

# An accelerated proximal gradient method for multiobjective optimization\*

Hiroki Tanabe<sup>1</sup>, Ellen H. Fukuda<sup>2</sup>, and Nobuo Yamashita<sup>2</sup>

<sup>1</sup>Yahoo Japan Corporation

<sup>2</sup>Kyoto University

{tanabe.hiroki.45n@kyoto-u.jp}, {ellen,nobuo}@i.kyoto-u.ac.jp

## Abstract

This paper presents an accelerated proximal gradient method for multiobjective optimization, in which each objective function is the sum of a continuously differentiable, convex function and a closed, proper, convex function. Extending first-order methods for multiobjective problems without scalarization has been widely studied, but providing accelerated methods with accurate proofs of convergence rates remains an open problem. Our proposed method is a multiobjective generalization of the accelerated proximal gradient method, also known as the Fast Iterative Shrinkage-Thresholding Algorithm (FISTA), for scalar optimization. The key to this successful extension is solving a subproblem with terms exclusive to the multiobjective case. This approach allows us to demonstrate the global convergence rate of the proposed method ( $O(1/k^2)$ ), using a merit function to measure the complexity. Furthermore, we present an efficient way to solve the subproblem via its dual representation, and we confirm the validity of the proposed method through some numerical experiments.

## 1 Introduction

Multiobjective optimization consists in minimizing (or maximizing) more than one objective function at once under possible constraints. In general, there is no single point that minimizes all objective functions simultaneously, so the concept of *Pareto optimality* becomes essential. We call a point Pareto optimal if there is no other point with the same or smaller objective function values and with at least one objective function value being strictly smaller.

One of the most popular strategies for solving multiobjective optimization problems is the *scalarization approach* [23, 24, 46]. It converts the original multiobjective problem into another, which has a parametrized scalar-valued objective function. If each objective function in the multiobjective optimization problem is convex, the converted single objective optimization problems typically become convex optimization. However, it can be challenging to choose the appropriate parameters (or weights) in advance. For example, Marler and Arora [30] discussed such difficulties in the weighted sum method. Another approach, which does not use scalarization, is based on metaheuristics [22] but lacks a theoretical proof of convergence to Pareto solutions.

To overcome those drawbacks, many descent algorithms for multiobjective optimization problems have been developed recently [21]. These algorithms decrease all objective functions at each iteration, offer the advantages of not requiring a priori parameter selection, and provide convergence guarantees under reasonable assumptions. For instance, Fliege and Svaiter [18] proposed the steepest descent method for differentiable multiobjective optimization problems. Other examples include the projected gradient [20, 26, 47], Newton's [17, 25], trust-region [12], and conjugate gradient methods [29]. Descent methods for infinite-dimensional vector optimization problems have also been studied, including the proximal point [7] and the inertial forward-backward methods [5]. Among these, methods that use only the first-order derivatives of the objective functions, such as the steepest descent and the projected gradient methods, are called *first-order methods*. Another well-known multiobjective first-order method is the proximal gradient [40], which works for composite problems, i.e., with each objective

\*A previous version of this manuscript can be seen in the Department of Applied Mathematics and Physics, Kyoto University's technical report system ([http://www.amp.i.kyoto-u.ac.jp/tecprep/ps\\_file/2022/2022-001.pdf](http://www.amp.i.kyoto-u.ac.jp/tecprep/ps_file/2022/2022-001.pdf)).

being the sum of a differentiable function and a convex but not necessarily differentiable one. This algorithm, as well as the steepest descent, is known to converge to Pareto solutions with rate  $O(1/k)$  [19, 41].

On the other hand, there are many studies related to the acceleration of single-objective first-order methods. After being established by Nesterov [34], researchers developed various accelerated schemes. In particular, the Fast Iterative Shrinkage-Thresholding Algorithm (FISTA) [3], an accelerated version of the proximal gradient method, has contributed to a wide range of research fields, including image and signal processing. These methods may increase the objective function values in some iterations, but overall they are known to converge faster than the original descent methods, both theoretically and experimentally.

However, in the multiobjective case, the studies associated with accelerated algorithms are still insufficient [16, 44]. In 2020, El Moudden and El Mouatasim [16] proposed an accelerated diagonal steepest descent method for multiobjective optimization, a natural extension of Nesterov's accelerated method for single-objective problems. They proved the global convergence rate of the algorithm ( $O(1/k^2)$ ) under the assumption that the sequence of the Lagrange multipliers of the subproblems is eventually fixed. Nevertheless, this assumption is restrictive because it indicates that the approach is essentially the same as the (single-objective) Nesterov's method, only applied to the minimization of a weighted sum of the objective functions.

Here, we propose a genuine accelerated proximal gradient method for multiobjective optimization. As it is usual, in each iteration, we solve a convex (scalar-valued) subproblem. While the accelerated and non-accelerated algorithms solve the same subproblem in the single-objective case, the subproblem of our accelerated method has terms that are not included in the non-accelerated version. However, we can ignore these terms in the single-objective case, and thus we can regard our proposed method as a generalization of FISTA. Moreover, under more natural assumptions, we prove the proposed method's global convergence rate ( $O(1/k^2)$ ) by using a merit function [42] to measure the complexity.

