \leq n - 1$ .

**Definition 2.5** (Principal angles). The principal angles  $\theta_k \in [0, \frac{\pi}{2}]$ ,  $k = 1, \dots, p$  be-

tween subspaces  $T_1$  and  $T_2$  are defined by, with  $u_0 = v_0 \stackrel{\text{def}}{=} 0$ , and

$$\begin{aligned} \cos(\theta_k) \stackrel{\text{def}}{=}} \langle u_k, v_k \rangle = \max \langle u, v \rangle \text{ s.t. } u \in T_1, v \in T_2, \|u\| = 1, \|v\| = 1, \\ \langle u, u_i \rangle = \langle v, v_i \rangle = 0, i = 0, \dots, k-1. \end{aligned}$$

The principal angles  $\theta_k$  are unique with  $0 \leq \theta_1 \leq \theta_2 \leq \dots \leq \theta_p \leq \pi/2$ .

**Definition 2.6** (Friedrichs angle). The Friedrichs angle  $\theta_F \in ]0, \frac{\pi}{2}]$  between the two subspaces  $T_1$  and  $T_2$  is

$$\begin{aligned} \cos(\theta_F(T_1, T_2)) \stackrel{\text{def}}{=} \max \langle u, v \rangle \text{ s.t. } u \in T_1 \cap (T_1 \cap T_2)^\perp, \|u\| = 1, \\ v \in T_2 \cap (T_1 \cap T_2)^\perp, \|v\| = 1. \end{aligned}$$

The following lemma shows the relation between the Friedrichs and principal angles whose proof can be found in [5, Proposition 3.3].

**Lemma 2.7.** *The Friedrichs angle is the  $(d+1)^{\text{th}}$  principal angle where  $d \stackrel{\text{def}}{=} \dim(T_1 \cap T_2)$ . Moreover,  $\theta_F(T_1, T_2) > 0$ .*

**Remark 2.** The principal angles can be obtained by the singular value decomposition (SVD). For instance, let  $X \in \mathbb{R}^{n \times p}$  and  $Y \in \mathbb{R}^{n \times q}$  form the orthonormal basis for the subspaces  $T_1$  and  $T_2$  respectively. Let  $USV^T$  be the SVD of  $X^T Y \in \mathbb{R}^{p \times q}$ , then  $\cos(\theta_k) = \sigma_k$ ,  $k = 1, 2, \dots, p$  and  $\sigma_k$  corresponds to the  $k^{\text{th}}$  largest singular value in  $\Sigma$ .

## 2.2. Partial smoothness

The concept of partial smoothness is first introduced in [31]. This concept, as well as that of identifiable surfaces [49], captures the essential features of the geometry of non-smoothness which are along the so-called active/identifiable manifold. Loosely speaking, a partly smooth function behaves smoothly as we move along the identifiable submanifold, and sharply if we move transversal to the manifold.

Let  $\mathcal{M}$  be a  $C^2$ -smooth embedded submanifold of  $\mathbb{R}^n$  around a point  $x$ . To lighten the notation, henceforth we state as  $C^2$ -manifold for short. The natural embedding of a submanifold  $\mathcal{M}$  into  $\mathbb{R}^n$  permits to define a Riemannian structure on  $\mathcal{M}$ , and we simply say  $\mathcal{M}$  is a Riemannian manifold.  $\mathcal{T}_{\mathcal{M}}(x)$  denotes the tangent space to  $\mathcal{M}$  at any point near  $x$  in  $\mathcal{M}$ . More materials on manifolds are given in Section A.

Below is the definition of partly smoothness associated to functions in  $\Gamma_0(\mathbb{R}^n)$ .

**Definition 2.8** (Partly smooth function). Let  $R \in \Gamma_0(\mathbb{R}^n)$ , and  $x \in \mathbb{R}^n$  such that  $\partial R(x) \neq \emptyset$ .  $R$  is then said to be partly smooth at  $x$  relative to a set  $\mathcal{M}$  containing  $x$  if

- (i) **Smoothness:**  $\mathcal{M}$  is a  $C^2$ -manifold around  $x$ ,  $R$  restricted to  $\mathcal{M}$  is  $C^2$  around  $x$ ;
- (ii) **Sharpness:** The tangent space  $\mathcal{T}_{\mathcal{M}}(x)$  coincides with  $T_x \stackrel{\text{def}}{=} \text{par}(\partial R(x))^\perp$ ;
- (iii) **Continuity:** The set-valued mapping  $\partial R$  is continuous at  $x$  relative to  $\mathcal{M}$ .

The class of partly smooth functions at  $x$  relative to  $\mathcal{M}$  is denoted as  $\text{PSF}_x(\mathcal{M})$ .

Owing to the results of [31], it can be shown that under transversality assumptions, the set of partly smooth functions is closed under addition and pre-composition by a linear operator. Popular examples of partly smooth functions are summarized in Section 6 whose details can be found in [34].

The next lemma gives expressions of the Riemannian gradient and Hessian (see Section A for definitions) of a partly smooth function.

**Lemma 2.9.** *If  $R \in \text{PSF}_x(\mathcal{M})$ , then for any  $x' \in \mathcal{M}$  near  $x$ ,*

$$\nabla_{\mathcal{M}}R(x') = P_{T_{x'}}(\partial R(x')).$$

In turn, for all  $h \in T_{x'}$ ,

$$\nabla_{\mathcal{M}}^2R(x')h = P_{T_{x'}}\nabla^2\tilde{R}(x')h + \mathfrak{W}_{x'}(h, P_{T_{x'}^\perp}\nabla\tilde{R}(x')),$$

where  $\tilde{R}$  is any smooth extension (representative) of  $R$  on  $\mathcal{M}$ , and  $\mathfrak{W}_x(\cdot, \cdot) : T_x \times T_x^\perp \rightarrow T_x$  is the Weingarten map of  $\mathcal{M}$  at  $x$ .

**Proof.** See [34, Fact 3.3]. □

### 3. Local linear convergence of Primal–Dual splitting methods

In this section, we present the main result of the paper, the local linear convergence analysis of Primal–Dual splitting methods. We start with the finite activity identification of the sequence  $(x_k, v_k)$  generated by the methods, from which we further show that the fixed-point iteration of Primal–Dual splitting methods locally can be linearized, and the linear convergence follows naturally.

#### 3.1. Finite activity identification

Let us first recall the result from [17,48], that under a proper renorming, Algorithm 1 can be written as Forward–Backward. Let  $\theta = 1$ , from the definition of the proximity operator (4), we have that (2) is equivalent to

$$-\begin{bmatrix} \nabla F & 0 \\ 0 & \nabla G^* \end{bmatrix} \begin{pmatrix} x_k \\ v_k \end{pmatrix} \in \begin{bmatrix} \partial R & L^* \\ -L & \partial J^* \end{bmatrix} \begin{pmatrix} x_{k+1} \\ v_{k+1} \end{pmatrix} + \begin{bmatrix} \text{Id}_n/\gamma_R & -L^* \\ -L & \text{Id}_m/\gamma_J \end{bmatrix} \begin{pmatrix} x_{k+1} - x_k \\ v_{k+1} - v_k \end{pmatrix}. \quad (10)$$

Let  $\mathcal{K} = \mathbb{R}^n \times \mathbb{R}^m$  be a product space,  $\text{Id}$  the identity operator on  $\mathcal{K}$ , and define the following variable and operators

$$z_k \stackrel{\text{def}}{=} \begin{pmatrix} x_k \\ v_k \end{pmatrix}, \quad \mathbf{A} \stackrel{\text{def}}{=} \begin{bmatrix} \partial R & L^* \\ -L & \partial J^* \end{bmatrix}, \quad \mathbf{B} \stackrel{\text{def}}{=} \begin{bmatrix} \nabla F & 0 \\ 0 & \nabla G^* \end{bmatrix}, \quad \mathbf{V} \stackrel{\text{def}}{=} \begin{bmatrix} \text{Id}_n/\gamma_R & -L^* \\ -L & \text{Id}_m/\gamma_J \end{bmatrix}. \quad (11)$$

It is easy to verify that  $\mathbf{A}$  is maximal monotone [12],  $\mathbf{B}$  is  $\min\{\beta_F, \beta_G\}$ -cocoercive. For  $\mathbf{V}$ , denote  $\nu = (1 - \sqrt{\gamma_J\gamma_R\|L\|^2}) \min\{\frac{1}{\gamma_J}, \frac{1}{\gamma_R}\}$ , then  $\mathbf{V}$  is symmetric and  $\nu$ -positive definite [19,48]. Define  $\mathcal{K}_{\mathbf{V}}$  the Hilbert space induced by  $\mathbf{V}$ .

Now (10) can be reformulated as

$$z_{k+1} = (\mathbf{V} + \mathbf{A})^{-1}(\mathbf{V} - \mathbf{B})(z_k) = (\text{Id} + \mathbf{V}^{-1}\mathbf{A})^{-1}(\text{Id} - \mathbf{V}^{-1}\mathbf{B})(z_k). \quad (12)$$

Clearly, (12) is the FB splitting on  $\mathcal{K}_{\mathbf{V}}$  [48]. When  $F = G^* = 0$ , it reduces to the metric PPA discussed in [13,28].



Before presenting the finite time activity identification under partial smoothness, we first recall the global convergence of Algorithm 1.

**Lemma 3.1** (Convergence of Algorithm 1). *Consider Algorithm 1 under assumptions (A.1)-(A.4). Let  $\theta = 1$  and choose  $\gamma_R, \gamma_J$  such that*

$$2 \min\{\beta_F, \beta_G\} \min\left\{\frac{1}{\gamma_J}, \frac{1}{\gamma_R}\right\} (1 - \sqrt{\gamma_J \gamma_R \|L\|^2}) > 1, \quad (13)$$

then there exists a Kuhn-Tucker pair  $(x^*, v^*)$  such that  $x^*$  solves  $(\mathcal{P}_P)$ ,  $v^*$  solves  $(\mathcal{P}_D)$ , and  $(x_k, v_k) \rightarrow (x^*, v^*)$ .

*Proof.* See [48, Corollary 4.2]. □

**Remark 3.**

- (i) Assumption (A.4) is important to ensure the existence of Kuhn-Tucker pairs. There are sufficient conditions which ensure that (A.4) can be satisfied. For instance, assuming that  $(\mathcal{P}_P)$  has at least one solution and some classical domain qualification condition is satisfied (see e.g. [18, Proposition 4.3]), assumption (A.4) can be shown to be in force.
- (ii) It is obvious from (13) that  $\gamma_J \gamma_R \|L\|^2 < 1$ , which is also the condition needed in [13] for convergence for the special case  $F = G^* = 0$ . The convergence condition in [21] differs from (13), however,  $\gamma_J \gamma_R \|L\|^2 < 1$  still is a key condition. The values of  $\gamma_J, \gamma_R$  can also be made varying along iterations, and convergence of the iteration remains under the rule provided in [19]. However, for the sake of brevity, we omit the details of this case here.
- (iii) Lemma 3.1 addresses global convergence of the iterates provided by Algorithm 1 only for the case  $\theta = 1$ . For the choices  $\theta \in [-1, 1[ \cup ]1, +\infty[$ , so far the corresponding convergence of the iteration cannot be obtained directly, and a correction step as proposed in [28] for  $\theta \in [-1, 1[$  is needed so that the iteration is a contraction. Unfortunately, such a correction step leads to a new iterative scheme, different from (2); see [28] for more details.

In a very recent paper [50], the authors also proved the convergence of the Primal-Dual splitting method of [48] for the case of  $\theta \in [-1, 1]$  with a proper modification of the iterates. Since the main focus of this work is to investigate local convergence behaviour, the analysis of global convergence of Algorithm 1 for any  $\theta \in [-1, +\infty[$  is beyond the scope of this paper. Thus, we will mainly consider the case  $\theta = 1$  in our analysis. Nevertheless, as we will see later, locally  $\theta > 1$  could give faster convergence rate compared to the choice  $\theta \in [-1, 1]$  for certain problems. This points out a future direction of research to design new Primal-Dual splitting methods.

Given a solution pair  $(x^*, v^*)$  of  $(\mathcal{P}_P)$ - $(\mathcal{P}_D)$ , we denote  $\mathcal{M}_{x^*}^R$  and  $\mathcal{M}_{v^*}^{J^*}$  the  $C^2$ -smooth manifolds that  $x^*$  and  $v^*$  live in respectively, and denote  $T_{x^*}^R, T_{v^*}^{J^*}$  the tangent spaces of  $\mathcal{M}_{x^*}^R, \mathcal{M}_{v^*}^{J^*}$  at  $x^*$  and  $v^*$  respectively.

**Theorem 3.2** (Finite activity identification). *Consider Algorithm 1 under assumptions (A.1)-(A.4). Let  $\theta = 1$  and choose  $\gamma_R, \gamma_J$  based on Lemma 3.1. Thus  $(x_k, v_k) \rightarrow (x^*, v^*)$ , where  $(x^*, v^*)$  is a Kuhn-Tucker pair that solves  $(\mathcal{P}_P)$ - $(\mathcal{P}_D)$ . If moreover*

$R \in \text{PSF}_{x^*}(\mathcal{M}_{x^*}^R)$  and  $J^* \in \text{PSF}_{v^*}(\mathcal{M}_{v^*}^{J^*})$ , and the non-degeneracy condition holds

$$\begin{aligned} -L^*v^* - \nabla F(x^*) &\in \text{ri}(\partial R(x^*)), \\ Lx^* - \nabla G^*(v^*) &\in \text{ri}(\partial J^*(v^*)). \end{aligned} \tag{ND}$$

Then, there exists a large enough  $K \in \mathbb{N}$  such that for all  $k \geq K$ ,

$$(x_k, v_k) \in \mathcal{M}_{x^*}^R \times \mathcal{M}_{v^*}^{J^*}.$$

Moreover,

- (i) if  $\mathcal{M}_{x^*}^R = x^* + T_{x^*}^R$ , then  $T_{x_k}^R = T_{x^*}^R$  and  $\bar{x}_k \in \mathcal{M}_{x^*}^R$  hold for  $k > K$ .
- (ii) If  $\mathcal{M}_{v^*}^{J^*} = v^* + T_{v^*}^{J^*}$ , then  $T_{v_k}^{J^*} = T_{v^*}^{J^*}$  holds for  $k > K$ .
- (iii) If  $R$  is locally polyhedral around  $x^*$ , then  $\forall k \geq K$ ,  $x_k \in \mathcal{M}_{x^*}^R = x^* + T_{x^*}^R$ ,  $T_{x_k}^R = T_{x^*}^R$ ,  $\nabla_{\mathcal{M}_{x^*}^R} R(x_k) = \nabla_{\mathcal{M}_{x^*}^R} R(x^*)$ , and  $\nabla_{\mathcal{M}_{x^*}^R}^2 R(x_k) = 0$ .
- (iv) If  $J^*$  is locally polyhedral around  $v^*$ , then  $\forall k \geq K$ ,  $v_k \in \mathcal{M}_{v^*}^{J^*} = v^* + T_{v^*}^{J^*}$ ,  $T_{v_k}^{J^*} = T_{v^*}^{J^*}$ ,  $\nabla_{\mathcal{M}_{v^*}^{J^*}} J^*(v_k) = \nabla_{\mathcal{M}_{v^*}^{J^*}} J^*(v^*)$ , and  $\nabla_{\mathcal{M}_{v^*}^{J^*}}^2 J^*(v_k) = 0$ .

**Remark 4.**

- (i) The non-degeneracy condition (ND) is a strengthened version of (1).
- (ii) In general, we have no identification guarantees for  $x_k$  and  $v_k$  if the proximity operators are computed approximately, even if the approximation errors are summable, in which case one can still prove global convergence. The reason behind this is that in the exact case, under condition (ND), the proximal mapping of the partly smooth function  $R$  and that of its restriction to  $\mathcal{M}_{x^*}^R$  locally agree nearby  $x^*$  (and similarly for  $J^*$  and  $v^*$ ). This property can be easily violated if approximate proximal mappings are involved, see [34] for an example.
- (iii) Theorem 3.2 only states the existence of  $K$  after which the identification of the sequences happens, but no bounds are provided. In [34,35], lower bounds of  $K$  for the FB and DR methods are delivered. Though similar lower bounds can be obtained for the Primal–Dual splitting method, the proof is not a simple adaptation of that in [34] even if the Primal–Dual splitting method can be cast as a FB splitting in an appropriate metric. Indeed, the estimate in [34] is intimately tied to the fact that FB was applied to an optimization problem, while in the context of Primal–Dual splitting, FB is applied to a monotone inclusion. Since, moreover, these lower-bounds are only of theoretical interest, we decided to forgo the corresponding details here for the sake of conciseness.