Furthermore, having the practical computational efficiency in mind, we derive a dual of the subproblem, which is convex and differentiable. Such a dual problem turns out to be easier to solve than the original one, especially when the number of objective functions is smaller than the dimension of the decision variables. We can also reconstruct the original subproblem's solution directly from the dual optimum. In addition, we implement the whole algorithm using this dual problem and confirm its effectiveness with numerical experiments.

The outline of this paper is as follows. In Section 2, we introduce some notations and concepts used in this paper. Section 3 recalls the proximal gradient method for multiobjective optimization proposed in [40]. We present the proposed accelerated proximal gradient method for multiobjective optimization in Section 4 and analyze its  $O(1/k^2)$  convergence rate in Section 5. Moreover, Section 6 introduces an efficient way to solve the subproblem via its dual form. Finally, we report some numerical results for test problems in Section 7, demonstrating that the proposed method is faster than the one without acceleration.

## 2 Preliminaries

All over this work, for any natural number  $d$ ,  $\mathbf{R}^d$  denotes the  $d$ -dimensional real space,  $\mathbf{R}_+^d \subseteq \mathbf{R}^d$  designates the nonnegative orthant of  $\mathbf{R}^d$ , i.e.,  $\mathbf{R}_+^d := \{v \in \mathbf{R}^d \mid v_i \geq 0, i = 1, \dots, d\}$ , and  $\Delta^d$  represents the standard simplex in  $\mathbf{R}^d$  given by

$$\Delta^d := \left\{ \lambda \in \mathbf{R}_+^d \mid \sum_{i=1}^d \lambda_i = 1 \right\}. \quad (1)$$

Then, we can consider the partial orders induced by  $\mathbf{R}_+^d$ : for all  $v^1, v^2 \in \mathbf{R}^d$ ,  $v^1 \leq v^2$  (alternatively,  $v^2 \geq v^1$ ) if  $v^2 - v^1 \in \mathbf{R}_+^d$  and  $v^1 < v^2$  (alternatively,  $v^2 > v^1$ ) if  $v^2 - v^1 \in \text{int } \mathbf{R}_+^d$ . In other words,  $v^1 \leq v^2$  and  $v^1 < v^2$  stand for  $v_i^1 \leq v_i^2$  and  $v_i^1 < v_i^2$  for all  $i = 1, \dots, d$ , respectively. Moreover, let  $\langle \cdot, \cdot \rangle$  be the Euclidean inner product in  $\mathbf{R}^d$ , i.e.,  $\langle u, v \rangle := \sum_{i=1}^d u_i v_i$ , and let  $\|\cdot\|$  be the Euclidean norm, i.e.,  $\|u\| := \sqrt{\langle u, u \rangle}$ . Furthermore, we define the  $\ell_1$ -norm and the  $\ell_\infty$ -norm by  $\|u\|_1 := \sum_{i=1}^d |u_i|$  and  $\|u\|_\infty := \max_{i=1, \dots, d} |u_i|$ , respectively. We now recall the obvious equality related to norm and inner product:

$$\|b - a\|^2 + 2\langle b - a, a - c \rangle = \|b - c\|^2 - \|a - c\|^2 \quad \text{for all } a, b, c \in \mathbf{R}^n. \quad (2)$$

On the other hand, for a closed, proper and convex function  $h: \mathbf{R}^n \rightarrow \mathbf{R} \cup \{+\infty\}$ , we call  $\eta \in \mathbf{R}^n$  a subgradient of  $h$  at  $x \in \mathbf{R}^n$  if

$$h(y) \geq h(x) + \langle \eta, y - x \rangle \quad \text{for all } y \in \mathbf{R}^n, \quad (3)$$

and we write  $\partial h(x)$  the subdifferential of  $h$  at  $x$ , i.e., the set of all subgradients of  $h$  at  $x$ . In addition, the subdifferential for a vector-valued function is the direct product of the subdifferentials of each component. We

also define the *Moreau envelope* or *Moreau-Yosida regularization* [33, 45] of  $h$  by

$$\mathcal{M}_h(x) := \min_{y \in \mathbf{R}^n} \left[ h(y) + \frac{1}{2} \|x - y\|^2 \right]. \quad (4)$$

The minimization problem in (4) has a unique solution because of the strong convexity of its objective function. We call this solution the *proximal operator* and write it as

$$\mathbf{prox}_h(x) := \operatorname{argmin}_{y \in \mathbf{R}^n} \left[ h(y) + \frac{1}{2} \|x - y\|^2 \right]. \quad (5)$$

**Remark 2.1.** (i) [2, Theorem 6.24] If  $h$  is the indicator function of a nonempty set  $S \subseteq \mathbf{R}^n$ , i.e.,

$$\chi_S(x) = \begin{cases} 0, & x \in S, \\ +\infty, & x \notin S, \end{cases} \quad (6)$$

then the proximal operator reduces to the projection onto  $S$ .

- (ii) [2, Theorem 6.42] The proximal operator of a closed, proper, and convex function  $h$  is non-expansive, i.e.,  $\|\mathbf{prox}_h(x) - \mathbf{prox}_h(y)\| \leq \|x - y\|$ . In other words,  $\mathbf{prox}_h$  is 1-Lipschitz continuous.
- (iii) [2, Theorem 6.60] Even if a closed, proper, and convex function  $h$  is non-differentiable, its Moreau envelope  $\mathcal{M}_h$  has a 1-Lipschitz continuous gradient as follows:  $\nabla \mathcal{M}_h(x) = x - \mathbf{prox}_h(x)$ .

We now focus on the following multiobjective optimization problem:

$$\min_{x \in \mathbf{R}^n} F(x) \quad (7)$$

with a vector-valued function  $F: \mathbf{R}^n \rightarrow (\mathbf{R} \cup \{+\infty\})^m$  with  $F := (F_1, \dots, F_m)^\top$ . We assume that each component  $F_i: \mathbf{R}^n \rightarrow \mathbf{R} \cup \{+\infty\}$  is defined by  $F_i(x) := f_i(x) + g_i(x)$  for all  $i = 1, \dots, m$  with convex and continuously differentiable functions  $f_i: \mathbf{R}^n \rightarrow \mathbf{R}, i = 1, \dots, m$  and closed, proper and convex functions  $g_i: \mathbf{R}^n \rightarrow \mathbf{R} \cup \{+\infty\}, i = 1, \dots, m$ . We also suppose that each  $\nabla f_i$  is Lipschitz continuous with constant  $L_i > 0$  and define  $L := \max_{i=1, \dots, m} L_i$ . From the so-called descent lemma [4, Proposition A.24], we have  $f_i(p) - f_i(q) \leq \langle \nabla f_i(q), p - q \rangle + (L/2) \|p - q\|^2$  for all  $p, q \in \mathbf{R}^n$  and  $i = 1, \dots, m$ , which gives

$$\begin{aligned} F_i(p) - F_i(r) &= f_i(p) - f_i(q) + g_i(p) + f_i(q) - F_i(r) \\ &\leq \langle \nabla f_i(q), p - q \rangle + g_i(p) + f_i(q) - F_i(r) + \frac{L}{2} \|p - q\|^2 \end{aligned} \quad (8)$$

for all  $p, q, r \in \mathbf{R}^n$  and  $i = 1, \dots, m$ .

Now, we introduce some concepts used in the multiobjective optimization problem (7). Recall that

$$X^* := \{x^* \in \mathbf{R}^n \mid \text{There does not exist } x \in \mathbf{R}^n \text{ such that } F(x) < F(x^*)\} \quad (9)$$

is the set of *weakly Pareto optimal* points for (7). We also define the effective domain of  $F$  by  $\operatorname{dom} F := \{x \in \mathbf{R}^n \mid F_i(x) < +\infty \text{ for all } i = 1, \dots, m\}$ , and we write the level set of  $F$  on  $c \in \mathbf{R}^m$  as

$$\mathcal{L}_F(c) := \{x \in \mathbf{R}^n \mid F(x) \leq c\}. \quad (10)$$

In addition, we express the image of  $A \subseteq \mathbf{R}^n$  and the inverse image of  $B \subseteq (\mathbf{R} \cup \{+\infty\})^m$  under  $F$  as  $F(A) := \{F(x) \in \mathbf{R}^m \mid x \in A\}$  and  $F^{-1}(B) := \{x \in \mathbf{R}^n \mid F(x) \in B\}$ , respectively.

Finally, let us recall the merit function  $u_0: \mathbf{R}^n \rightarrow \mathbf{R} \cup \{+\infty\}$  proposed in [42]:

$$u_0(x) := \sup_{z \in \mathbf{R}^n} \min_{i=1, \dots, m} [F_i(x) - F_i(z)], \quad (11)$$

which returns zero at optimal solutions and strictly positive values otherwise. The following theorem shows that  $u_0$  is a merit function in the Pareto sense.

**Theorem 2.1.** [42, Theorem 3.1] Let  $u_0$  be defined by (11). Then, we get  $u_0(x) \geq 0$  for all  $x \in \mathbf{R}^n$ . Moreover,  $x \in \mathbf{R}^n$  is weakly Pareto optimal for (7) if and only if  $u_0(x) = 0$ .

Note that when  $m = 1$ , we have  $u_0(x) = F_1(x) - F_1^*$ , where  $F_1^*$  is the optimal objective value. This is clearly a merit function for scalar-valued optimization.

### 3 Proximal gradient methods for multiobjective optimization

Let us now recall the proximal gradient method for (7), an extension of the classical proximal gradient method, proposed by Tanabe, Fukuda, and Yamashita [40]. We explain how to generate the sequence of iterates, and afterward, we show the algorithm and its convergence rate.