**Proof of Theorem 3.2.** From the definition of proximity operator (4) and the updating of  $x_{k+1}$  in (2), we have

$$\frac{1}{\gamma_R}(x_k - \gamma_R \nabla F(x_k) - \gamma_R L^* v_k - x_{k+1}) \in \partial R(x_{k+1}),$$

which yields

$$\begin{aligned} &\text{dist}(-L^*v^* - \nabla F(x^*), \partial R(x_{k+1})) \\ &\leq \left\| -L^*v^* - \nabla F(x^*) - \frac{1}{\gamma_R}(x_k - \gamma_R \nabla F(x_k) - \gamma_R L^* v_k - x_{k+1}) \right\| \\ &\leq \left( \frac{1}{\gamma_R} + \frac{1}{\beta_F} \right) \|x_k - x^*\| + \|L\| \|v_k - v^*\| \rightarrow 0. \end{aligned}$$

Then similarly for  $v_{k+1}$ , we have

$$\frac{1}{\gamma_J}(v_k - \gamma_J \nabla G^*(v_k) + \gamma_J L \bar{x}_{k+1} - v_{k+1}) \in \partial J^*(v_{k+1}),$$

and

$$\begin{aligned} & \text{dist}(Lx^* - \nabla G^*(v^*), \partial J^*(v_{k+1})) \\ & \leq \|Lx^* - \nabla G^*(v^*) - \frac{1}{\gamma_J}(v_k - \gamma_J \nabla G^*(v_k) + \gamma_J L \bar{x}_{k+1} - v_{k+1})\| \\ & \leq \left(\frac{1}{\gamma_J} + \frac{1}{\beta_G}\right) \|v_k - v^*\| + \|L\|((1 + \theta)\|x_{k+1} - x^*\| + \theta\|x_k - x^*\|) \rightarrow 0. \end{aligned}$$

By assumption,  $J^* \in \Gamma_0(\mathbb{R}^m)$ ,  $R \in \Gamma_0(\mathbb{R}^n)$ , hence they are sub-differentially continuous at every point in their respective domains [42, Example 13.30], and in particular at  $v^*$  and  $x^*$ . It then follows that  $J^*(v_k) \rightarrow J^*(v^*)$  and  $R(x_k) \rightarrow R(x^*)$ . Altogether with the non-degeneracy condition (ND), shows that the conditions of [27, Theorem 5.3] are fulfilled for  $\langle -Lx^* + \nabla G^*(v^*), \cdot \rangle + J^*$  and  $\langle L^*v^* + \nabla F(x^*), \cdot \rangle + R$ , and the finite identification claim follows.

- (i) In this case,  $\mathcal{M}_{x^*}^R$  is an affine subspace, it is straight to have  $\bar{x}_k \in \mathcal{M}_{x^*}^R$ . Then as  $R$  is partly smooth at  $x^*$  relative to  $\mathcal{M}_{x^*}^R$ , the sharpness property holds for all nearby points in  $\mathcal{M}_{x^*}^R$  [31, Proposition 2.10]. Thus for  $k$  large enough, we have indeed  $\mathcal{T}_{x_k}(\mathcal{M}_{x^*}^R) = T_{x_k}^R = T_{x^*}^R$  as claimed.
- (ii) Similar to (i).
- (iii) It is immediate to verify that a locally polyhedral function around  $x^*$  is indeed partly smooth relative to the affine subspace  $x^* + T_{x^*}^R$ , and thus, the first claim follows from (i). For the rest, it is sufficient to observe that by polyhedrality, for any  $x \in \mathcal{M}_{x^*}^R$  near  $x^*$ ,  $\partial R(x) = \partial R(x^*)$ . Therefore, combining local normal sharpness [31, Proposition 2.10] and Lemma A.2 yields the second conclusion.
- (iv) Similar to (iii). □

### 3.2. Locally linearized iteration

Relying on the identification result, now we are able to show that the globally nonlinear fixed-point iteration (12) can be locally linearized along the manifolds  $\mathcal{M}_{x^*}^R \times \mathcal{M}_{v^*}^{J^*}$ . As a result, the convergence rate of the iteration essentially boils down to analyzing the spectral properties of the matrix obtained in the linearization.

Given a Kuhn-Tucker pair  $(x^*, v^*)$ , define the following two functions

$$\bar{R}(x) \stackrel{\text{def}}{=} R(x) + \langle x, L^*v^* + \nabla F(x^*) \rangle, \quad \bar{J}^*(y) \stackrel{\text{def}}{=} J^*(y) - \langle y, Lx^* - \nabla G^*(v^*) \rangle. \quad (14)$$

We have the following lemma.

**Lemma 3.3.** *Let  $(x^*, v^*)$  be a Kuhn-Tucker pair such that  $R \in \text{PSF}_{x^*}(\mathcal{M}_{x^*}^R)$ ,  $J^* \in \text{PSF}_{v^*}(\mathcal{M}_{v^*}^{J^*})$ . Denote the Riemannian Hessians of  $\bar{R}$  and  $\bar{J}^*$  as*

$$H_{\bar{R}} \stackrel{\text{def}}{=} \gamma_R P_{T_{x^*}^R} \nabla_{\mathcal{M}_{x^*}^R}^2 \bar{R}(x^*) P_{T_{x^*}^R} \quad \text{and} \quad H_{\bar{J}^*} \stackrel{\text{def}}{=} \gamma_J P_{T_{v^*}^{J^*}} \nabla_{\mathcal{M}_{v^*}^{J^*}}^2 \bar{J}^*(v^*) P_{T_{v^*}^{J^*}}. \quad (15)$$

Then  $H_{\bar{R}}$  and  $H_{\bar{J}^*}$  are symmetric positive semi-definite under either of the following conditions:

- (i) (ND) holds.

(ii)  $\mathcal{M}_{x^*}^R$  and  $\mathcal{M}_{v^*}^{J^*}$  are affine subspaces.  
Define,

$$W_{\bar{R}} \stackrel{\text{def}}{=} (\text{Id}_n + H_{\bar{R}})^{-1} \quad \text{and} \quad W_{\bar{J}^*} \stackrel{\text{def}}{=} (\text{Id}_m + H_{\bar{J}^*})^{-1}, \quad (16)$$

then both  $W_{\bar{R}}$  and  $W_{\bar{J}^*}$  are firmly non-expansive.

**Proof.** See [35, Lemma 6.1]. □

For the smooth functions  $F$  and  $G^*$ , in addition to (A.1) and (A.2), for the rest of the paper, we assume that

(A.5)  $F$  and  $G^*$  locally are  $C^2$ -smooth around  $x^*$  and  $v^*$  respectively.  
Now define the restricted Hessians of  $F$  and  $G^*$ ,

$$H_F \stackrel{\text{def}}{=} \text{P}_{T_{x^*}^R} \nabla^2 F(x^*) \text{P}_{T_{x^*}^R} \quad \text{and} \quad H_{G^*} \stackrel{\text{def}}{=} \text{P}_{T_{v^*}^{J^*}} \nabla^2 G^*(v^*) \text{P}_{T_{v^*}^{J^*}}. \quad (17)$$

Denote  $\bar{H}_F \stackrel{\text{def}}{=} \text{Id}_n - \gamma_R H_F$ ,  $\bar{H}_{G^*} \stackrel{\text{def}}{=} \text{Id}_m - \gamma_J H_{G^*}$ ,  $\bar{L} \stackrel{\text{def}}{=} \text{P}_{T_{v^*}^{J^*}} L \text{P}_{T_{x^*}^R}$  and

$$M_{\text{PD}} \stackrel{\text{def}}{=} \begin{bmatrix} W_{\bar{R}} \bar{H}_F & -\gamma_R W_{\bar{R}} \bar{L}^* \\ \gamma_J (1 + \theta) W_{\bar{J}^*} \bar{L} W_{\bar{R}} \bar{H}_F - \theta \gamma_J W_{\bar{J}^*} \bar{L} & W_{\bar{J}^*} \bar{H}_{G^*} - \gamma_R \gamma_J (1 + \theta) W_{\bar{J}^*} \bar{L} W_{\bar{R}} \bar{L}^* \end{bmatrix}. \quad (18)$$

We have the following proposition.

**Proposition 3.4** (Local linearization). *Suppose that Algorithm 1 is run under the identification conditions of Theorem 3.2, and moreover assumption (A.5) holds. Then for all  $k$  large enough,*

$$z_{k+1} - z^* = M_{\text{PD}} (z_k - z^*) + o(\|z_k - z^*\|). \quad (19)$$

**Remark 5.**

- (i) For the case of varying  $(\gamma_J, \gamma_R)$  along iteration, *i.e.*  $\{(\gamma_{J,k}, \gamma_{R,k})\}_k$ . According to the result of [34], (19) remains true if these parameters converge to some constants such that condition (13) still holds.
- (ii) Taking  $\bar{H}_{G^*} = \text{Id}_m$  (*i.e.*  $G^* = 0$ ) in (18), one gets the linearized iteration associated to the Primal–Dual splitting method of [21]. If we further let  $\bar{H}_F = \text{Id}_n$ , this will correspond to the linearized version of the method in [13].

**Proof of Proposition 3.4.** From the update of  $x_k$  in (2), we have

$$\begin{aligned} x_k - \gamma_R \nabla F(x_k) - \gamma_R L^* v_k - x_{k+1} &\in \gamma_R \partial R(x_{k+1}), \\ -\gamma_R \nabla F(x^*) - \gamma_R L^* v^* &\in \gamma_R \partial R(x^*). \end{aligned}$$

Denote  $\tau_k^R$  the parallel translation from  $T_{x_k}^R$  to  $T_{x^*}^R$ . Then project on to corresponding tangent spaces and apply parallel translation,

$$\begin{aligned} \gamma_R \tau_k^R \nabla_{\mathcal{M}_{x^*}^R} R(x_{k+1}) &= \tau_k^R \text{P}_{T_{x^*}^R} x_{k+1} (x_k - \gamma_R \nabla F(x_k) - \gamma_R L^* v_k - x_{k+1}) \\ &= \text{P}_{T_{x^*}^R} (x_k - \gamma_R \nabla F(x_k) - \gamma_R L^* v_k - x_{k+1}) \\ &\quad + (\tau_k^R \text{P}_{T_{x^*}^R} x_{k+1} - \text{P}_{T_{x^*}^R}) (x_k - \gamma_R \nabla F(x_k) - \gamma_R L^* v_k - x_{k+1}), \\ \gamma_R \nabla_{\mathcal{M}_{x^*}^R} R(x^*) &= \text{P}_{T_{x^*}^R} (-\gamma_R \nabla F(x^*) - \gamma_R L^* v^*), \end{aligned}$$

which leads to

$$\begin{aligned}
& \gamma_R \tau_k^R \nabla_{\mathcal{M}_{x^*}^R} R(x_{k+1}) - \gamma_R \nabla_{\mathcal{M}_{x^*}^R} R(x^*) \\
&= \mathbb{P}_{T_{x^*}^R}((x_k - \gamma_R \nabla F(x_k) - \gamma_R L^* v_k - x_{k+1}) - (x^* - \gamma_R \nabla F(x^*) - \gamma_R L^* v^* - x^*)) \\
&\quad + \underbrace{(\tau_k^R \mathbb{P}_{T_{x^*}^R} x_{k+1} - \mathbb{P}_{T_{x^*}^R})}_{\text{Term 1}}(-\gamma_R \nabla F(x^*) - \gamma_R L^* v^*) \\
&\quad + \underbrace{(\tau_k^R \mathbb{P}_{T_{x^*}^R} x_{k+1} - \mathbb{P}_{T_{x^*}^R})}_{\text{Term 2}}((x_k - \gamma_R \nabla F(x_k) - \gamma_R L^* v_k - x_{k+1}) + (\gamma_R \nabla F(x^*) + \gamma_R L^* v^*)).
\end{aligned} \tag{20}$$

Moving **Term 1** to the other side leads to

$$\begin{aligned}
& \gamma_R \tau_k^R \nabla_{\mathcal{M}_{x^*}^R} R(x_{k+1}) - \gamma_R \nabla_{\mathcal{M}_{x^*}^R} R(x^*) - (\tau_k^R \mathbb{P}_{T_{x^*}^R} x_{k+1} - \mathbb{P}_{T_{x^*}^R})(-\gamma_R \nabla F(x^*) - \gamma_R L^* v^*) \\
&= \gamma_R \tau_k^R (\nabla_{\mathcal{M}_{x^*}^R} R(x_{k+1}) + (L^* v^* + \nabla F(x^*))) - \gamma_R (\nabla_{\mathcal{M}_{x^*}^R} R(x^*) + (L^* v^* + \nabla F(x^*))) \\
&= \gamma_R \mathbb{P}_{T_{x^*}^R} \nabla_{\mathcal{M}_{x^*}^R}^2 \bar{R}(x^*) \mathbb{P}_{T_{x^*}^R} (x_{k+1} - x^*) + o(\|x_{k+1} - x^*\|),
\end{aligned}$$

where Lemma A.2 is applied. Since  $x_{k+1} = \text{prox}_{\gamma_R R}(x_k - \gamma_R \nabla F(x_k) - \gamma_R L^* v_k)$ ,  $\text{prox}_{\gamma_R R}$  is firmly non-expansive and  $\text{Id}_n - \gamma_R \nabla F$  is non-expansive (under the parameter setting), then

$$\begin{aligned}
& \|(x_k - \gamma_R \nabla F(x_k) - \gamma_R L^* v_k - x_{k+1}) - (x^* - \gamma_R \nabla F(x^*) - \gamma_R L^* v^* - x^*)\| \\
&\leq \|(\text{Id}_n - \gamma_R \nabla F)(x_k) - (\text{Id}_n - \gamma_R \nabla F)(x^*)\| + \gamma_R \|L^* v_k - L^* v^*\| \\
&\leq \|x_k - x^*\| + \gamma_R \|L\| \|v_k - v^*\|.
\end{aligned} \tag{21}$$

Therefore, for **Term 2**, owing to Lemma A.1, we have

$$\begin{aligned}
& (\tau_k^R \mathbb{P}_{T_{x^*}^R} x_{k+1} - \mathbb{P}_{T_{x^*}^R})((x_k - \gamma_R \nabla F(x_k) - \gamma_R L^* v_k - x_{k+1}) - (x^* - \gamma_R \nabla F(x^*) - \gamma_R L^* v^* - x^*)) \\
&= o(\|x_k - x^*\| + \gamma_R \|L\| \|v_k - v^*\|).
\end{aligned}$$

Therefore, from (20), and apply  $x_k - x^* = \mathbb{P}_{T_{x^*}^R}(x_k - x^*) + o(x_k - x^*)$  [32, Lemma 5.1] to  $(x_{k+1} - x^*)$  and  $(x_k - x^*)$ , we get

$$\begin{aligned}
& (\text{Id}_n + H_{\bar{R}})(x_{k+1} - x^*) + o(\|x_{k+1} - x^*\|) \\
&= (x_k - x^*) - \gamma_R \mathbb{P}_{T_{x^*}^R}(\nabla F(x_k) - \nabla F(x^*)) - \gamma_R \mathbb{P}_{T_{x^*}^R} L^*(v_k - v^*) \\
&\quad + o(\|x_k - x^*\| + \gamma_R \|L\| \|v_k - v^*\|).
\end{aligned}$$

Then apply Taylor expansion to  $\nabla F$ , and apply [32, Lemma 5.1] to  $(v_k - v^*)$ ,

$$\begin{aligned}
& (\text{Id}_n + H_{\bar{R}})(x_{k+1} - x^*) \\
&= (\text{Id}_n - \gamma_R H_F)(x_k - x^*) - \gamma_R \bar{L}^*(v_k - v^*) + o(\|x_k - x^*\| + \gamma_R \|L\| \|v_k - v^*\|).
\end{aligned} \tag{22}$$

Then invert  $(\text{Id}_n + H_{\bar{R}})$  and apply [32, Lemma 5.1], we get

$$x_{k+1} - x^* = W_{\bar{R}} \bar{H}_F(x_k - x^*) - \gamma_R W_{\bar{R}} \bar{L}^*(v_k - v^*) + o(\|x_k - x^*\| + \gamma_R \|L\| \|v_k - v^*\|). \tag{23}$$

Now from the update of  $v_{k+1}$

$$\begin{aligned} v_k - \gamma_J \nabla G^*(v_k) + \gamma_J L \bar{x}_{k+1} - v_{k+1} &\in \gamma_J \partial J^*(v_{k+1}), \\ v^* - \gamma_J \nabla G^*(v^*) + \gamma_J L x^* - v^* &\in \gamma_J \partial J^*(v^*). \end{aligned}$$