For given  $x \in \text{dom } F$  and  $\ell > 0$ , we consider the following minimization problem:

$$\min_{z \in \mathbf{R}^n} \varphi_\ell(z; x), \quad (12)$$

where  $\varphi_\ell(z; x) := \max_{i=1, \dots, m} [\langle \nabla f_i(x), z - x \rangle + g_i(z) - g_i(x)] + (\ell/2) \|z - x\|^2$ . The convexity of  $g_i$  implies that  $z \mapsto \varphi_\ell(z; x)$  is strongly convex, so the problem (12) always has a unique solution. Let us write such a solution as  $p_\ell(x)$  and let  $\theta_\ell(x)$  be its optimal function value, i.e.,

$$p_\ell(x) := \underset{z \in \mathbf{R}^n}{\text{argmin}} \varphi_\ell(z; x) \quad \text{and} \quad \theta_\ell(x) := \min_{z \in \mathbf{R}^n} \varphi_\ell(z; x). \quad (13)$$

The following proposition shows that  $p_\ell(x)$  and  $\theta_\ell(x)$  helps to characterize the weak Pareto optimality of (7).

**Proposition 3.1.** *Let  $p_\ell$  and  $\theta_\ell$  be defined by (13). Then, the statements below hold.*

- (i) *The following three conditions are equivalent: (a)  $x$  is weakly Pareto optimal; (b)  $p_\ell(x) = x$ ; (c)  $\theta_\ell(x) = 0$ .*
- (ii) *The mappings  $p_\ell$  and  $\theta_\ell$  are both continuous.*

*Proof.* It is clear from [40, Lemma 3.2] and the convexity of  $f_i$ . □

From Proposition 3.1, we can treat  $\|p_\ell(x) - x\|_\infty < \varepsilon$  for some  $\varepsilon > 0$  as a stopping criteria. Moreover, if  $\ell > L/2$  then we have  $F_i(p_\ell(x)) \leq F_i(x)$  for all  $x \in \text{dom } F$  and  $i = 1, \dots, m$  [41]. Now, we state below the proximal gradient method for (7).

---

**Algorithm 1** Proximal gradient method for multiobjective optimization [40]

---

**Input:**  $x^0 \in \text{dom } F$ ,  $\ell > L/2$ ,  $\varepsilon > 0$

**Output:**  $x^*$ : A weakly Pareto optimal point

- 1:  $k \leftarrow 0$
  - 2: **loop**
  - 3:  $x^{k+1} \leftarrow p_\ell(x^k)$ , where  $p_\ell$  is defined by (13)
  - 4: **if**  $\|x^{k+1} - x^k\|_\infty \geq \varepsilon$  **then return**  $x^{k+1}$
  - 5: **end if**
  - 6:  $k \leftarrow k + 1$
  - 7: **end loop**
- 

When  $\ell \geq L$ , Algorithm 1 is known to generate  $\{x^k\}$  such that  $\{u_0(x^k)\}$  converges to zero with rate  $O(1/k)$  under the following assumption. Note that this assumption is not particularly strong, as suggested in [41, Remark 5.2].

**Assumption 3.1.** [41, Assumption 5.1] *Let  $X^*$  and  $\mathcal{L}_F$  be defined by (9) and (10), respectively. Then, for all  $x \in \mathcal{L}_F(F(x^0))$ , there exists  $x^* \in X^*$  such that  $F(x^*) \leq F(x)$  and*

$$R := \sup_{F^* \in F(X^* \cap \mathcal{L}_F(F(x^0)))} \inf_{z \in F^{-1}(\{F^*\})} \|z - x^0\|^2 < +\infty. \quad (14)$$

**Theorem 3.1.** [41, Theorem 5.2] *Assume that  $\ell \geq L$ . Then, under Assumption 3.1, Algorithm 1 generates a sequence  $\{x^k\}$  such that  $u_0(x^k) \leq (\ell R)/(2k)$  for all  $k \geq 1$ .*

At the end of this section, we note some remarks about Algorithm 1.

**Remark 3.1.** (i) *Since  $x \in \text{dom } F$  implies  $p_\ell(x) \in \text{dom } F$ , Algorithm 1 is well-defined.*

(ii) If  $g_i = 0$ , Algorithm 1 corresponds to the steepest descent method [18]:

$$x^{k+1} := \min_{z \in \mathbf{R}^n} \left[ \max_{i=1, \dots, m} \langle \nabla f_i(x^k), z - x^k \rangle + \frac{\ell}{2} \|z - x^k\|^2 \right].$$

On the other hand, when  $f_i = 0$ , it matches the proximal point method [7]:

$$x^{k+1} := \min_{z \in \mathbf{R}^n} \left\{ \max_{i=1, \dots, m} [g_i(z) - g_i(x)] + \frac{\ell}{2} \|z - x^k\|^2 \right\}.$$

Furthermore, when  $g_i$  is the indicator function (6) of a convex set  $S \subseteq \mathbf{R}^n$ , it coincides with the projected gradient method [26]:

$$x^{k+1} := \min_{z \in S - x^k} \left[ \max_{i=1, \dots, m} \langle \nabla f_i(x^k), z - x^k \rangle + \frac{\ell}{2} \|z - x^k\|^2 \right].$$

(iii) When it is difficult to estimate the Lipschitz constant  $L$ , we can set the initial value of  $\ell$  appropriately. Then, at each iteration we increase  $\ell$  by multiplying it with some prespecified scalar, until  $F_i(p_\ell(x^k)) - F_i(x^k) \leq \theta_\ell(x^k)$  is satisfied for all  $i = 1, \dots, m$ . If  $L$  is finite, the number of times that  $\ell$  is increased is at most a constant.

## 4 An accelerated proximal gradient method for multiobjective optimization

This section proposes an accelerated version of the proximal gradient method for multiobjective optimization. Similarly to the non-accelerated version given in the last section, a subproblem is considered in each iteration. More specifically, the proposed method solves the following subproblem for given  $x \in \text{dom } F$ ,  $y \in \mathbf{R}^n$ , and  $\ell \geq L$ :

$$\min_{z \in \mathbf{R}^n} \varphi_\ell^{\text{acc}}(z; x, y), \quad (15)$$

where

$$\varphi_\ell^{\text{acc}}(z; x, y) := \max_{i=1, \dots, m} [\langle \nabla f_i(y), z - y \rangle + g_i(z) + f_i(y) - F_i(x)] + \frac{\ell}{2} \|z - y\|^2.$$

Note that when  $y = x$ , (15) is reduced to the subproblem (12) of the proximal gradient method. Note also that when  $m = 1$ , the subproblem becomes

$$\min_{z \in \mathbf{R}^n} \langle \nabla f_1(y), z - y \rangle + g_1(z) + \frac{\ell}{2} \|z - y\|^2, \quad (16)$$

which is the subproblem of the single-objective FISTA [3]. The distinctive feature of our proposal (15) is the term  $f_i(y) - F_i(x)$ , whereas the easy analogy from the single-objective subproblem (16) is

$$\min_{z \in \mathbf{R}^n} \max_{i=1, \dots, m} [\langle \nabla f_i(y), z - y \rangle + g_i(z)] + \frac{\ell}{2} \|z - y\|^2. \quad (17)$$

By putting such a term, the inside of the max operator approximates  $F_i(z) - F_i(x)$  rather than  $F_i(z) - F_i(y)$ . This is a negligible difference in the single-objective case, but deeply affects the proof in the multi-objective case.

Since  $g_i$  is convex for all  $i = 1, \dots, m$ ,  $z \mapsto \varphi_\ell^{\text{acc}}(z; x, y)$  is strongly convex. Thus, the subproblem (15) has a unique optimal solution  $p_\ell^{\text{acc}}(x, y)$  and takes the optimal function value  $\theta_\ell^{\text{acc}}(x, y)$ , i.e.,

$$p_\ell^{\text{acc}}(x, y) := \operatorname{argmin}_{z \in \mathbf{R}^n} \varphi_\ell^{\text{acc}}(z; x, y) \quad \text{and} \quad \theta_\ell^{\text{acc}}(x, y) := \min_{z \in \mathbf{R}^n} \varphi_\ell^{\text{acc}}(z; x, y). \quad (18)$$

Moreover, the optimality condition of (15) implies that for all  $x \in \text{dom } F$  and  $y \in \mathbf{R}^n$  there exists  $\eta(x, y) \in \partial g(p_\ell^{\text{acc}}(x, y))$  and a Lagrange multiplier  $\lambda(x, y) \in \mathbf{R}^m$  such that

$$\sum_{i=1}^m \lambda_i(x, y) [\nabla f_i(y) + \eta_i(x, y)] = -\ell [p_\ell^{\text{acc}}(x, y) - y] \quad (19a)$$

$$\lambda(x, y) \in \Delta^m, \quad \lambda_j(x, y) = 0 \quad \text{for all } j \notin \mathcal{I}(x, y), \quad (19b)$$

where  $\Delta^m$  denotes the standard simplex (1) and

$$\mathcal{I}(x, y) := \operatorname{argmax}_{i=1, \dots, m} [\langle \nabla f_i(y), p_\ell^{\text{acc}}(x, y) - y \rangle + g_i(p_\ell^{\text{acc}}(x, y)) + f_i(y) - F_i(x)]. \quad (20)$$

Now, we introduce a relation useful for the subsequent analysis.

**Lemma 4.1.** *Let  $p_\ell^{\text{acc}}$  and  $\theta_\ell^{\text{acc}}$  be defined by (18). Then, we have*

$$\begin{aligned} & -\frac{\ell}{2} [\|p_\ell^{\text{acc}}(x, y) - z\|^2 - \|y - z\|^2] \\ & \geq \theta_\ell^{\text{acc}}(x, y) + \sum_{i=1}^m \lambda_i(x, y) [\langle \nabla f_i(y), y - z \rangle - g_i(z) - f_i(y) + F_i(x)] \end{aligned}$$

for all  $x, z \in \operatorname{dom} F$  and  $y \in \mathbf{R}^n$ .