Denote  $\tau_{k+1}^{J^*}$  the parallel translation from  $T_{v_{k+1}}^{J^*}$  to  $T_{v^*}^{J^*}$ , then

$$\begin{aligned} \gamma_J \tau_{k+1}^{J^*} \nabla_{\mathcal{M}_{v^*}^{J^*}} J^*(v_{k+1}) &= \tau_{k+1}^{J^*} P_{T_{v_{k+1}}^{J^*}} (v_k - \gamma_J \nabla G^*(v_k) + \gamma_J L \bar{x}_{k+1} - v_{k+1}) \\ &= P_{T_{v^*}^{J^*}} (v_k - \gamma_J \nabla G^*(v_k) + \gamma_J L \bar{x}_{k+1} - v_{k+1}) \\ &\quad + (\tau_{k+1}^{J^*} P_{T_{v_{k+1}}^{J^*}} - P_{T_{v^*}^{J^*}}) (v_k - \gamma_J \nabla G^*(v_k) + \gamma_J L \bar{x}_{k+1} - v_{k+1}), \\ \gamma_J \nabla_{\mathcal{M}_{v^*}^{J^*}} J^*(v^*) &= P_{T_{v^*}^{J^*}} (v^* - \gamma_J \nabla G^*(v^*) + \gamma_J L x^* - v^*) \end{aligned}$$

which leads to

$$\begin{aligned} &\gamma_J \tau_{k+1}^{J^*} \nabla_{\mathcal{M}_{v^*}^{J^*}} J^*(v_{k+1}) - \gamma_J \nabla_{\mathcal{M}_{v^*}^{J^*}} J^*(v^*) \\ &= P_{T_{v^*}^{J^*}} ((v_k - \gamma_J \nabla G^*(v_k) + \gamma_J L \bar{x}_{k+1} - v_{k+1}) - (v^* - \gamma_J \nabla G^*(v^*) + \gamma_J L x^* - v^*)) \\ &\quad + (\tau_{k+1}^{J^*} P_{T_{v_{k+1}}^{J^*}} - P_{T_{v^*}^{J^*}}) (v_k - \gamma_J \nabla G^*(v_k) + \gamma_J L \bar{x}_{k+1} - v_{k+1}) \\ &= P_{T_{v^*}^{J^*}} ((v_k - \gamma_J \nabla G^*(v_k) + \gamma_J L \bar{x}_{k+1} - v_{k+1}) - (v^* - \gamma_J \nabla G^*(v^*) + \gamma_J L x^* - v^*)) \\ &\quad + \underbrace{(\tau_{k+1}^{J^*} P_{T_{v_{k+1}}^{J^*}} - P_{T_{v^*}^{J^*}}) (\gamma_J L x^* - \gamma_J \nabla G^*(v^*))}_{\text{Term 3}} \\ &\quad + \underbrace{(\tau_{k+1}^{J^*} P_{T_{v_{k+1}}^{J^*}} - P_{T_{v^*}^{J^*}}) ((v_k - \gamma_J \nabla G^*(v_k) + \gamma_J L \bar{x}_{k+1} - v_{k+1}) + \gamma_J (\nabla G^*(v^*) - L x^*))}_{\text{Term 4}}. \end{aligned} \tag{24}$$

Similarly to the previous analysis, for **Term 3**, move to the lefthand side of the inequality and apply Lemma A.2,

$$\begin{aligned} &\gamma_J \tau_{k+1}^{J^*} \nabla_{\mathcal{M}_{v^*}^{J^*}} J^*(v_{k+1}) - \gamma_J \nabla_{\mathcal{M}_{v^*}^{J^*}} J^*(v^*) - (\tau_{k+1}^{J^*} P_{T_{v_{k+1}}^{J^*}} - P_{T_{v^*}^{J^*}}) (\gamma_J L x^* - \gamma_J \nabla G^*(v^*)) \\ &= \gamma_J \tau_{k+1}^{J^*} (\nabla_{\mathcal{M}_{v^*}^{J^*}} J^*(v_{k+1}) - (L x^* - \nabla G^*(v^*))) - \gamma_J (\nabla_{\mathcal{M}_{v^*}^{J^*}} J^*(v^*) - (L x^* - \nabla G^*(v^*))) \\ &= \gamma_J P_{T_{v^*}^{J^*}} \nabla_{\mathcal{M}_{v^*}^{J^*}}^2 J^*(v^*) P_{T_{v^*}^{J^*}} (v_{k+1} - v^*) + o(\|v_{k+1} - v^*\|). \end{aligned}$$

Since  $\theta \leq 1$ , we have

$$\begin{aligned} \|\bar{x}_{k+1} - x^*\| &\leq (1 + \theta) \|x_{k+1} - x^*\| + \theta \|x_k - x^*\| \\ &\leq 2(\|x_k - x^*\| + \gamma_R \|L\| \|v_k - v^*\|) + \|x_k - x^*\| = 3\|x_k - x^*\| + 2\gamma_R \|L\| \|v_k - v^*\|. \end{aligned}$$

Then for **Term 4**, since  $\gamma_J \gamma_R \|L\|^2 < 1$ ,  $\text{prox}_{\gamma_J J^*}$  is firmly non-expansive and  $\text{Id}_m - \gamma_J \nabla G^*$  is non-expansive, we have

$$\begin{aligned} &(\tau_{k+1}^{J^*} P_{T_{v_{k+1}}^{J^*}} - P_{T_{v^*}^{J^*}}) ((v_k - \gamma_J \nabla G^*(v_k) + \gamma_J L \bar{x}_{k+1} - v_{k+1}) - (v^* - \gamma_J \nabla G^*(v^*) + \gamma_J L x^* - v^*)) \\ &= o(\|v_k - v^*\| + \gamma_J \|L\| \|x_k - x^*\|). \end{aligned}$$

Therefore, from (24), apply [32, Lemma 5.1] to  $(v_{k+1} - v^*)$  and  $(v_k - v^*)$ , we get

$$\begin{aligned} & (\text{Id}_m + H_{\bar{J}^*})(v_{k+1} - v^*) \\ &= (\text{Id}_m - \gamma_J H_{G^*})(v_k - v^*) + \gamma_J \bar{L}(\bar{x}_{k+1} - x^*) + o(\|v_k - v^*\| + \gamma_J \|L\| \|x_k - x^*\|). \end{aligned} \quad (25)$$

Then similar to (23), we get from (25)

$$\begin{aligned} v_{k+1} - v^* &= W_{\bar{J}^*} \bar{H}_{G^*}(v_k - v^*) + \gamma_J W_{\bar{J}^*} \bar{L}(\bar{x}_{k+1} - x^*) + o(\|v_k - v^*\| + \gamma_J \|L\| \|x_k - x^*\|) \\ &= W_{\bar{J}^*} \bar{H}_{G^*}(v_k - v^*) + (1 + \theta) \gamma_J W_{\bar{J}^*} \bar{L}(x_{k+1} - x^*) - \theta \gamma_J W_{\bar{J}^*} \bar{L}(x_k - x^*) \\ &\quad + o(\|v_k - v^*\| + \gamma_J \|L\| \|x_k - x^*\|) \\ &= W_{\bar{J}^*} \bar{H}_{G^*}(v_k - v^*) - \theta \gamma_J W_{\bar{J}^*} \bar{L}(x_k - x^*) \\ &\quad + (1 + \theta) \gamma_J W_{\bar{J}^*} \bar{L}(W_{\bar{R}} \bar{H}_F(x_k - x^*) - \gamma_R W_{\bar{R}} \bar{L}^*(v_k - v^*)) \\ &\quad + o(\|x_k - x^*\| + \gamma_R \|L\| \|v_k - v^*\|) + o(\|v_k - v^*\| + \gamma_J \|L\| \|x_k - x^*\|) \\ &= (W_{\bar{J}^*} \bar{H}_{G^*} - (1 + \theta) \gamma_J \gamma_R W_{\bar{J}^*} \bar{L} W_{\bar{R}} \bar{L}^*)(v_k - v^*) \\ &\quad + ((1 + \theta) \gamma_J W_{\bar{J}^*} \bar{L} W_{\bar{R}} \bar{H}_F - \theta \gamma_J W_{\bar{J}^*} \bar{L})(x_k - x^*) \\ &\quad + o(\|x_k - x^*\| + \gamma_R \|L\| \|v_k - v^*\|) + o(\|v_k - v^*\| + \gamma_J \|L\| \|x_k - x^*\|). \end{aligned} \quad (26)$$

Now we consider the small  $o$ -terms. For the 2 small  $o$ -terms in (22) and (25). First, let  $a_1, a_2$  be two constants, then we have

$$|a_1| + |a_2| = \sqrt{(|a_1| + |a_2|)^2} \leq \sqrt{2(a_1^2 + a_2^2)} = \sqrt{2} \left\| \begin{pmatrix} a_1 \\ a_2 \end{pmatrix} \right\|.$$

Denote  $b = \max\{1, \gamma_J \|L\|, \gamma_R \|L\|\}$ , then

$$\begin{aligned} & (\|v_k - v^*\| + \gamma_J \|L\| \|x_k - x^*\|) + (\|x_k - x^*\| + \gamma_R \|L\| \|v_k - v^*\|) \\ & \leq 2b(\|x_k - x^*\| + \|v_k - v^*\|) \leq 2\sqrt{2}b \left\| \begin{pmatrix} x_k - x^* \\ v_k - v^* \end{pmatrix} \right\|. \end{aligned}$$

Combining this with (23) and (26), and ignoring the constants of the small  $o$ -term leads to the claimed result.  $\square$

Now we need to study the spectral properties of  $M_{\text{PD}}$ . Let  $p \stackrel{\text{def}}{=} \dim(T_{x^*}^R), q \stackrel{\text{def}}{=} \dim(T_{v^*}^{J^*})$  be the dimensions of the tangent spaces  $T_{x^*}^R$  and  $T_{v^*}^{J^*}$  respectively, define  $S_{x^*}^R \stackrel{\text{def}}{=} (T_{x^*}^R)^\perp$  and  $S_{v^*}^{J^*} \stackrel{\text{def}}{=} (T_{v^*}^{J^*})^\perp$ . Assume that  $q \geq p$  (alternative situations are discussed in Remark 6). Let

$$\bar{L} = X \Sigma_{\bar{L}} Y^*$$

be the singular value decomposition of  $\bar{L}$ , and define the rank as  $l \stackrel{\text{def}}{=} \text{rank}(\bar{L})$ . Clearly, we have  $l \leq p$ . Denote  $M_{\text{PD}}^k$  the  $k^{\text{th}}$  power of  $M_{\text{PD}}$ .

**Lemma 3.5** (Convergence property of  $M_{\text{PD}}$ ). *The following holds for the matrix  $M_{\text{PD}}$ :*

(i) *If  $\theta = 1$ , then there exists a finite matrix  $M_{\text{PD}}^\infty$  to which  $M_{\text{PD}}^k$  converges, i.e.*

$$M_{\text{PD}}^\infty \stackrel{\text{def}}{=} \lim_{k \rightarrow \infty} M_{\text{PD}}^k. \quad (27)$$

Moreover,

$$\forall k \in \mathbb{N}, M_{\text{PD}}^k - M_{\text{PD}}^\infty = (M_{\text{PD}} - M_{\text{PD}}^\infty)^k \text{ and } \rho(M_{\text{PD}} - M_{\text{PD}}^\infty) < 1.$$

Given any  $\rho \in ]\rho(M_{\text{PD}} - M_{\text{PD}}^\infty), 1[$ , there is  $K$  large enough such that for all  $k \geq K$ ,

$$\|M_{\text{PD}}^k - M_{\text{PD}}^\infty\| = O(\rho^k).$$

- (ii) If  $F = G^* = 0$ , and  $R, J^*$  are locally polyhedral around  $(x^*, v^*)$ . Then given any  $\theta \in ]0, 1[$ ,  $M_{\text{PD}}$  is convergent with

$$M_{\text{PD}}^\infty = \begin{bmatrix} Y & \\ & X \end{bmatrix} \begin{bmatrix} 0_l & & & \\ & \text{Id}_{n-l} & & \\ & & 0_l & \\ & & & \text{Id}_{m-l} \end{bmatrix} \begin{bmatrix} Y^* & \\ & X^* \end{bmatrix}. \quad (28)$$

Moreover, all the eigenvalues of  $M_{\text{PD}} - M_{\text{PD}}^\infty$  are complex with the spectral radius

$$\rho(M_{\text{PD}} - M_{\text{PD}}^\infty) = \sqrt{1 - \theta \gamma_R \gamma_J \sigma_{\min}^2} < 1, \quad (29)$$

where  $\sigma_{\min}$  is the smallest non-zero singular value of  $\bar{L}$ .

**Remark 6.** We discuss in short several possible cases of (28) when  $F = G^* = 0$  and  $R, J^*$  are locally polyhedral around  $(x^*, v^*)$ .

- (i) If  $L = \text{Id}$ , then  $\bar{L} = P_{T_{v^*}^{J^*}} P_{T_{x^*}^R}$  and  $\sigma_{\min}$  is the cosine value of the biggest principal angle (yet strictly smaller than  $\pi/2$ ) between tangent spaces  $T_{x^*}^R$  and  $T_{v^*}^{J^*}$ .
- (ii) For the spectral radius formula in (29), let us consider the case of Arrow–Hurwicz scheme [3], i.e.  $\theta = 0$ . Let  $R, J^*$  be locally polyhedral, and  $\Sigma_{\bar{L}} = (\sigma_j)_{\{j=1, \dots, l\}}$  be the singular values of  $\bar{L}$ , then the eigenvalues of  $M_{\text{PD}}$  are

$$\rho_j = \frac{1}{2} \left( (2 - \gamma_R \gamma_J \sigma_j^2) \pm \sqrt{\gamma_R \gamma_J \sigma_j^2 (\gamma_R \gamma_J \sigma_j^2 - 4)} \right), \quad j \in \{1, \dots, l\}, \quad (30)$$

which apparently are complex ( $\gamma_R \gamma_J \sigma_j^2 \leq \gamma_R \gamma_J \|L\|^2 < 1$ ). Moreover,

$$|\rho_j| = \frac{1}{2} \sqrt{(2 - \gamma_R \gamma_J \sigma_j^2)^2 - \gamma_R \gamma_J \sigma_j^2 (\gamma_R \gamma_J \sigma_j^2 - 4)} = 1.$$

This implies that  $M_{\text{PD}}$  has multiple eigenvalues with absolute values all equal to 1, then owing to the result of [5], we have  $M_{\text{PD}}$  is not convergent.

Furthermore, for  $\theta \in [-1, 0[$ , we have  $1 - \theta \gamma_R \gamma_J \sigma_{\min}^2 > 1$  meaning that  $M_{\text{PD}}$  is not convergent, this implies that the correction step proposed in [28] is necessary for  $\theta \in [-1, 0]$ . Discussion on  $\theta > 1$  is left to Section 4.

**Proof of Proposition 3.5.**

- (i) When  $\theta = 1$ ,  $M_{\text{PD}}$  becomes

$$M_{\text{PD}} = \begin{bmatrix} W_{\bar{R}} \bar{H}_F & -\gamma_R W_{\bar{R}} \bar{L}^* \\ 2\gamma_J W_{\bar{J}^*} \bar{L} W_{\bar{R}} \bar{H}_F - \gamma_J W_{\bar{J}^*} \bar{L} & W_{\bar{J}^*} \bar{H}_{G^*} - 2\gamma_R \gamma_J W_{\bar{J}^*} \bar{L} W_{\bar{R}} \bar{L}^* \end{bmatrix} \quad (31)$$

Next we show that  $M_{\text{PD}}$  is averaged non-expansive.



First define the following matrices

$$\mathbf{A} = \begin{bmatrix} H_R/\gamma_R & \bar{L}^* \\ -\bar{L} & H_{J^*}/\gamma_J \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} H_F & 0 \\ 0 & H_{G^*} \end{bmatrix}, \quad \mathbf{V} = \begin{bmatrix} \text{Id}_n/\gamma_R & -\bar{L}^* \\ -\bar{L} & \text{Id}_m/\gamma_J \end{bmatrix}, \quad (32)$$

where we have  $\mathbf{A}$  is maximal monotone [12],  $\mathbf{B}$  is  $\min\{\beta_F, \beta_G\}$ -cocoercive, and  $\mathbf{V}$  is  $\nu$ -positive definite with  $\nu = (1 - \sqrt{\gamma_J \gamma_R \|L\|^2}) \min\{\frac{1}{\gamma_R}, \frac{1}{\gamma_J}\}$ .

Now we have

$$\mathbf{V} + \mathbf{A} = \begin{bmatrix} \frac{\text{Id}_n + H_R}{\gamma_R} & 0 \\ -2\bar{L} & \frac{\text{Id}_m + H_{J^*}}{\gamma_J} \end{bmatrix} \Rightarrow (\mathbf{V} + \mathbf{A})^{-1} = \begin{bmatrix} \gamma_R W_{\bar{R}} & 0 \\ 2\gamma_J \gamma_R W_{\bar{J}^*} \bar{L} W_{\bar{R}} & \gamma_J W_{\bar{J}^*} \end{bmatrix},$$

and  $\mathbf{V} - \mathbf{B} = \begin{bmatrix} \frac{1}{\gamma_R} \bar{H}_F & -\bar{L}^* \\ -\bar{L} & \frac{1}{\gamma_J} \bar{H}_{G^*} \end{bmatrix}$ . As a result, we get

$$\begin{aligned} (\mathbf{V} + \mathbf{A})^{-1}(\mathbf{V} - \mathbf{B}) &= \begin{bmatrix} \gamma_R W_{\bar{R}} & 0 \\ 2\gamma_J \gamma_R W_{\bar{J}^*} \bar{L} W_{\bar{R}} & \gamma_J W_{\bar{J}^*} \end{bmatrix} \begin{bmatrix} \frac{1}{\gamma_R} \bar{H}_F & -\bar{L}^* \\ -\bar{L} & \frac{1}{\gamma_J} \bar{H}_{G^*} \end{bmatrix} \\ &= \begin{bmatrix} W_{\bar{R}} \bar{H}_F & -\gamma_R W_{\bar{R}} \bar{L}^* \\ 2\gamma_J W_{\bar{J}^*} \bar{L} W_{\bar{R}} \bar{H}_F - \gamma_J W_{\bar{J}^*} \bar{L} & W_{\bar{J}^*} \bar{H}_{G^*} - 2\gamma_J \gamma_R W_{\bar{J}^*} \bar{L} W_{\bar{R}} \bar{L}^* \end{bmatrix}, \end{aligned}$$

which is exactly (31).