*Proof.* Let  $x, z \in \operatorname{dom} F$  and  $y \in \mathbf{R}^n$ . From (19a) and the definition (3) of the subgradient, we get

$$\begin{aligned} & -\ell \langle p_\ell^{\text{acc}}(x, y) - y, p_\ell^{\text{acc}}(x, y) - z \rangle \\ & \geq \sum_{i=1}^m \lambda_i(x, y) [\langle \nabla f_i(y), p_\ell^{\text{acc}}(x, y) - z \rangle + g_i(p_\ell^{\text{acc}}(x, y)) - g_i(z)] \\ & = \sum_{i=1}^m \lambda_i(x, y) [\langle \nabla f_i(y), p_\ell^{\text{acc}}(x, y) - y \rangle + g_i(p_\ell^{\text{acc}}(x, y)) + f_i(y) - F_i(x)] \\ & \quad + \sum_{i=1}^m \lambda_i(x, y) [\langle \nabla f_i(y), y - z \rangle - g_i(z) - f_i(y) + F_i(x)] \\ & = \max_{i=1, \dots, m} [\langle \nabla f_i(y), p_\ell^{\text{acc}}(x, y) - y \rangle + g_i(p_\ell^{\text{acc}}(x, y)) + f_i(y) - F_i(x)] \\ & \quad + \sum_{i=1}^m \lambda_i(x, y) [\langle \nabla f_i(y), y - z \rangle - g_i(z) - f_i(y) + F_i(x)], \end{aligned}$$

where the second equality comes from (19b) and (20). Adding  $(\ell/2)\|p_\ell^{\text{acc}}(x, y) - y\|^2$  to both sides and the definition (18) of  $p_\ell^{\text{acc}}$  and  $\theta_\ell^{\text{acc}}$  lead to

$$\begin{aligned} & -\frac{\ell}{2} [2\langle p_\ell^{\text{acc}}(x, y) - y, p_\ell^{\text{acc}}(x, y) - z \rangle - \|p_\ell^{\text{acc}}(x, y) - y\|^2] \\ & \geq \theta_\ell^{\text{acc}}(x, y) + \sum_{i=1}^m \lambda_i(x, y) [\langle \nabla f_i(y), y - z \rangle - g_i(z) - f_i(y) + F_i(x)]. \end{aligned}$$

The left-hand side of this inequality is equal to  $-(\ell/2)[2\langle p_\ell^{\text{acc}}(x, y) - y, y - z \rangle + \|p_\ell^{\text{acc}}(x, y) - y\|^2]$ . Hence, applying (2) with  $(a, b, c) := (y, z, p_\ell^{\text{acc}}(x, y))$ , we get the desired inequality.  $\square$

We also note that by taking  $z = y$  in the objective function of (15), we have

$$\theta_\ell^{\text{acc}}(x, y) \leq \varphi_\ell^{\text{acc}}(y; x, y) = \max_{i=1, \dots, m} \{F_i(y) - F_i(x)\} \quad (21)$$

for all  $x \in \operatorname{dom} F$  and  $y \in \mathbf{R}^n$ . Moreover, from (8) with  $p = z, q = y, r = x$ , and the fact that  $\ell \geq L$ , it follows that

$$\theta_\ell^{\text{acc}}(x, y) \geq \max_{i=1, \dots, m} \{F_i(p_\ell^{\text{acc}}(x, y)) - F_i(x)\}$$

for all  $x \in \operatorname{dom} F$  and  $y \in \mathbf{R}^n$ . We now characterize weak Pareto optimality in terms of the mappings  $p_\ell^{\text{acc}}$  and  $\theta_\ell^{\text{acc}}$ , similarly to Proposition 3.1 for the proximal gradient method.

**Proposition 4.1.** *Let  $p_\ell^{\text{acc}}(x, y)$  and  $\theta_\ell^{\text{acc}}(x, y)$  be defined by (18). Then, the statements below hold.*

- (i) *The following three conditions are equivalent: (a)  $y \in \mathbf{R}^n$  is weakly Pareto optimal for (7); (b)  $p_\ell^{\text{acc}}(x, y) = y$  for some  $x \in \mathbf{R}^n$ ; (c)  $\theta_\ell^{\text{acc}}(x, y) = \max_{i=1, \dots, m} [F_i(y) - F_i(x)]$  for some  $x \in \mathbf{R}^n$ .*

(ii) The mappings  $p_\ell^{\text{acc}}$  and  $\theta_\ell^{\text{acc}}$  are locally Hölder continuous with exponent  $1/2$  and locally Lipschitz continuous, respectively, i.e., for any bounded set  $W \subseteq \mathbf{R}^n$ , there exists  $M_p > 0$  and  $M_\theta > 0$  such that

$$\begin{aligned} \|p_\ell^{\text{acc}}(\hat{x}, \hat{y}) - p_\ell^{\text{acc}}(\check{x}, \check{y})\| &\leq M_p \|(\hat{x}, \hat{y}) - (\check{x}, \check{y})\|^{1/2}, \\ |\theta_\ell^{\text{acc}}(\hat{x}, \hat{y}) - \theta_\ell^{\text{acc}}(\check{x}, \check{y})| &\leq M_\theta \|(\hat{x}, \hat{y}) - (\check{x}, \check{y})\| \end{aligned}$$

for all  $\hat{x}, \hat{y}, \check{x}, \check{y} \in W$ .