From Lemma 2.4 we know that  $M_{\text{PD}} : \mathcal{K}_{\mathbf{V}} \rightarrow \mathcal{K}_{\mathbf{V}}$  is averaged non-expansive, hence it is convergent [6]. Then we have the induced matrix norm

$$\lim_{k \rightarrow \infty} \|M_{\text{PD}}^k - M_{\text{PD}}^\infty\|_{\mathcal{V}} = \lim_{k \rightarrow \infty} \|M_{\text{PD}} - M_{\text{PD}}^\infty\|_{\mathcal{V}}^k = 0.$$

Since we are in the finite dimensional space and  $\mathbf{V}$  is an isomorphism, then the above limit implies that

$$\lim_{k \rightarrow \infty} \|M_{\text{PD}} - M_{\text{PD}}^\infty\|^k = 0,$$

which means that  $\rho(M_{\text{PD}} - M_{\text{PD}}^\infty) < 1$ . The rest of the proof is classical using the spectral radius formula, see *e.g.* [5, Theorem 2.12(i)].

- (ii) When  $R$  and  $J^*$  are locally polyhedral, then  $W_{\bar{R}} = \text{Id}_n, W_{\bar{J}^*} = \text{Id}_m$ , altogether with  $F = 0, G = 0$ , for any  $\theta \in [0, 1]$ , we have

$$M_{\text{PD}} = \begin{bmatrix} \text{Id}_n & -\gamma_R \bar{L}^* \\ \gamma_J \bar{L} & \text{Id}_m - \gamma_R \gamma_J (1 + \theta) \bar{L} \bar{L}^* \end{bmatrix}. \quad (33)$$

With the SVD of  $\bar{L}$ , for  $M_{\text{PD}}$ , we have

$$\begin{aligned} M_{\text{PD}} &= \begin{bmatrix} \text{Id}_n & -\gamma_R \bar{L}^* \\ \gamma_J \bar{L} & \text{Id}_m - (1 + \theta) \gamma_R \gamma_J \bar{L} \bar{L}^* \end{bmatrix} = \begin{bmatrix} Y Y^* & -\gamma_R Y \Sigma_{\bar{L}}^* X^* \\ \gamma_J X \Sigma_{\bar{L}} Y^* & X X^* - (1 + \theta) \gamma_R \gamma_J X \Sigma_{\bar{L}}^2 X^* \end{bmatrix} \\ &= \begin{bmatrix} Y & \\ & X \end{bmatrix} \underbrace{\begin{bmatrix} \text{Id}_n & -\gamma_R \Sigma_{\bar{L}}^* \\ \gamma_J \Sigma_{\bar{L}} & \text{Id}_m - (1 + \theta) \gamma_R \gamma_J \Sigma_{\bar{L}}^2 \end{bmatrix}}_{M_\Sigma} \begin{bmatrix} Y^* & \\ & X^* \end{bmatrix}. \end{aligned}$$

Since we assume that  $\text{rank}(\bar{L}) = l \leq p$ , then  $\Sigma_{\bar{L}}$  can be represented as

$$\Sigma_{\bar{L}} = \begin{bmatrix} \Sigma_l & 0_{l,n-l} \\ 0_{m-l,l} & 0_{m-l,n-l} \end{bmatrix},$$

where  $\Sigma_l = (\sigma_j)_{j=1,\dots,l}$ . Back to  $M_\Sigma$ , we have

$$M_\Sigma = \begin{bmatrix} \text{Id}_l & 0_{l,n-l} & -\gamma_R \Sigma_l & 0_{l,m-l} \\ 0_{n-l,l} & \text{Id}_{n-l} & 0_{n-l,l} & 0_{n-l,m-l} \\ \gamma_J \Sigma_l & 0_{l,n-l} & \text{Id}_l - (1+\theta)\gamma_R \gamma_J \Sigma_l^2 & 0_{l,m-l} \\ 0_{m-l,l} & 0_{m-l,n-l} & 0_{m-l,l} & \text{Id}_{m-l} \end{bmatrix}.$$

Let's study the eigenvalues of  $M_\Sigma$ ,

$$\begin{aligned} & |M_\Sigma - \rho \text{Id}_{m+n}| \\ &= \left| \begin{bmatrix} (1-\rho)\text{Id}_l & 0_{l,n-l} & -\gamma_R \Sigma_l & 0_{l,m-l} \\ 0_{n-l,l} & (1-\rho)\text{Id}_{n-l} & 0_{n-l,l} & 0_{n-l,m-l} \\ \gamma_J \Sigma_l & 0_{l,n-l} & (1-\rho)\text{Id}_l - (1+\theta)\gamma_R \gamma_J \Sigma_l^2 & 0_{l,m-l} \\ 0_{m-l,l} & 0_{m-l,n-l} & 0_{m-l,l} & (1-\rho)\text{Id}_{m-l} \end{bmatrix} \right| \\ &= (1-\rho)^{m+n-2l} \left| \begin{bmatrix} (1-\rho)\text{Id}_l & -\gamma_R \Sigma_l \\ \gamma_J \Sigma_l & (1-\rho)\text{Id}_l - (1+\theta)\gamma_R \gamma_J \Sigma_l^2 \end{bmatrix} \right|. \end{aligned}$$

Since  $(-\gamma_R \Sigma_l)((1-\rho)\text{Id}_l) = ((1-\rho)\text{Id}_l)(-\gamma_R \Sigma_l)$ , then by [43, Theorem 3], we have

$$\begin{aligned} |M_\Sigma - \rho \text{Id}_{m+n}| &= (1-\rho)^{m+n-2l} \left| \begin{bmatrix} (1-\rho)\text{Id}_l & -\gamma_R \Sigma_l \\ \gamma_J \Sigma_l & (1-\rho)\text{Id}_l - (1+\theta)\gamma_R \gamma_J \Sigma_l^2 \end{bmatrix} \right| \\ &= (1-\rho)^{m+n-2l} \left| [(1-\rho)((1-\rho)\text{Id}_l - (1+\theta)\gamma_R \gamma_J \Sigma_l^2) + \gamma_R \gamma_J \Sigma_l \Sigma_l] \right| \\ &= (1-\rho)^{m+n-2l} \left| [(1-\rho)^2 \text{Id}_l - (1-\rho)(1+\theta)\gamma_R \gamma_J \Sigma_l^2 + \gamma_R \gamma_J \Sigma_l \Sigma_l] \right| \\ &= (1-\rho)^{m+n-2l} \prod_{j=1}^l (\rho^2 - (2 - (1+\theta)\gamma_J \gamma_R \sigma_j^2)\rho + (1 - \theta\gamma_J \gamma_R \sigma_j^2)). \end{aligned}$$

For the eigenvalues  $\rho$ , clearly, except the 1's, we have for  $j = 1, \dots, l$

$$\rho_j = \frac{(2 - (1+\theta)\gamma_J \gamma_R \sigma_j^2) \pm \sqrt{(1+\theta)^2 \gamma_J^2 \gamma_R^2 \sigma_j^4 - 4\gamma_J \gamma_R \sigma_j^2}}{2}.$$

Since  $\gamma_J \gamma_R \sigma_j^2 \leq \gamma_J \gamma_R \|L\|^2 < 1$ , then  $\rho_j$  are complex and

$$|\rho_j| = \frac{1}{2} \sqrt{(2 - (1+\theta)\gamma_J \gamma_R \sigma_j^2)^2 - ((1+\theta)^2 \gamma_J^2 \gamma_R^2 \sigma_j^4 - 4\gamma_J \gamma_R \sigma_j^2)} = \sqrt{1 - \theta\gamma_J \gamma_R \sigma_j^2} < 1.$$

As a result, we also obtain the  $M_{\text{PD}}^\infty$ , which reads

$$M_{\text{PD}}^\infty = \begin{bmatrix} Y & \\ & X \end{bmatrix} \begin{bmatrix} 0_l & & & \\ & \text{Id}_{n-l} & & \\ & & 0_l & \\ & & & \text{Id}_{m-l} \end{bmatrix} \begin{bmatrix} Y^* & \\ & X^* \end{bmatrix}.$$

□

**Corollary 3.6.** *Suppose that Algorithm 1 is run under the identification conditions of Theorem 3.2, and moreover assumption (A.5) holds. Then the following holds*

(i) *the linearized iteration (19) is equivalent to*

$$(\mathbf{Id} - M_{\text{PD}}^\infty)(\mathbf{z}_{k+1} - \mathbf{z}^*) = M_{\text{PD}}(\mathbf{Id} - M_{\text{PD}}^\infty)(\mathbf{z}_k - \mathbf{z}^*) + o(\|\mathbf{Id} - M_{\text{PD}}^\infty\| \|\mathbf{z}_k - \mathbf{z}^*\|). \quad (34)$$

(ii) *If moreover  $R, J^*$  are locally polyhedral around  $(x^*, v^*)$ , and  $F, G^*$  are quadratic, then  $M_{\text{PD}}^\infty(\mathbf{z}_k - \mathbf{z}^*) = 0$  for all  $k$  large enough, and (34) becomes*

$$\mathbf{z}_{k+1} - \mathbf{z}^* = (M_{\text{PD}} - M_{\text{PD}}^\infty)(\mathbf{z}_k - \mathbf{z}^*). \quad (35)$$

**Proof.** See [35, Corollary 5.1]. □

### 3.3. Local linear convergence

Finally, we are able to present the local linear convergence result.

**Theorem 3.7** (Local linear convergence). *Suppose that Algorithm 1 is run under the identification conditions of Theorem 3.2, and moreover assumption (A.5) holds. Then:*

(i) *given any  $\rho \in ]\rho(M_{\text{PD}} - M_{\text{PD}}^\infty), 1[$ , there exists a  $K$  large enough such that  $\forall k \geq K$ ,*

$$\|(\mathbf{Id} - M_{\text{PD}}^\infty)(\mathbf{z}_k - \mathbf{z}^*)\| = O(\rho^{k-K}). \quad (36)$$

(ii) *If moreover,  $R, J^*$  are locally polyhedral around  $(x^*, v^*)$ , and  $F, G^*$  are quadratic, then there exists a  $K$  large enough such that for all  $k \geq K$ , we have directly*

$$\|\mathbf{z}_k - \mathbf{z}^*\| = O(\rho^{k-K}), \quad (37)$$

for  $\rho \in ]\rho(M_{\text{PD}} - M_{\text{PD}}^\infty), 1[$ .

**Proof.** See [35, Theorem 5.1]. □

#### Remark 7.

- (i) Similar to Proposition 3.4 and Remark 5, the above result remains hold if  $(\gamma_J, \gamma_R)$  are varying yet convergent. However, the local rate convergence of  $\|\mathbf{z}_k - \mathbf{z}^*\|$  will depends on how fast  $\{(\gamma_{J,k}, \gamma_{R,k})\}_k$  converge, that means, if they converge at a sublinear rate, then the convergence rate of  $\|\mathbf{z}_k - \mathbf{z}^*\|$  will eventually become sublinear. See [35, Section 8.3] for the case of Douglas–Rachford splitting method.
- (ii) When  $F = G^* = 0$  and both  $R$  and  $J^*$  are locally polyhedral around the  $(x^*, v^*)$ , then the convergence rate of the Primal–Dual splitting method is controlled by  $\theta$  and  $\gamma_J \gamma_R$  as shown in (29); see the upcoming section for a detailed discussion.

For general situations (*i.e.*  $F, G^*$  are nontrivial and  $R, J^*$  are general partly smooth functions), the factors that contribute to the local convergence rate are much more complicated, such as the Riemannian Hessians of  $R, J^*$ .

## 4. Discussions

In this part, we present several discussions on the obtained local linear convergence result, including the effects of  $\theta \geq 1$ , local oscillation and connections with FB and DR methods.

To make the discussion easier to deliver, for the rest of this section we focus on the case where  $F = G^* = 0$ , *i.e.* the Primal–Dual splitting method of [13], and moreover  $R, J^*$  are locally polyhedral around the Kuhn–Tucker pair  $(x^*, v^*)$ . Under such setting, the matrix defined in (18) becomes

$$M_{\text{PD}} \stackrel{\text{def}}{=} \begin{bmatrix} \text{Id}_n & -\gamma_R \bar{L}^* \\ \gamma_J \bar{L} & \text{Id}_m - (1 + \theta)\gamma_J \gamma_R \bar{L} \bar{L}^* \end{bmatrix}. \quad (38)$$

#### 4.1. Choice of $\theta$

Owing to Lemma 3.5, the matrix  $M_{\text{PD}}$  in (38) is convergent for  $\theta \in ]0, 1]$ , see Eq. (28), with the spectral radius

$$\rho(M_{\text{PD}} - M_{\text{PD}}^\infty) = \sqrt{1 - \theta\gamma_R \gamma_J \sigma_{\min}^2} < 1, \quad (39)$$

with  $\sigma_{\min}$  being the smallest non-zero singular value of  $\bar{L}$ .

In general, given a solution pair  $(x^*, v^*)$ ,  $\sigma_{\min}$  is fixed, hence the spectral radius  $\rho(M_{\text{PD}} - M_{\text{PD}}^\infty)$  is simply controlled by  $\theta$  and the product  $\gamma_J \gamma_R$ . To make the local convergence rate as faster as possible, it is obvious that we need to make the value of  $\theta\gamma_J \gamma_R$  as big as possible. Recall in the global convergence of Primal–Dual splitting method or the result from [13], that  $\gamma_J \gamma_R \|L\|^2 < 1$ . Denote  $\sigma_{\max}$  the biggest singular value of  $\bar{L}$ . It is then straightforward that  $\gamma_J \gamma_R \sigma_{\max}^2 \leq \gamma_J \gamma_R \|L\|^2 < 1$  and moreover

$$\begin{aligned} \rho(M_{\text{PD}} - M_{\text{PD}}^\infty) &= \sqrt{1 - \theta\gamma_R \gamma_J \sigma_{\min}^2} \\ &> \sqrt{1 - \theta(\sigma_{\min}/\|L\|)^2} \geq \sqrt{1 - \theta(\sigma_{\min}/\sigma_{\max})^2}. \end{aligned} \quad (40)$$

If we define  $\text{cnd} \stackrel{\text{def}}{=} \sigma_{\max}/\sigma_{\min}$  the condition number of  $\bar{L}$ , then we have

$$\rho(M_{\text{PD}} - M_{\text{PD}}^\infty) > \sqrt{1 - \theta(1/\text{cnd})^2}.$$

To this end, it is clear that  $\theta = 1$  gives the best convergence rate for  $\theta \in [-1, 1]$ . Next let us look at what happens locally if we choose  $\theta > 1$ . The spectral radius formula (39) implies that bigger value of  $\theta$  yields smaller spectral radius  $\rho(M_{\text{PD}} - M_{\text{PD}}^\infty)$ . Therefore, locally we should choose  $\theta$  as big as possible. However, there is an upper bound of  $\theta$  which is discussed below.

Following Remark 6, let  $\Sigma_{\bar{L}} = (\sigma_j)_{\{j=1, \dots, l\}}$  be the singular values of  $\bar{L}$ , let  $\rho_j$  be the eigenvalue of  $M_{\text{PD}} - M_{\text{PD}}^\infty$ , we have known that  $\rho_j$  is complex with

$$\rho_j = \frac{1}{2} \left( (2 - (1 + \theta)\gamma_R \gamma_J \sigma_j^2) \pm \sqrt{(1 + \theta)^2 \gamma_R^2 \gamma_J^2 \sigma_j^4 - 4\gamma_R \gamma_J \sigma_j^2} \right), \quad |\rho_j| = \sqrt{1 - \theta\gamma_R \gamma_J \sigma_j^2}.$$

Now let  $\theta > 1$ , to ensure  $|\rho_j|$  make sense for all  $j \in \{1, \dots, l\}$ , there must holds

$$1 - \theta\gamma_R \gamma_J \sigma_{\max}^2 \geq 0 \iff \theta \leq \frac{1}{\gamma_R \gamma_J \sigma_{\max}^2},$$

which means that  $\theta$  indeed is bounded from above.