*Proof.* Part (i) : From (21) and the fact that  $\theta_\ell^{\text{acc}}(x, y) = \varphi_\ell^{\text{acc}}(p_\ell^{\text{acc}}(x, y); x, y)$ , the equivalence between (b) and (c) is apparent. Now, let us show that (a) and (b) are equivalent. When  $y$  is weakly Pareto optimal, we can immediately see from Proposition 3.1 that  $p_\ell^{\text{acc}}(x, y) = p_\ell(y) = y$  by letting  $x = y$ . Conversely, suppose that  $p_\ell^{\text{acc}}(x, y) = y$  for some  $x \in \mathbf{R}^n$ . Let  $z \in \mathbf{R}^n$  and  $\alpha \in (0, 1)$ . The optimality of  $p_\ell^{\text{acc}}(x, y) = y$  for (15) gives

$$\begin{aligned} \max_{i=1, \dots, m} [F_i(y) - F_i(x)] &\leq \varphi_\ell^{\text{acc}}(y + \alpha(z - y); x, y) \\ &= \max_{i=1, \dots, m} [\langle \nabla f_i(y), \alpha(z - y) \rangle + g_i(y + \alpha(z - y)) + f_i(y) - F_i(x)] \\ &\quad + \frac{\ell}{2} \|\alpha(z - y)\|^2. \end{aligned}$$

Thus, from the convexity of  $f_i$ , we get

$$\max_{i=1, \dots, m} [F_i(y) - F_i(x)] \leq \max_{i=1, \dots, m} [F_i(y + \alpha(z - y)) - F_i(x)] + \frac{\ell}{2} \|\alpha(z - y)\|^2.$$

Moreover, the convexity of  $F_i$  yields

$$\begin{aligned} \max_{i=1, \dots, m} [F_i(y) - F_i(x)] &\leq \max_{i=1, \dots, m} [\alpha F_i(z) + (1 - \alpha)F_i(y) - F_i(x)] + \frac{\ell}{2} \|\alpha(z - y)\|^2 \\ &\leq \alpha \max_{i=1, \dots, m} [F_i(z) - F_i(y)] + \max_{i=1, \dots, m} \{F_i(y) - F_i(x)\} + \frac{\ell}{2} \|\alpha(z - y)\|^2. \end{aligned}$$

Therefore, we get

$$\max_{i=1, \dots, m} [F_i(z) - F_i(y)] \geq -\frac{\ell\alpha}{2} \|z - y\|^2.$$

Taking  $\alpha \searrow 0$ , we obtain  $\max_{i=1, \dots, m} [F_i(z) - F_i(y)] \geq 0$ , which implies the weak Pareto optimality of  $y$ .

Part (ii) : Take  $\hat{x}, \hat{y}, \check{x}, \check{y} \in W$ . Adding the two inequalities of Lemma 4.1 with  $(x, y, z) := (\hat{x}, \hat{y}, p_\ell^{\text{acc}}(\check{x}, \check{y}))$ ,  $(\check{x}, \check{y}, p_\ell^{\text{acc}}(\hat{x}, \hat{y}))$  gives

$$\begin{aligned} &-\ell \|p_\ell^{\text{acc}}(\hat{x}, \hat{y}) - p_\ell^{\text{acc}}(\check{x}, \check{y})\|^2 + \frac{\ell}{2} \|p_\ell^{\text{acc}}(\check{x}, \check{y}) - \hat{y}\|^2 + \frac{\ell}{2} \|p_\ell^{\text{acc}}(\hat{x}, \hat{y}) - \check{y}\|^2 \\ &\geq \theta_\ell^{\text{acc}}(\hat{x}, \hat{y}) + \theta_\ell^{\text{acc}}(\check{x}, \check{y}) \\ &\quad + \sum_{i=1}^m \lambda_i(\hat{x}, \hat{y}) [\langle \nabla f_i(\hat{y}), \hat{y} - p_\ell^{\text{acc}}(\check{x}, \check{y}) \rangle - g_i(p_\ell^{\text{acc}}(\check{x}, \check{y})) - f_i(\hat{y}) + F_i(\hat{x})] \\ &\quad + \sum_{i=1}^m \lambda_i(\check{x}, \check{y}) [\langle \nabla f_i(\check{y}), \check{y} - p_\ell^{\text{acc}}(\hat{x}, \hat{y}) \rangle - g_i(p_\ell^{\text{acc}}(\hat{x}, \hat{y})) - f_i(\check{y}) + F_i(\check{x})]. \end{aligned}$$

From the definition (18) of  $p_\ell^{\text{acc}}$  and  $\theta_\ell^{\text{acc}}$  and (19b), we have