Unfortunately, since  $\bar{L} = P_{T_{x^*}^R} L P_{T_{v^*}^{J^*}}$ , the upper bound can be only obtained if we had the solution pair  $(x^*, v^*)$ . However, in practice one can use back-tracking or the Armijo-Goldstein-rule to find the proper  $\theta$ . See Section 6.4 for an illustration of online searching of  $\theta$ . It should be noted that such updating rule can also be applied to  $\gamma_J, \gamma_R$  since we have  $\|\bar{L}\| \leq \|L\|$ . Moreover, it should be noted that in practice one can choose to enlarge either  $\theta$  or  $\gamma_J \gamma_R$  as they will have very similar acceleration outcome.

**Remark 8.** It should be noted that the above discussion on the effect of  $\theta > 1$  may only valid for the case  $F = 0, G^* = 0$ , *i.e.* the Primal-Dual splitting method of [13]. If  $F$  and/or  $G^*$  are not vanished, then locally,  $\theta < 1$  may give faster convergence rate.

## 4.2. Oscillations

For the inertial Forward-Backward and FISTA [7] methods, it is shown in [34] that they locally oscillate when the inertia momentum are too high (see [34, Section 4.4] for more details). When solving certain type of problems (*i.e.*  $F = G^* = 0$  and  $R, J^*$  are locally polyhedral around the solution pair  $(x^*, v^*)$ ), the Primal-Dual splitting method also locally oscillates (see Figure 6 for an illustration). As revealed in the proof of Lemma 3.5, all the eigenvalues of  $M_{\text{PD}} - M_{\text{PD}}^\infty$  in (38) are complex. This means that locally the sequences generated by the Primal-Dual splitting iteration may oscillate.

For  $\sigma_{\min}$ , the smallest non-zero singular of  $\bar{L}$ , one of its corresponding eigenvalues of  $M_{\text{PD}}$  reads

$$\rho_{\sigma_{\min}} = \frac{1}{2} \left( (2 - (1 + \theta)\gamma_J \gamma_R \sigma_{\min}^2) + \sqrt{(1 + \theta)^2 \gamma_R^2 \gamma_J^2 \sigma_{\min}^4 - 4\gamma_J \gamma_R \sigma_{\min}^2} \right),$$

and  $(1 + \theta)^2 \gamma_R^2 \gamma_J^2 \sigma_{\min}^4 - 4\gamma_J \gamma_R \sigma_{\min}^2 < 0$ . Denote  $\omega$  the argument of  $\rho_{\sigma_{\min}}$ , then

$$\cos(\omega) = \frac{2 - (1 + \theta)\gamma_J \gamma_R \sigma_{\min}^2}{\sqrt{1 - \theta\gamma_J \gamma_R \sigma_{\min}^2}}. \quad (41)$$

The oscillation period of the sequence  $\|z_k - z^*\|$  is then exactly  $\frac{\pi}{\omega}$ . See Figure 6 for an illustration.

## 4.3. Relations with FB and DR/ADMM

In this part, we discuss the relation between the obtained result and our previous work on local linear convergence of Forward-Backward [32,34] and Douglas-Rachford/ADMM [35,36].

### 4.3.1. Forward-Backward splitting

For problem  $(\mathcal{P}_P)$ , when  $J = G^* = 0$ , Algorithm 1 reduces to, denoting  $\gamma = \gamma_R$  and  $\beta = \beta_F$ ,

$$x_{k+1} = \text{prox}_{\gamma R}(x_k - \gamma \nabla F(x_k)), \quad \gamma \in ]0, 2\beta[, \quad (42)$$

which is the non-relaxed FB splitting [37] with constant step-size.

Let  $x^* \in \text{Argmin}(R + F)$  be a global minimizer to which  $\{x_k\}_{k \in \mathbb{N}}$  of (42) converges, the non-degeneracy condition (ND) for identification then becomes  $-\nabla F(x^*) \in$

$\text{ri}(\partial R(x^*))$  which recovers the conditions of [34, Theorem 4.11]. Following the notations of Section 3, define  $M_{\text{FB}} \stackrel{\text{def}}{=} W_{\bar{R}}(\text{Id}_n - \gamma H_F)$ , we have for all  $k$  large enough

$$x_{k+1} - x^* = M_{\text{FB}}(x_k - x^*) + o(\|x_k - x^*\|).$$

From Theorem 3.7, we obtain the following result for the FB splitting method, for the case  $\gamma$  being fixed. Denote  $\mathcal{M}_{x^*}$  the manifold that  $x^*$  lives in.

**Corollary 4.1.** *For problem  $(\mathcal{P}_P)$ , let  $J = G^* = 0$  and suppose that (A.1) holds and  $\text{Argmin}(R + F) \neq \emptyset$ , and the FB iteration (42) creates a sequence  $x_k \rightarrow x^* \in \text{Argmin}(R + F)$  such that  $R \in \text{PSF}_{x^*}(\mathcal{M}_{x^*})$ ,  $F$  is  $C^2$  near  $x^*$ , and condition  $-\nabla F(x^*) \in \text{ri}(\partial R(x^*))$  holds. Then*

- (i) *given any  $\rho \in ]\rho(M_{\text{FB}} - M_{\text{FB}}^\infty), 1[$ , there exists a  $K$  large enough such that for all  $k \geq K$ ,*

$$\|(\text{Id} - M_{\text{FB}}^\infty)(x_k - x^*)\| = O(\rho^{k-K}). \quad (43)$$

- (ii) *If moreover,  $R$  are locally polyhedral around  $x^*$ , there exists a  $K$  large enough such that for all  $k \geq K$ , we have directly*

$$\|x_k - x^*\| = O(\rho^{k-K}), \quad (44)$$

for  $\rho \in ]\rho(M_{\text{FB}} - M_{\text{FB}}^\infty), 1[$ .

**Proof.** Owing to [6],  $M_{\text{FB}}$  is  $\frac{2\beta}{4\beta - \gamma}$ -averaged non-expansive, hence convergent. The convergence rates in (43) and (44) are straightforward from Theorem 3.7.  $\square$

The result in [32,34] required a so-called restricted injectivity (RI) condition, which means that  $H_F$  should be positive definite. Moreover, this RI condition was removed only when  $J$  is locally polyhedral around  $x^*$  (e.g. see [34, Theorem 4.9]). Here, we show that neither RI condition nor polyhedrality are needed to show local linear convergence. As such, the analysis on this paper generalizes that of [34] even for FB splitting. However, the price of removing those conditions is that the obtained convergence rate is on a different criterion (i.e.  $\|(\text{Id} - M_{\text{FB}}^\infty)(x_k - x^*)\|$ ) other than the sequence itself directly (i.e.  $\|x_k - x^*\|$ ).

#### 4.3.2. Douglas–Rachford splitting and ADMM

Let  $F = G^* = 0$  and  $L = \text{Id}$ , then problem  $(\mathcal{P}_P)$  becomes

$$\min_{x \in \mathbb{R}^n} R(x) + J(x).$$

For the above problem, below we briefly show that DR splitting is the limiting case of Primal–Dual splitting by letting  $\gamma_R \gamma_J = 1$ . First, for the Primal–Dual splitting scheme of Algorithm 1, let  $\theta = 1$  and change the order of updating the variables, we obtain the following iteration

$$\begin{cases} v_{k+1} = \text{prox}_{\gamma_J J^*}(v_k + \gamma_J \bar{x}_k) \\ x_{k+1} = \text{prox}_{\gamma_R R}(x_k - \gamma_R v_{k+1}) \\ \bar{x}_{k+1} = 2x_{k+1} - x_k. \end{cases} \quad (45)$$

Then apply the Moreau's identity (6) to  $\text{prox}_{\gamma_J J^*}$ , let  $\gamma_J = 1/\gamma_R$  and define  $z_{k+1} = x_k - \gamma_R v_{k+1}$ , iteration (45) becomes

$$\begin{cases} u_{k+1} = \text{prox}_{\gamma_R J}(2x_k - z_k) \\ z_{k+1} = z_k + u_{k+1} - x_k \\ x_{k+1} = \text{prox}_{\gamma_R R}(z_{k+1}), \end{cases} \quad (46)$$

which is the non-relaxed DR splitting method [24]. At convergence, we have  $u_k, x_k \rightarrow x^* = \text{prox}_{\gamma_R R}(z^*)$  where  $z^*$  is a fixed point of the iteration. See also [13, Section 4.2].

Specializing the derivation of (18) to (45) and (46), we obtain the following two linearized fixed-point operator for (45) and (46) respectively

$$M_{\text{PD}} = \begin{bmatrix} \text{Id}_n & -\gamma_R \mathbf{P}_{T_{x^*}^R} \mathbf{P}_{T_{v^*}^{J^*}} \\ \gamma_J \mathbf{P}_{T_{v^*}^{J^*}} \mathbf{P}_{T_{x^*}^R} & \text{Id}_n - 2\gamma_J \gamma_R \mathbf{P}_{T_{v^*}^{J^*}} \mathbf{P}_{T_{x^*}^R} \mathbf{P}_{T_{v^*}^{J^*}} \end{bmatrix},$$

$$M_{\text{DR}} = \begin{bmatrix} \text{Id}_n & -\gamma_R \mathbf{P}_{T_{x^*}^R} \mathbf{P}_{T_{v^*}^{J^*}} \\ \frac{1}{\gamma_R} \mathbf{P}_{T_{v^*}^{J^*}} \mathbf{P}_{T_{x^*}^R} & \text{Id}_n - 2\mathbf{P}_{T_{v^*}^{J^*}} \mathbf{P}_{T_{x^*}^R} \mathbf{P}_{T_{v^*}^{J^*}} \end{bmatrix}.$$

Owing to (ii) of Lemma 3.5,  $M_{\text{PD}}, M_{\text{DR}}$  are convergent. Let  $\omega$  be the largest principal angle (yet smaller than  $\pi/2$ ) between tangent spaces  $T_{x^*}^R$  and  $T_{v^*}^{J^*}$ , then we have the spectral radius of  $M_{\text{PD}} - M_{\text{PD}}^\infty$  reads ((i) of Remark 6),

$$\begin{aligned} \rho(M_{\text{PD}} - M_{\text{PD}}^\infty) &= \sqrt{1 - \gamma_J \gamma_R \cos^2(\omega)} \\ &\geq \sqrt{1 - \cos^2(\omega)} = \sin(\omega) = \cos(\pi/2 - \omega). \end{aligned} \quad (47)$$

Suppose that the Kuhn-Tucker pair  $(x^*, v^*)$  is unique, and moreover that  $R$  and  $J$  are polyhedral. Therefore, we have that if  $J^*$  is locally polyhedral near  $v^*$  along  $v^* + T_{v^*}^{J^*}$ , then  $J$  is locally polyhedral near  $x^*$  around  $x^* + T_{x^*}^J$ , and moreover there holds  $T_{x^*}^J = (T_{v^*}^{J^*})^\perp$ . As a result, the principal angles between  $T_{x^*}^R, T_{v^*}^{J^*}$  and the ones between  $T_{x^*}^R, T_{x^*}^J$  are complementary, which means that  $\pi/2 - \omega$  is the Friedrichs angle between tangent spaces  $T_{x^*}^R, T_{x^*}^J$ . Thus, following (47), we have

$$\rho(M_{\text{PD}} - M_{\text{PD}}^\infty) = \sqrt{1 - \gamma_J \gamma_R \cos^2(\omega)} \geq \cos(\pi/2 - \omega) = \rho(M_{\text{DR}} - M_{\text{DR}}^\infty).$$

We emphasise the fact that such connection can be drawn only for the polyhedral case, which justifies the different analysis carried in [35,36]. In addition, for DR we were able to characterize situations where finite convergence provably occurs, while this is not (yet) the case for Primal-Dual splitting even for  $R$  and  $J^*$  being locally polyhedral around  $(x^*, v^*)$  and  $F = G^* = 0$  but  $L$  is non-trivial.

## 5. Multiple infimal convolutions

In this section, we consider problem  $(\mathcal{P}_P)$  with more than one infimal convolution. Let  $m \geq 1$  be a positive integer. Consider the problem of solving

$$\min_{x \in \mathbb{R}^n} R(x) + F(x) + \sum_{i=1}^m (J_i \sharp G_i)(L_i x), \quad (\mathcal{P}_P^m)$$

where (A.1) holds for  $R$  and  $F$ , and for every  $i = 1, \dots, m$  the followings are hold:

(A'.2)  $J_i, G_i \in \Gamma_0(\mathbb{R}^{m_i})$ , with  $G_i$  being differentiable and  $\beta_{G_i}$ -strongly convex for  $\beta_{G_i} > 0$ .

(A'.3)  $L_i : \mathbb{R}^n \rightarrow \mathbb{R}^{m_i}$  is a linear operator.

(A'.4) The condition  $0 \in \text{ran}(\partial R + \nabla F + \sum_{i=1}^m L_i^*(\partial J_i \square \partial G_i)L_i)$  holds.

The dual problem of ( $\mathcal{P}_P^m$ ) reads,

$$\min_{v_1 \in \mathbb{R}^{m_1}, \dots, v_m \in \mathbb{R}^{m_m}} (R^* \uplus F^*) \left( - \sum_{i=1}^m L_i^* v_i \right) + \sum_{i=1}^m (J_i^*(v_i) + G_i^*(v_i)). \quad (\mathcal{P}_D^m)$$

Problem ( $\mathcal{P}_P^m$ ) is considered in [19,48], and a Primal–Dual splitting algorithm is proposed there which is an extension of Algorithm 1 using a product space trick, see Algorithm 2 hereafter for details. In both schemes in [19,48],  $\theta$  is set as 1.

---

**Algorithm 2:** A Primal–Dual splitting method

---

**Initial:** Choose  $\gamma_R, (\gamma_{J_i})_i > 0$ . For  $k = 0$ ,  $x_0 \in \mathbb{R}^n$ ,  $v_{i,0} \in \mathbb{R}^{m_i}$ ,  $i \in \{1, \dots, m\}$ ;

**repeat**

$$\left[ \begin{array}{l} x_{k+1} = \text{prox}_{\gamma_R R}(x_k - \gamma_R \nabla F(x_k) - \gamma_R \sum_i L_i^* v_{i,k}) \\ \bar{x}_{k+1} = 2x_{k+1} - x_k \\ \text{For } i = 1, \dots, m \\ \left[ \begin{array}{l} v_{i,k+1} = \text{prox}_{\gamma_{J_i} J_i^*}(v_{i,k} - \gamma_{J_i} \nabla G_i^*(v_{i,k}) + \gamma_{J_i} L_i \bar{x}_{k+1}), \end{array} \right. \end{array} \right. \quad (48)$$

$k = k + 1$ ;

**until** convergence;

---

### 5.1. Product space

The following result is taken from [19]. Define the product space  $\mathcal{K} = \mathbb{R}^n \times \mathbb{R}^{m_1} \times \dots \times \mathbb{R}^{m_m}$ , and let  $\mathbf{Id}$  be the identity operator on  $\mathcal{K}$ . Define the following operators

$$\mathbf{A} \stackrel{\text{def}}{=} \begin{bmatrix} \partial R & L_1^* & \cdots & L_m^* \\ -L_1 & \partial J_1 & & \\ \vdots & & \ddots & \\ -L_m & & & \partial J_m \end{bmatrix}, \mathbf{B} \stackrel{\text{def}}{=} \begin{bmatrix} \nabla F & & & \\ & \nabla G_1^* & & \\ & & \ddots & \\ & & & \nabla G_m^* \end{bmatrix}, \mathbf{V} \stackrel{\text{def}}{=} \begin{bmatrix} \frac{\text{Id}_n}{\gamma_R} & -L_1^* & \cdots & -L_m^* \\ -L_1 & \frac{\text{Id}_{m_1}}{\gamma_{J_1}} & & \\ \vdots & & \ddots & \\ -L_m & & & \frac{\text{Id}_{m_m}}{\gamma_{J_m}} \end{bmatrix}. \quad (49)$$

Then  $\mathbf{A}$  is maximal monotone,  $\mathbf{B}$  is  $\min\{\beta_F, \beta_{G_1}, \dots, \beta_{G_m}\}$ -cocoercive, and  $\mathbf{V}$  is symmetric and  $\nu$ -positive definite with  $\nu = (1 - \sqrt{\gamma_R \sum_i \gamma_{J_i} \|L_i\|^2}) \min\{\frac{1}{\gamma_R}, \frac{1}{\gamma_{J_1}}, \dots, \frac{1}{\gamma_{J_m}}\}$ . Define  $\mathbf{z}_k = (x_k, v_{1,k}, \dots, v_{m,k})^T$ , then it can be shown that (48) is equivalent to

$$\mathbf{z}_{k+1} = (\mathbf{V} + \mathbf{A})^{-1}(\mathbf{V} - \mathbf{B})\mathbf{z}_k = (\mathbf{Id} + \mathbf{V}^{-1}\mathbf{A})^{-1}(\mathbf{Id} - \mathbf{V}^{-1}\mathbf{B})\mathbf{z}_k. \quad (50)$$



## 5.2. Local convergence analysis

Let  $(x^*, v_1^*, \dots, v_m^*)$  be a Kuhn-Tucker pair. Define the following functions

$$\bar{J}_i^*(v_i) \stackrel{\text{def}}{=} J_i^*(v_i) - \langle v_i, L_i x^* - \nabla G_i^*(v_i^*) \rangle, \quad v_i \in \mathbb{R}^{m_i}, \quad i \in \{1, \dots, m\}, \quad (51)$$

and the Riemannian Hessian of each  $\bar{J}_i^*$ ,

$$H_{\bar{J}_i^*} \stackrel{\text{def}}{=} \gamma_{J_i} \text{P}_{T_{v_i^*}^{J_i^*}} \nabla^2_{\mathcal{M}_{v_i^*}^{J_i^*}} \bar{J}_i^*(v_i^*) \text{P}_{T_{v_i^*}^{J_i^*}} \quad \text{and} \quad W_{\bar{J}_i^*} \stackrel{\text{def}}{=} (\text{Id}_{m_i} + H_{\bar{J}_i^*})^{-1}, \quad i \in \{1, \dots, m\}. \quad (52)$$

For each  $i \in \{1, \dots, m\}$ , owing to Lemma 3.3, we have that  $W_{\bar{J}_i^*}$  is firmly non-expansive if the non-degeneracy condition  $(\text{ND}_m)$  holds. Now suppose

(A'.5)  $F$  locally is  $C^2$ -smooth around  $x^*$  and  $G_i^*$  locally is  $C^2$  around  $v_i^*$ .

Define the restricted Hessian  $H_{G_i^*} \stackrel{\text{def}}{=} \text{P}_{T_{v_i^*}^{J_i^*}} \nabla^2 G_i^*(v_i^*) \text{P}_{T_{v_i^*}^{J_i^*}}$ . Define  $\bar{H}_{G_i^*} \stackrel{\text{def}}{=} \text{Id}_{m_i} - \gamma_{J_i^*} H_{G_i^*}$ ,

$\bar{L}_i \stackrel{\text{def}}{=} \text{P}_{T_{v_i^*}^{J_i^*}} L_i \text{P}_{T_{x^*}^R}$ , and the matrix

$$M_{\text{PD}} \stackrel{\text{def}}{=} \begin{bmatrix} \gamma_{J_1^*} W_{J_1^*} \bar{L}_1 (2W_{\bar{R}} \bar{H}_F - \text{Id}_n) & W_{J_1^*} (\bar{H}_{G_1^*} - 2\gamma_{J_1^*} \gamma_{\bar{R}} \bar{L}_1 W_{\bar{R}} \bar{L}_1^*) & \cdots & -\gamma_{\bar{R}} W_{\bar{R}} \bar{L}_m^* \\ \vdots & \vdots & \ddots & \vdots \\ \gamma_{J_m^*} W_{J_m^*} \bar{L}_m (2W_{\bar{R}} \bar{H}_F - \text{Id}_n) & \cdots & & W_{J_m^*} (\bar{H}_{G_m^*} - 2\gamma_{J_m^*} \gamma_{\bar{R}} \bar{L}_m W_{\bar{R}} \bar{L}_m^*) \end{bmatrix}. \quad (53)$$

Using the same strategy of the proof of Lemma 3.5, one can show that  $M_{\text{PD}}$  is convergent, which again is denoted as  $M_{\text{PD}}^\infty$ , and  $\rho(M_{\text{PD}} - M_{\text{PD}}^\infty) < 1$ .

**Corollary 5.1.** *Consider Algorithm 2 under assumptions (A.1) and (A'.2)-(A'.5). Choose  $\gamma_R, (\gamma_{J_i})_i > 0$  such that*

$$2 \min\{\beta_F, \beta_{G_1}, \dots, \beta_{G_m}\} \min\left\{\frac{1}{\gamma_R}, \frac{1}{\gamma_{J_1}}, \dots, \frac{1}{\gamma_{J_m}}\right\} (1 - \sqrt{\gamma_R \sum_i \gamma_{J_i} \|L_i\|^2}) > 1. \quad (54)$$

Then  $(x_k, v_{1,k}, \dots, v_{m,k}) \rightarrow (x^*, v_1^*, \dots, v_m^*)$ , where  $x^*$  solves  $(\mathcal{P}_P^m)$  and  $(v_1^*, \dots, v_m^*)$  solve  $(\mathcal{P}_D^m)$ . If moreover  $R \in \text{PSF}_{x^*}(\mathcal{M}_{x^*}^R)$  and  $J_i^* \in \text{PSF}_{v_i^*}(\mathcal{M}_{v_i^*}^{J_i^*})$ ,  $i \in \{1, \dots, m\}$ , and the non-degeneracy condition holds

$$\begin{aligned} -\sum_i L_i^* v_i^* - \nabla F(x^*) &\in \text{ri}(\partial R(x^*)) \\ L_i x^* - \nabla G_i^*(v_i^*) &\in \text{ri}(\partial J_i^*(v_i^*)), \quad \forall i \in \{1, \dots, m\}. \end{aligned} \quad (\text{ND}_m)$$

Then,

(i) there exists  $K > 0$  such that for all  $k \geq K$ ,

$$(x_k, v_{1,k}, \dots, v_{m,k}) \in \mathcal{M}_{x^*}^R \times \mathcal{M}_{v_1^*}^{J_1^*} \times \cdots \times \mathcal{M}_{v_m^*}^{J_m^*}.$$

(ii) Given any  $\rho \in ]\rho(M_{\text{PD}} - M_{\text{PD}}^\infty), 1[$ , there exists a  $K$  large enough such that for all  $k \geq K$ ,

$$\|(\text{Id} - M_{\text{PD}}^\infty)(z_k - z^*)\| = O(\rho^{k-K}). \quad (55)$$

If moreover,  $R, J_1^*, \dots, J_m^*$  are locally polyhedral around  $(x^*, v_1^*, \dots, v_m^*)$ , then we have directly have  $\|z_k - z^*\| = O(\rho^{k-K})$ .

**Proof.** Owing to [19], condition (54) guarantees the convergence of the algorithm.

- (i) the identification result follows naturally from Theorem 3.2.
- (ii) the result follows Proposition 3.4, Corollary 3.6 and Theorem 3.7. First, for the update of  $x_k$  of (48), we have

$$\begin{aligned} x_{k+1} - x^* &= W_{\bar{R}} \bar{H}_F(x_k - x^*) - \gamma_R W_{\bar{R}} \sum_i \bar{L}_i^*(v_{i,k} - v_i^*) + o(\|x_k - x^*\| + \gamma_R \sum_i \|L_i\| \|v_{i,k} - v_i^*\|). \end{aligned} \quad (56)$$

Then the update of  $v_{i,k+1}$ , for each  $i = 1, \dots, m$ , similar to (26), we get

$$\begin{aligned} v_{i,k+1} - v_i^* &= (W_{\bar{J}_i^*} \bar{H}_{G_i^*} - (1 + \theta) \gamma_{J_i} \gamma_R W_{\bar{J}_i^*} \bar{L}_i W_{\bar{R}} \bar{L}_i^*)(v_{i,k} - v_i^*) \\ &\quad + (2\gamma_{J_i} W_{\bar{J}_i^*} \bar{L}_i W_{\bar{R}} \bar{H}_F - \gamma_{J_i} W_{\bar{J}_i^*} \bar{L}_i)(x_k - x^*) \\ &\quad + o(\|x_k - x^*\| + \gamma_R \sum_i \|L_i\| \|v_{i,k} - v_i^*\|) + o(\|v_{i,k} - v_i^*\| + \gamma_{J_i} \|L_i\| \|x_k - x^*\|). \end{aligned} \quad (57)$$

Now consider the small  $o$ -terms. For the 2 small  $o$ -terms in (22) and (25). First, let  $a_0, a_1, \dots, a_m$  be  $m + 1$  constants, then we have

$$\sum_{i=0}^m |a_i| = \sqrt{(\sum_{i=0}^m |a_i|)^2} \leq \sqrt{(m+1) \sum_{i=0}^m |a_i|^2} = \sqrt{m+1} \|(a_0, \dots, a_m)^T\|.$$

Denote  $b = \max\{1, \sum_i \sigma_i \|L_i\|, \gamma_R \|L_1\|, \dots, \gamma_R \|L_m\|\}$ , then

$$\begin{aligned} &\sum_i (\|v_{i,k} - v_i^*\| + \sigma_i \|L_i\| \|x_k - x^*\|) + (\|x_k - x^*\| + \gamma_R \sum_i \|L_i\| \|v_{i,k} - v_i^*\|) \\ &\leq 2b(\|x_k - x^*\| + \sum_i \|v_{i,k} - v_i^*\|) \leq 2b\sqrt{m+1} \|z_k - z^*\|. \end{aligned}$$

Combining this with (56) and (57), and ignoring the constants of the small  $o$ -term, we have that the fixed-point iteration (50) is equivalent to

$$z_{k+1} - z^* = M_{\text{PD}}(z_k - z^*) + o(\|z_k - z^*\|).$$

The rest of the proof follows the proof of Theorem 3.7.  $\square$

## 6. Numerical experiments

In this section, we verify our theoretical results on several concrete examples arising from fields including inverse problem, signal/image processing and machine learning.

### 6.1. Examples of partly smooth function

Table 1 provides some examples of partly smooth functions that will be used in this section, whose more details can be found in [34, Section 5] and the references therein.

**Table 1.** Examples of partly smooth functions. For  $x \in \mathbb{R}^n$  and some subset of indices  $b \subset \{1, \dots, n\}$ ,  $x_b$  is the restriction of  $x$  to the entries indexed in  $b$ .  $D_{\text{DIF}}$  stands for the finite differences operator.

Function	Expression	Partial smooth manifold
$\ell_1$ -norm	$\ x\ _1 = \sum_{i=1}^n  x_i $	$\mathcal{M} = T_x = \{z \in \mathbb{R}^n : I_z \subseteq I_x\}, I_x = \{i : x_i \neq 0\}$
$\ell_{1,2}$ -norm	$\sum_{i=1}^m \ x_{b_i}\ $	$\mathcal{M} = T_x = \{z \in \mathbb{R}^n : I_z \subseteq I_x\}, I_x = \{i : x_{b_i} \neq 0\}$
$\ell_\infty$ -norm	$\max_{i=\{1, \dots, n\}}  x_i $	$\mathcal{M} = T_x = \{z \in \mathbb{R}^n : z_{I_x} \in \mathbb{R}\text{sign}(x_{I_x})\}, I_x = \{i :  x_i  = \ x\ _\infty\}$
TV semi-norm	$\ x\ _{\text{TV}} = \ D_{\text{DIF}}x\ _1$	$\mathcal{M} = T_x = \{z \in \mathbb{R}^n : I_{D_{\text{DIF}}z} \subseteq I_{D_{\text{DIF}}x}\}, I_{D_{\text{DIF}}x} = \{i : (D_{\text{DIF}}x)_i \neq 0\}$
Nuclear norm	$\ x\ _* = \sum_{i=1}^r \sigma(x)$	$\mathcal{M} = \{z \in \mathbb{R}^{n_1 \times n_2} : \text{rank}(z) = \text{rank}(x) = r\}, \sigma(x)$ singular values of $x$

The  $\ell_1, \ell_\infty$ -norms and the anisotropic TV semi-norm are polyhedral functions, hence their Riemannian Hessian are simply 0. The  $\ell_{1,2}$ -norm is not polyhedral yet partly smooth relative to a subspace, the nuclear norm is partly smooth relative to the set of fixed-rank matrices, which on the other hand is curved, the Riemannian of these two functions are non-trivial and can be found in [47] and references therein.

## 6.2. Linear inverse problems

Given an object  $x_{\text{ob}} \in \mathbb{R}^n$ , often times we can not access it directly, but through the observation model,

$$b = \mathcal{K}x_{\text{ob}}, \quad (58)$$

where  $b \in \mathbb{R}^m$  is the observation,  $\mathcal{K} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is some linear operator. A more complicated situation is when the observation is contaminated by noise, namely,  $b = \mathcal{K}x_{\text{ob}} + w$ , where  $w \in \mathbb{R}^m$  is the noise.

The operator  $\mathcal{K}$  usually is ill-conditioned or even singular, hence recovering or approximating  $x_{\text{ob}}$  from (58) in general is ill-posed. However, usually some prior knowledge of  $x_{\text{ob}}$  can be available. Thus, a popular approach to recover  $x_{\text{ob}}$  from  $b$  is via regularization, by solving

$$\min_{x \in \mathbb{R}^n} R(x) + J(\mathcal{K}x - b), \quad (59)$$

where

- $R \in \Gamma_0(\mathbb{R}^n)$  is the regularizer based on the prior information, e.g.  $\ell_1, \ell_{1,2}, \ell_\infty$ -norms, nuclear norm;
- $J \in \Gamma_0(\mathbb{R}^m)$  enforces fidelity to the observations. Typically  $J = \iota_0$  when there is no noise, i.e.  $w = 0$ .

Clearly, (59) is a special instance of  $(\mathcal{P}_P)$  with  $F = G^* = 0$ . Thus Algorithm 1 can be applied to solve it.

We consider problem (59) with  $R$  being  $\ell_1, \ell_{1,2}, \ell_\infty$ -norms, and nuclear norm.  $\mathcal{K} \in \mathbb{R}^{m \times n}$  is generated uniformly at random with independent zero-mean standard Gaussian entries. The settings of the experiments are:

**$\ell_1$ -norm**  $(m, n) = (48, 128), \|x_{\text{ob}}\|_0 = 8;$

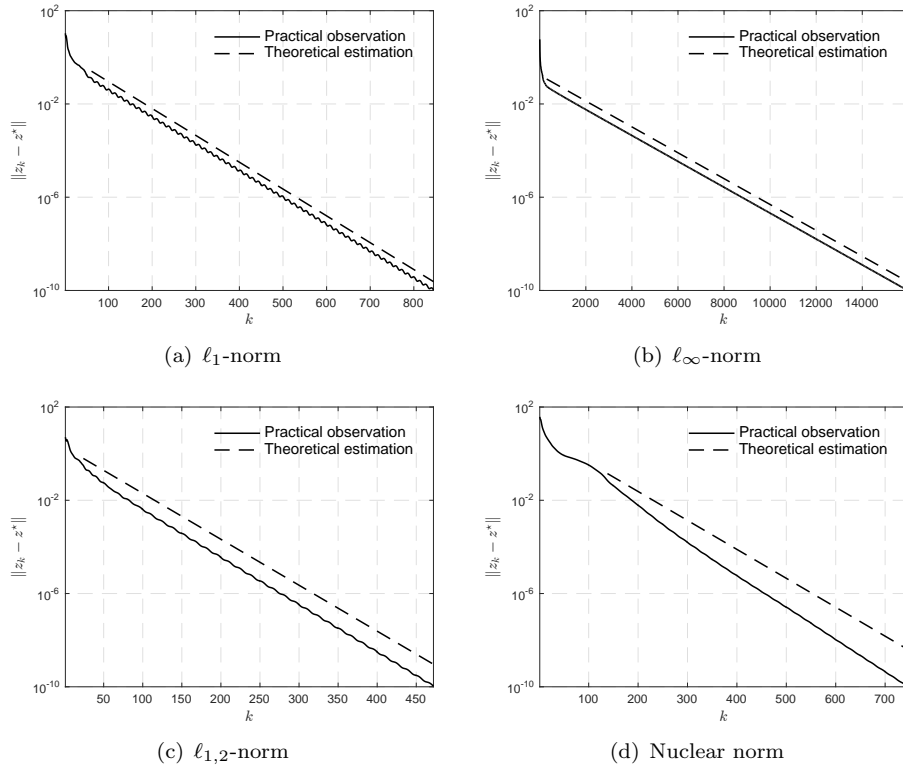
**$\ell_{1,2}$ -norm**  $(m, n) = (48, 128), x_{\text{ob}}$  has 3 non-zero blocks of size 4;

**$\ell_\infty$ -norm**  $(m, n) = (63, 64), |I(x_{\text{ob}})| = 8;$

**Nuclear norm**  $(m, n) = (500, 1024), x_{\text{ob}} \in \mathbb{R}^{32 \times 32}$  and  $\text{rank}(x_{\text{ob}}) = 4.$

Figure 1 displays the profile of  $\|z_k - z^*\|$  as a function of  $k$ , and the starting point of the dashed line is the iteration number at which the active partial smoothness manifold of  $\mathcal{M}_{x^*}^R$  is identified (recall that  $\mathcal{M}_{x^*}^{J^*} = \{0\}$  which is trivially identified from the first iteration). One can easily see that for the  $\ell_1$  and  $\ell_\infty$ -norms, Theorem 3.7 applies and

our estimates are very tight, meaning that the dashed and solid lines has the same slope. For the case of  $\ell_{1,2}$ -norm and nuclear norm, though not optimal, our estimates are very tight.



**Figure 1.** Observed (solid) and predicted (dashed) convergence profiles of Algorithm 1 in terms of  $\|z_k - z^*\|$ . (a)  $\ell_1$ -norm. (b)  $\ell_\infty$ -norm. (c)  $\ell_{1,2}$ -norm. (d) Nuclear norm. The starting point of the dashed line is the iteration at which the active manifold of  $J$  is identified.

### 6.3. Total variation based denoising

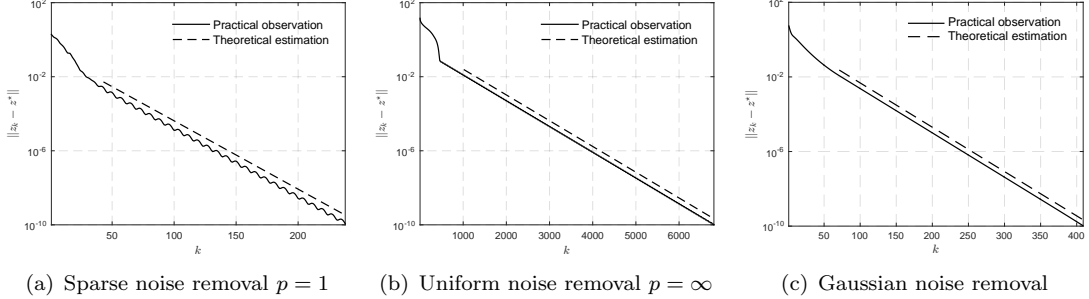
In this part, we consider several total variation based denoising examples, for the first two examples, we consider the observation  $b = x_{\text{ob}} + w$ , where  $x_{\text{ob}}$  is a piecewise-constant vector, and  $w$  is an unknown noise supposed to be either uniform or sparse. The goal is to recover  $x_{\text{ob}}$  from  $b$  using the prior information on  $x_{\text{ob}}$  (*i.e.* piecewise-smooth) and  $w$  (uniform or sparse). To achieve this goal, a popular and natural approach in the signal processing literature is to solve

$$\min_{x \in \mathbb{R}^n} \|D_{\text{DIF}} x\|_1 \quad \text{subject to} \quad \|b - x\|_p \leq \tau, \quad (60)$$

where  $p = +\infty$  for uniform noise, and  $p = 1$  for sparse noise, and  $\tau > 0$  is a parameter depending on the noise level.

Problem (60) can also be formulated into the form of  $(\mathcal{P}_P)$ . Indeed, one can take  $R = \iota_{\|b - \cdot\|_p \leq \tau}$ ,  $J = \|\cdot\|_1$ ,  $F = G^* = 0$ , and  $L = D_{\text{DIF}}$  is the finite difference operator (with appropriate boundary conditions). The proximity operators of  $R$  and  $J$  can be computed easily. Clearly, both two indicator functions are polyhedral, and their proximal operator are simple to compute.

For both examples, we set  $n = 128$  and  $x_{\text{ob}}$  is such that  $D_{\text{DIF}}x_{\text{ob}}$  has 8 nonzero entries. For  $p = +\infty$ ,  $w$  is generated uniformly in  $[-1, 1]$ , and for  $p = 1$ ,  $w$  is sparse with 16 nonzero entries. The corresponding local convergence profiles are depicted in Figure 2(a)-(b). Owing to polyhedrality, our rate predictions are again optimal.



**Figure 2.** Observed (solid) and predicted (dashed) convergence profiles of Primal-Dual (2) in terms of  $\|z_k - z^*\|$ . (a) Sparse noise removal. (b) Uniform noise removal. (c) Gaussian noise removal. The starting point of the dashed line is the iteration at which the active manifold of  $J$  is identified.

We also consider an underdetermined linear regression problem  $b = \mathcal{K}x_{\text{ob}} + w$ . We assume that the vector  $x_{\text{ob}}$  is group sparse and each non-zero group is piecewise constant. This regression problem can then be approached by solving

$$\min_{x \in \mathbb{R}^n} \mu_1 \|x\|_{1,2} + \frac{1}{2} \|\mathcal{K}x - b\|^2 + \mu_2 \|D_{\text{DIF}}x\|_1,$$

where  $\mu_1, \mu_2 > 0$ ,  $\|\cdot\|_{1,2}$  promotes group sparsity, and  $\|D_{\text{DIF}}\cdot\|_1$  promotes piece-wise constancy. This is again in the form of  $(\mathcal{P}_P)$ , where  $R = \mu_1 \|\cdot\|_{1,2}$ ,  $F = \frac{1}{2} \|\mathcal{K}\cdot - b\|^2$ ,  $J = \mu_2 \|\cdot\|_1$ ,  $G^* = 0$ , and  $L = D_{\text{DIF}}$ . For this example, we set  $x_{\text{ob}} \in \mathbb{R}^{128}$  with 2 piecewise constant non-zeros blocks of size 8. The result is shown in Figure 2(c), the estimate is again very sharp.

#### 6.4. Choices of $\theta$ and $\gamma_J, \gamma_R$

In this part, we present a comparison on different choices of  $\theta$  and  $\gamma_J, \gamma_R$  to see their influences on the finite identification and local linear convergence rate. Two examples are consider for these comparisons, problem (59) with  $R$  being  $\ell_1$ -norm and  $\ell_{1,2}$ -norm.

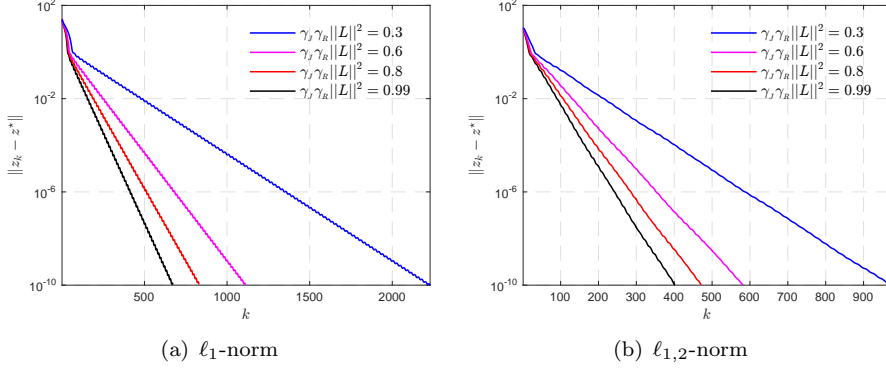
**Fixed  $\theta$**  For this comparison, we fix  $\theta = 1$  and consider 4 different cases of  $\gamma_J \gamma_R \|L\|^2$ :

$$\gamma_J \gamma_R \|L\|^2 \in \{0.3, 0.6, 0.8, 0.99\},$$

with  $\gamma_J = \gamma_R$ . The result is shown in Figure 3, and we have the following observations:

- (i) The smaller the value of  $\gamma_J \gamma_R \|L\|^2$ , the slower the iteration converges;
- (ii) Bigger value of  $\gamma_R$  leads to faster identification (since  $J^*$  is globally  $C^2$ -smooth, so only the identification of  $R$  for this case).

**Fixed  $\gamma_J \gamma_R \|L\|^2$**  Now we turn to the opposite direction, fix  $\gamma_J \gamma_R \|L\|^2$  and change  $\theta$ . In the test, we fixed  $\gamma_J \gamma_R \|L\|^2 = 0.9$  and  $\gamma_J = \gamma_R$ , 5 different choices of  $\theta$  are considered,



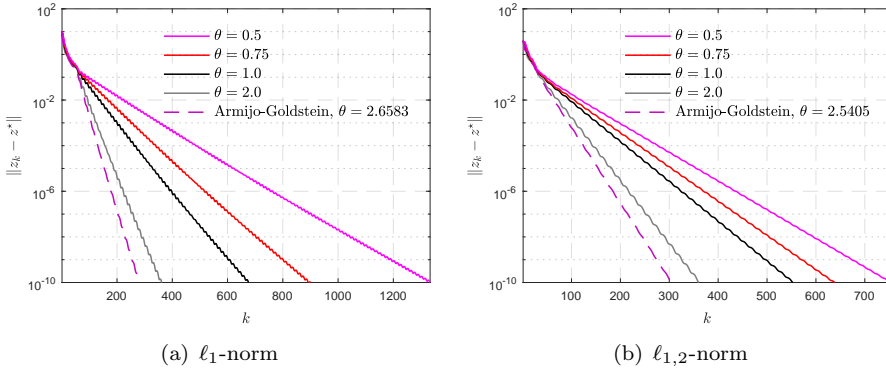
**Figure 3.** Comparison of the choice of  $\gamma_J, \gamma_R$  when  $\theta$  is fixed.

which are

$$\theta \in \{0.5, 0.75, 1.0, 2.0\},$$

plus the last one with Armijo–Goldstein-rule for updating  $\theta$  adaptively. Although there is no convergence guarantee for  $\theta = 2.0$ , in the test it converges and we choose to put it here as an illustration of the effects of  $\theta > 1$ . The result is shown in Figure 4, and we have the following observations

- (i) Similar to the previous one, the smaller the value of  $\theta$ , the slower the iteration converges. Also, the Armijo–Goldstein-rule is the fastest of all.
- (ii) Interestingly, the value of  $\theta$  has no impacts to the identification of the iteration.



**Figure 4.** Comparison of the choice of  $\theta$  when  $\gamma_J, \gamma_R$  are fixed.

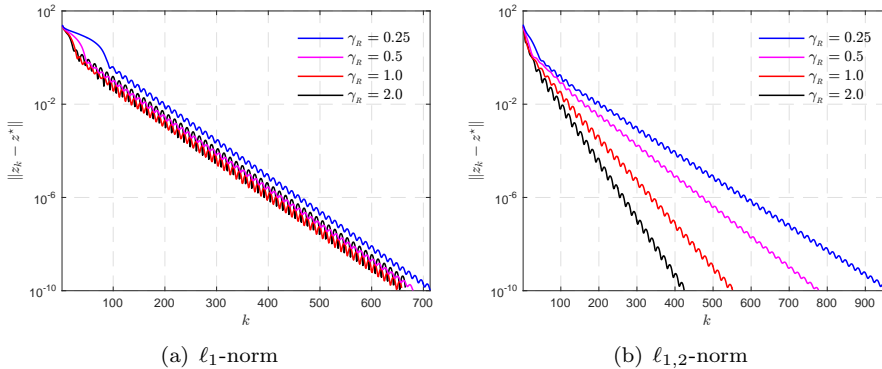
**Fixed  $\theta$  and  $\gamma_J \gamma_R$**  For the above comparisons, we fix  $\gamma_J = \gamma_R$ , so for this comparison, we compare the different choices of them. We fix  $\theta = 1$  and  $\gamma_J \gamma_R \|L\|^2 = 0.99$ , then we choose

$$\gamma_R \in \{0.25, 0.5, 1, 2\} \text{ and } \gamma_J = \frac{0.99}{\gamma_R \|L\|^2}.$$

Figure 5 shows the comparison result, we also have two observations:

- (i) For the  $\ell_1$ -norm, since both functions are polyhedral, local convergence rate are the same for all choices of  $\gamma_R$ , see (29) for the expression of the rate. The only

- difference is the identification speed,  $\gamma_R = 0.25$  gives the slowest identification, however it uses almost the same number of iterations reaching the given accuracy.
- (ii) For the  $\ell_{1,2}$ -norm, on the other hand, the choice of  $\gamma_R$  affects both the identification and local convergence rate. It can be observed that bigger  $\gamma_R$  leads to faster local rate, however, it does not mean that the bigger the better. In fact, too big value will slow down the convergence.



**Figure 5.** Comparison of fixed  $\theta$  and  $\gamma_J \gamma_R$ , but varying  $\gamma_R$ .

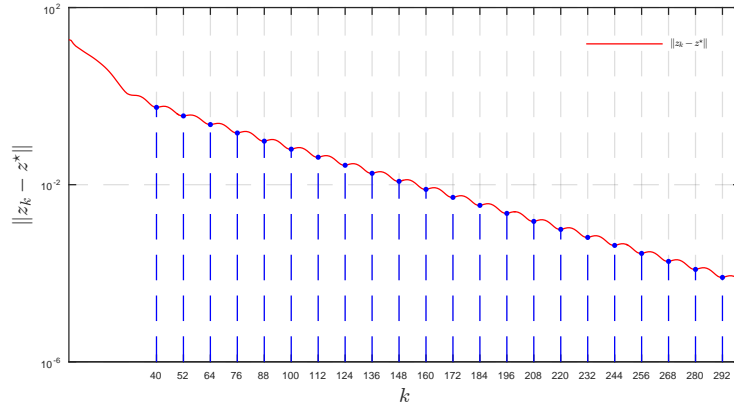
### 6.5. Oscillation of the method

We dedicate the last part of this section to demonstrate the oscillation behaviour of the Primal–Dual splitting method when dealing with polyhedral functions. As we have seen from the above experiments, oscillation of  $\|z_k - z^*\|$  happens for all examples whose involved functions  $R, J^*$  are polyhedral, even for the non-polyhedral  $\ell_{1,2}$ -norm (for the  $\ell_\infty$ -norm, the oscillation is not visible since the oscillation period is too small compared to the number of iterations).

Now to verify our discussion in Section 4, we consider problem (59) with  $R$  being  $\ell_1$ -norm, and the result is shown in Figure 6. As revealed in (41), the argument of the leading eigenvalue of  $M_{PD} - M_{PD}^\infty$  is controlled by  $\theta \gamma_J \gamma_R$ , so is the oscillation period. Therefore, the value  $\gamma_J \gamma_R$  is tuned such that the oscillation period is an integer, and  $\pi/\omega = 12$  for the example we tested. Figure 6 shows graphically the observed oscillation, apparently the oscillation pattern coincides well with the theoretical estimation.

## 7. Discussion and conclusion

In this paper, we studied local convergence properties of a class of Primal–Dual splitting methods when the involved non-smooth functions are moreover partly smooth. In particular, we demonstrated that these methods identify the active manifolds in finite time and then converge locally linearly at a rate that we characterized precisely. We also built connections of the presented result to our previous work on Forward–Backward splitting and Douglas–Rachford splitting/ADMM. Though we focused on one class of Primal–Dual splitting methods, there are other Primal–Dual splitting schemes, such as those in [12,18,25], to which our analysis and conclusions can be straightforwardly extended.



**Figure 6.** Oscillation behaviour of the Primal–Dual splitting method when dealing with polyhedral functions.

## Acknowledgments

This work has been partly supported by the European Research Council (ERC project SIGMA-Vision). JF was partly supported by Institut Universitaire de France. The authors would like to thank Russell Luke for helpful discussions. JL was partly supported by Leverhulme Trust project “Breaking the non-convexity barrier”, the EPSRC grant “EP/M00483X/1”, EPSRC centre “EP/N014588/1”, the Cantab Capital Institute for the Mathematics of Information, and the Global Alliance project “Statistical and Mathematical Theory of Imaging”.

## References

- [1] P.-A. Absil, R. Mahony, and J. Trumpf. An extrinsic look at the Riemannian Hessian. In *Geometric Science of Information*, pages 361–368. Springer, 2013.
- [2] A. Agarwal, S. Negahban, and M. J. Wainwright. Fast global convergence of gradient methods for high-dimensional statistical recovery. *The Annals of Statistics*, 40(5):2452–2482, 2012.
- [3] K.-J. Arrow, L. Hurwicz, H. Uzawa, H.-B. Chenery, S.-M. Johnson, S. Karlin, and T. Marschak. *Studies in linear and non-linear programming*. 1959.
- [4] H. Bauschke, J.Y.B. Cruz, T.A. Nghia, H.M. Phan, and X. Wang. The rate of linear convergence of the Douglas–Rachford algorithm for subspaces is the cosine of the Friedrichs angle. *J. of Approx. Theo.*, 185(63–79), 2014.
- [5] H. H. Bauschke, J. Y. Bello Cruz, T. T. A. Nghia, H. M. Pha, and X. Wang. Optimal rates of linear convergence of relaxed alternating projections and generalized Douglas–Rachford methods for two subspaces. *Numerical Algorithms*, 73(1):33–76, 2016.
- [6] H. H. Bauschke and P. L. Combettes. *Convex analysis and monotone operator theory in Hilbert spaces*. Springer, 2011.
- [7] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- [8] Daniel Boley. Local linear convergence of the alternating direction method of multipliers on quadratic or linear programs. *SIAM Journal on Optimization*, 23(4):2183–2207, 2013.
- [9] R. I. Boţ, and A. Hendrich. Convergence analysis for a Primal-Dual monotone + skew splitting algorithm with applications to total variation minimization. Tech. Rep., arXiv:1211.1706, 2012.
- [10] R. I. Boţ, E. R. Csetnek, A. Heinrich, and C. Hendrich. On the convergence rate im-



- provement of a primal-dual splitting algorithm for solving monotone inclusion problems. *Mathematical Programming*, 150(2):251–279, 2015.
- [11] K. Bredies and D. A. Lorenz. Linear convergence of iterative soft-thresholding. *Journal of Fourier Analysis and Applications*, 14(5-6):813–837, 2008.
- [12] L. M. Briceno-Arias and P. L. Combettes. A monotone+ skew splitting model for composite monotone inclusions in duality. *SIAM Journal on Optimization*, 21(4):1230–1250, 2011.
- [13] A. Chambolle and T. Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *Journal of Mathematical Imaging and Vision*, 40(1):120–145, 2011.
- [14] A. Chambolle and T. Pock. On the ergodic convergence rates of a first-order primal-dual algorithm. *Mathematical Programming*, pages 1–35, 2015.
- [15] I. Chavel. *Riemannian geometry: a modern introduction*, volume 98. Cambridge University Press, 2006.
- [16] P. Chen, J. Huang, and X. Zhang. A primal-dual fixed point algorithm for convex separable minimization with applications to image restoration. *Inverse Problems*, 29(2):025011, 2013.
- [17] P. L. Combettes, L. Condat, J.-C. Pesquet, and B. C. Vu. A Forward-Backward view of some Primal-Dual optimization methods in image recovery. In *Image Processing (ICIP), 2014 IEEE International Conference on*, pages 4141–4145. IEEE, 2014.
- [18] P. L. Combettes and J. C. Pesquet. Primal-dual splitting algorithm for solving inclusions with mixtures of composite, lipschitzian, and parallel-sum type monotone operators. *Set-Valued and variational analysis*, 20(2):307–330, 2012.
- [19] P. L. Combettes and B. C. Vũ. Variable metric Forward-Backward splitting with applications to monotone inclusions in duality. *Optimization*, 63(9):1289–1318, 2014.
- [20] P. L. Combettes and I. Yamada. Compositions and convex combinations of averaged non-expansive operators. *Journal of Mathematical Analysis and Applications*, 425(1):55–70, 2015.
- [21] L. Condat. A primal-dual splitting method for convex optimization involving Lipschitzian, proximable and linear composite terms. *Journal of Optimization Theory and Applications*, pages 1–20, 2012.
- [22] D. Davis. Convergence rate analysis of primal-dual splitting schemes. *SIAM Journal on Optimization*, 25(3):1912–1943, 2015.
- [23] L. Demanet and X. Zhang. Eventual linear convergence of the Douglas-Rachford iteration for basis pursuit. *Mathematics of Computation*, 85(297):209–238, 2016.
- [24] J. Douglas and H. H. Rachford. On the numerical solution of heat conduction problems in two and three space variables. *Transactions of the American mathematical Society*, 82(2):421–439, 1956.
- [25] E. Esser, X. Zhang, and T. F. Chan. A general framework for a class of first order primal-dual algorithms for convex optimization in imaging science. *SIAM Journal on Imaging Sciences*, 3(4):1015–1046, 2010.
- [26] E. Hale, W. Yin, and Y. Zhang. Fixed-point continuation for  $\ell_1$ -minimization: Methodology and convergence. *SIAM Journal on Optimization*, 19(3):1107–1130, 2008.
- [27] W. L. Hare and A. S. Lewis. Identifying active constraints via partial smoothness and prox-regularity. *Journal of Convex Analysis*, 11(2):251–266, 2004.
- [28] B. He and X. Yuan. Convergence analysis of primal-dual algorithms for a saddle-point problem: from contraction perspective. *SIAM Journal on Imaging Sciences*, 5(1):119–149, 2012.
- [29] K. Hou, Z. Zhou, A. M.-C. So, and Z. Q. Luo. On the linear convergence of the proximal gradient method for trace norm regularization. In *Advances in Neural Information Processing Systems*, pages 710–718, 2013.
- [30] J. M. Lee. *Smooth manifolds*. Springer, 2003.
- [31] A. S. Lewis. Active sets, nonsmoothness, and sensitivity. *SIAM Journal on Optimization*, 13(3):702–725, 2003.

- [32] J. Liang, J. Fadili, and G. Peyré. Local linear convergence of Forward–Backward under partial smoothness. In *Advances in Neural Information Processing Systems*, pages 1970–1978, 2014.
- [33] J. Liang, J. Fadili, and G. Peyré. Convergence rates with inexact non-expansive operators. *Mathematical Programming*, 159(1):403–434, September 2016.
- [34] J. Liang, J. Fadili, and G. Peyré. Activity identification and local linear convergence of Forward–Backward-type methods. *SIAM Journal on Optimization*, 27(1):408–437, 2017.
- [35] J. Liang, J. Fadili, and G. Peyré. Local convergence properties of Douglas–Rachford and alternating direction method of multipliers. *Journal of Optimization Theory and Applications*, 172(3):874–913, 2017.
- [36] J. Liang, J. Fadili, G. Peyré, and R. Luke. Activity identification and local linear convergence of Douglas–Rachford/ADMM under partial smoothness. In *International Conference on Scale Space and Variational Methods in Computer Vision*, pages 642–653. Springer, 2015.
- [37] P. L. Lions and B. Mercier. Splitting algorithms for the sum of two nonlinear operators. *SIAM Journal on Numerical Analysis*, 16(6):964–979, 1979.
- [38] S. A. Miller and J. Malick. Newton methods for nonsmooth convex minimization: connections among-lagrangian, riemannian newton and sqp methods. *Mathematical programming*, 104(2-3):609–633, 2005.
- [39] J. J. Moreau. Proximité et dualité dans un espace hilbertien. *Bulletin de la Société mathématique de France*, 93:273–299, 1965.
- [40] N. Ogura and I. Yamada. Non-strictly convex minimization over the fixed point set of an asymptotically shrinking non-expansive mapping. *Numerical Functional Analysis and Optimization*, 23:113–137, 2002.
- [41] R. T. Rockafellar. *Convex analysis*, volume 28. Princeton university press, 1997.
- [42] R. T. Rockafellar and R. Wets. *Variational analysis*, volume 317. Springer Verlag, 1998.
- [43] J. R. Silvester. Determinants of block matrices. *The Mathematical Gazette*, pages 460–467, 2000.
- [44] T. Sun, R. Barrio, H. Jiang, and L. Cheng. Local linear convergence of a primal-dual algorithm for the augmented convex models. *Journal of Scientific Computing*, 69(3):1301–1315, 2016.
- [45] S. Tao, D. Boley, and S. Zhang. Local linear convergence of ISTA and FISTA on the lasso problem. *SIAM Journal on Optimization*, 26(1):313–336, 2016.
- [46] P. Tseng. A modified Forward–Backward splitting method for maximal monotone mappings. *SIAM Journal on Control and Optimization*, 38(2):431–446, 2000.
- [47] S. Vaïter, C. Deledalle, J. M. Fadili, G. Peyré, and C. Dossal. The degrees of freedom of partly smooth regularizers. *Annals of the Institute of Statistical Mathematics*, 2015. to appear.
- [48] B. C. Vũ. A splitting algorithm for dual monotone inclusions involving cocoercive operators. *Advances in Computational Mathematics*, pages 1–15, 2011.
- [49] S. J. Wright. Identifiable surfaces in constrained optimization. *SIAM Journal on Control and Optimization*, 31(4):1063–1079, 1993.
- [50] L. M. Briceño-Arias and D. Davis. Forward–Backward-half forward algorithm with non self-adjoint linear operators for solving monotone inclusions. *arXiv preprint arXiv:1703.03436*, 2017.

## Appendix A. Riemannian Geometry

Let  $\mathcal{M}$  be a  $C^2$ -smooth embedded submanifold of  $\mathbb{R}^n$  around a point  $x$ . With some abuse of terminology, we shall state  $C^2$ -manifold instead of  $C^2$ -smooth embedded submanifold of  $\mathbb{R}^n$ . The natural embedding of a submanifold  $\mathcal{M}$  into  $\mathbb{R}^n$  permits to define a Riemannian structure and to introduce geodesics on  $\mathcal{M}$ , and we simply say

$\mathcal{M}$  is a Riemannian manifold. We denote respectively  $\mathcal{T}_{\mathcal{M}}(x)$  and  $\mathcal{N}_{\mathcal{M}}(x)$  the tangent and normal space of  $\mathcal{M}$  at point near  $x$  in  $\mathcal{M}$ .

**Exponential map** Geodesics generalize the concept of straight lines in  $\mathbb{R}^n$ , preserving the zero acceleration characteristic, to manifolds. Roughly speaking, a geodesic is locally the shortest path between two points on  $\mathcal{M}$ . We denote by  $\mathbf{g}(t; x, h)$  the value at  $t \in \mathbb{R}$  of the geodesic starting at  $\mathbf{g}(0; x, h) = x \in \mathcal{M}$  with velocity  $\dot{\mathbf{g}}(t; x, h) = \frac{d\mathbf{g}}{dt}(t; x, h) = h \in \mathcal{T}_{\mathcal{M}}(x)$  (which is uniquely defined). For every  $h \in \mathcal{T}_{\mathcal{M}}(x)$ , there exists an interval  $I$  around 0 and a unique geodesic  $\mathbf{g}(t; x, h) : I \rightarrow \mathcal{M}$  such that  $\mathbf{g}(0; x, h) = x$  and  $\dot{\mathbf{g}}(0; x, h) = h$ . The mapping  $\text{Exp}_x : \mathcal{T}_{\mathcal{M}}(x) \rightarrow \mathcal{M}$ ,  $h \mapsto \text{Exp}_x(h) = \mathbf{g}(1; x, h)$  is called Exponential map. Given  $x, x' \in \mathcal{M}$ , the direction  $h \in \mathcal{T}_{\mathcal{M}}(x)$  we are interested in is such that  $\text{Exp}_x(h) = x' = \mathbf{g}(1; x, h)$ .

**Parallel translation** Given two points  $x, x' \in \mathcal{M}$ , let  $\mathcal{T}_{\mathcal{M}}(x), \mathcal{T}_{\mathcal{M}}(x')$  be their corresponding tangent spaces. Define  $\tau : \mathcal{T}_{\mathcal{M}}(x) \rightarrow \mathcal{T}_{\mathcal{M}}(x')$  the parallel translation along the unique geodesic joining  $x$  to  $x'$ , which is isomorphism and isometry w.r.t. the Riemannian metric.

**Riemannian gradient and Hessian** For a vector  $v \in \mathcal{N}_{\mathcal{M}}(x)$ , the Weingarten map of  $\mathcal{M}$  at  $x$  is the operator  $\mathfrak{W}_x(\cdot, v) : \mathcal{T}_{\mathcal{M}}(x) \rightarrow \mathcal{T}_{\mathcal{M}}(x)$  defined by  $\mathfrak{W}_x(\cdot, v) = -\text{P}_{\mathcal{T}_{\mathcal{M}}(x)} dV[h]$  where  $V$  is any local extension of  $v$  to a normal vector field on  $\mathcal{M}$ . The definition is independent of the choice of the extension  $V$ , and  $\mathfrak{W}_x(\cdot, v)$  is a symmetric linear operator which is closely tied to the second fundamental form of  $\mathcal{M}$ , see [15, Proposition II.2.1].

Let  $J$  be a real-valued function which is  $C^2$  along the  $\mathcal{M}$  around  $x$ . The covariant gradient of  $J$  at  $x' \in \mathcal{M}$  is the vector  $\nabla_{\mathcal{M}} J(x') \in \mathcal{T}_{\mathcal{M}}(x')$  defined by

$$\langle \nabla_{\mathcal{M}} J(x'), h \rangle = \frac{d}{dt} J(\text{P}_{\mathcal{M}}(x' + th)) \Big|_{t=0}, \quad \forall h \in \mathcal{T}_{\mathcal{M}}(x'),$$

where  $\text{P}_{\mathcal{M}}$  is the projection operator onto  $\mathcal{M}$ . The covariant Hessian of  $J$  at  $x'$  is the symmetric linear mapping  $\nabla_{\mathcal{M}}^2 J(x')$  from  $\mathcal{T}_{\mathcal{M}}(x')$  to itself which is defined as

$$\langle \nabla_{\mathcal{M}}^2 J(x') h, h \rangle = \frac{d^2}{dt^2} J(\text{P}_{\mathcal{M}}(x' + th)) \Big|_{t=0}, \quad \forall h \in \mathcal{T}_{\mathcal{M}}(x'). \quad (\text{A1})$$

This definition agrees with the usual definition using geodesics or connections [38]. Now assume that  $\mathcal{M}$  is a Riemannian embedded submanifold of  $\mathbb{R}^n$ , and that a function  $J$  has a  $C^2$ -smooth restriction on  $\mathcal{M}$ . This can be characterized by the existence of a  $C^2$ -smooth extension (representative) of  $J$ , *i.e.* a  $C^2$ -smooth function  $\tilde{J}$  on  $\mathbb{R}^n$  such that  $\tilde{J}$  agrees with  $J$  on  $\mathcal{M}$ . Thus, the Riemannian gradient  $\nabla_{\mathcal{M}} J(x')$  is also given by

$$\nabla_{\mathcal{M}} J(x') = \text{P}_{\mathcal{T}_{\mathcal{M}}(x')} \nabla \tilde{J}(x'), \quad (\text{A2})$$

and  $\forall h \in \mathcal{T}_{\mathcal{M}}(x')$ , the Riemannian Hessian reads

$$\begin{aligned} \nabla_{\mathcal{M}}^2 J(x') h &= \text{P}_{\mathcal{T}_{\mathcal{M}}(x')} d(\nabla_{\mathcal{M}} J)(x')[h] = \text{P}_{\mathcal{T}_{\mathcal{M}}(x')} d(x' \mapsto \text{P}_{\mathcal{T}_{\mathcal{M}}(x')} \nabla_{\mathcal{M}} \tilde{J})[h] \\ &= \text{P}_{\mathcal{T}_{\mathcal{M}}(x')} \nabla^2 \tilde{J}(x') h + \mathfrak{W}_{x'}(h, \text{P}_{\mathcal{N}_{\mathcal{M}}(x')} \nabla \tilde{J}(x')), \end{aligned} \quad (\text{A3})$$

where the last equality comes from [1, Theorem 1]. When  $\mathcal{M}$  is an affine or linear subspace of  $\mathbb{R}^n$ , then obviously  $\mathcal{M} = x + \mathcal{T}_{\mathcal{M}}(x)$ , and  $\mathfrak{W}_{x'}(h, \mathbb{P}_{\mathcal{N}_{\mathcal{M}}(x')} \nabla \tilde{J}(x')) = 0$ , hence (A3) reduces to  $\nabla_{\mathcal{M}}^2 J(x') = \mathbb{P}_{\mathcal{T}_{\mathcal{M}}(x')} \nabla^2 \tilde{J}(x') \mathbb{P}_{\mathcal{T}_{\mathcal{M}}(x')}$ ; See [15,30] for more materials on differential and Riemannian manifolds.

We have the following proposition characterizing the parallel translation and the Riemannian Hessian of two close points in  $\mathcal{M}$ .

**Lemma A.1.** *Let  $x, x'$  be two close points in  $\mathcal{M}$ , denote  $\mathcal{T}_{\mathcal{M}}(x), \mathcal{T}_{\mathcal{M}}(x')$  be the tangent spaces of  $\mathcal{M}$  at  $x, x'$  respectively, and  $\tau : \mathcal{T}_{\mathcal{M}}(x') \rightarrow \mathcal{T}_{\mathcal{M}}(x)$  be the parallel translation along the unique geodesic joining from  $x$  to  $x'$ , then for the parallel translation we have, given any bounded vector  $v \in \mathbb{R}^n$*

$$(\tau \mathbb{P}_{\mathcal{T}_{\mathcal{M}}(x')} - \mathbb{P}_{\mathcal{T}_{\mathcal{M}}(x)})v = o(\|v\|). \quad (\text{A4})$$

The Riemannian Taylor expansion of  $J \in C^2(\mathcal{M})$  at  $x$  for  $x'$  reads,

$$\tau \nabla_{\mathcal{M}} J(x') = \nabla_{\mathcal{M}} J(x) + \nabla_{\mathcal{M}}^2 J(x) \mathbb{P}_{\mathcal{T}_{\mathcal{M}}(x)}(x' - x) + o(\|x' - x\|). \quad (\text{A5})$$

**Proof.** See [34, Lemma B.1 and B.2]. □

**Lemma A.2.** *Let  $\mathcal{M}$  be a  $C^2$ -smooth manifold,  $\bar{x} \in \mathcal{M}$ ,  $R \in \text{PSF}_{\bar{x}}(\mathcal{M})$  and  $\bar{u} \in \partial R(\bar{x})$ . Let  $\tilde{R}$  be a smooth representative of  $R$  on  $\mathcal{M}$  near  $x$ , then given any  $h \in T_{\bar{x}}$ ,*

(i) *when  $\mathcal{M}$  is a general smooth manifold, if there holds  $\bar{u} \in \text{ri}(\partial R(\bar{x}))$ , define the function  $\bar{R}(x) = R(x) - \langle x, \bar{u} \rangle$ , then*

$$\langle \mathbb{P}_{T_{\bar{x}}} \nabla_{\mathcal{M}}^2 \bar{R}(\bar{x}) \mathbb{P}_{T_{\bar{x}}} h, h \rangle \geq 0. \quad (\text{A6})$$

(ii) *if  $\mathcal{M}$  is affine/linear, then we have directly,*

$$\langle \mathbb{P}_{T_{\bar{x}}} \nabla_{\mathcal{M}}^2 \tilde{R}(\bar{x}) \mathbb{P}_{T_{\bar{x}}} h, h \rangle \geq 0. \quad (\text{A7})$$

**Proof.** See [34, Lemma 4.3]. □