$$\begin{aligned}
& -\ell \|p_\ell^{\text{acc}}(\hat{x}, \hat{y}) - p_\ell^{\text{acc}}(\tilde{x}, \tilde{y})\|^2 \\
& \geq \sum_{i=1}^m \lambda_i(\tilde{x}, \tilde{y}) [\langle \nabla f_i(\hat{y}), p_\ell^{\text{acc}}(\hat{x}, \hat{y}) - \hat{y} \rangle + g_i(p_\ell^{\text{acc}}(\hat{x}, \hat{y})) + f_i(\hat{y}) - F_i(\hat{x})] \\
& \quad + \sum_{i=1}^m \lambda_i(\hat{x}, \hat{y}) [\langle \nabla f_i(\tilde{y}), p_\ell^{\text{acc}}(\tilde{x}, \tilde{y}) - \tilde{y} \rangle + g_i(p_\ell^{\text{acc}}(\tilde{x}, \tilde{y})) + f_i(\tilde{y}) - F_i(\tilde{x})] \\
& \quad + \sum_{i=1}^m \lambda_i(\hat{x}, \hat{y}) [\langle \nabla f_i(\hat{y}), \hat{y} - p_\ell^{\text{acc}}(\tilde{x}, \tilde{y}) \rangle - g_i(p_\ell^{\text{acc}}(\tilde{x}, \tilde{y})) - f_i(\hat{y}) + F_i(\hat{x})] \\
& \quad + \sum_{i=1}^m \lambda_i(\tilde{x}, \tilde{y}) [\langle \nabla f_i(\tilde{y}), \tilde{y} - p_\ell^{\text{acc}}(\hat{x}, \hat{y}) \rangle - g_i(p_\ell^{\text{acc}}(\hat{x}, \hat{y})) - f_i(\tilde{y}) + F_i(\tilde{x})] \\
& \quad - \frac{\ell}{2} \left[ \|p_\ell^{\text{acc}}(\hat{x}, \hat{y}) - \tilde{y}\|^2 - \|p_\ell^{\text{acc}}(\hat{x}, \hat{y}) - \hat{y}\|^2 \right. \\
& \quad \quad \left. + \|p_\ell^{\text{acc}}(\tilde{x}, \tilde{y}) - \hat{y}\|^2 - \|p_\ell^{\text{acc}}(\tilde{x}, \tilde{y}) - \tilde{y}\|^2 \right] \\
& = \sum_{i=1}^m \lambda_i(\hat{x}, \hat{y}) [\langle \nabla f_i(\hat{y}), \hat{y} - \tilde{y} \rangle + \langle \nabla f_i(\hat{y}) - \nabla f_i(\tilde{y}), \tilde{y} - p_\ell^{\text{acc}}(\tilde{x}, \tilde{y}) \rangle \\
& \quad \quad \quad - f_i(\hat{y}) + f_i(\tilde{y}) + F_i(\hat{x}) - F_i(\tilde{x})] \\
& \quad + \sum_{i=1}^m \lambda_i(\tilde{x}, \tilde{y}) [\langle \nabla f_i(\tilde{y}), \tilde{y} - \hat{y} \rangle + \langle \nabla f_i(\tilde{y}) - \nabla f_i(\hat{y}), \hat{y} - p_\ell^{\text{acc}}(\hat{x}, \hat{y}) \rangle \\
& \quad \quad \quad - f_i(\tilde{y}) + f_i(\hat{y}) + F_i(\tilde{x}) - F_i(\hat{x})] \\
& \quad - \ell \langle p_\ell^{\text{acc}}(\hat{x}, \hat{y}) - p_\ell^{\text{acc}}(\tilde{x}, \tilde{y}), \hat{y} - \tilde{y} \rangle.
\end{aligned}$$

Thus, (19b) and Cauchy-Schwarz inequalities applied in each inner product that appears in the right-hand side of the above expression imply

$$\begin{aligned}
& -\ell \|p_\ell^{\text{acc}}(\hat{x}, \hat{y}) - p_\ell^{\text{acc}}(\tilde{x}, \tilde{y})\|^2 \\
& \geq -2 \max_{i=1, \dots, m} \|\nabla f_i(\hat{y})\| \|\hat{y} - \tilde{y}\| \\
& \quad - \left[ \|\hat{y} - p_\ell^{\text{acc}}(\hat{x}, \hat{y})\| + \|\tilde{y} - p_\ell^{\text{acc}}(\tilde{x}, \tilde{y})\| \right] \max_{i=1, \dots, m} \|\nabla f_i(\hat{y}) - \nabla f_i(\tilde{y})\| \\
& \quad - 2 \max_{i=1, \dots, m} |f_i(\hat{y}) - f_i(\tilde{y})| - 2 \max_{i=1, \dots, m} |F_i(\hat{x}) - F_i(\tilde{x})| \\
& \quad - \ell \|p_\ell^{\text{acc}}(\hat{x}, \hat{y}) - p_\ell^{\text{acc}}(\tilde{x}, \tilde{y})\| \|\hat{y} - \tilde{y}\|.
\end{aligned}$$

Let us now show that each term of the right-hand side of the above inequality is bounded by a positive constant multiple of  $-\|\hat{x} - \tilde{x}\|$  or  $-\|\hat{y} - \tilde{y}\|$ . The first term is direct because the boundedness of  $W$  implies  $\max_{i=1, \dots, m} \|\nabla f_i(\hat{y})\| < +\infty$ . Since  $W$  is bounded and the objective function of (15) is strongly convex,  $p_\ell^{\text{acc}}(x, y)$  also belongs to some bounded set for all  $x, y \in W$ , thus  $\|\hat{y} - p_\ell^{\text{acc}}(\hat{x}, \hat{y})\| < +\infty$  and  $\|\tilde{y} - p_\ell^{\text{acc}}(\tilde{x}, \tilde{y})\| < +\infty$ . Thus, the Lipschitz continuity of  $\nabla f_i$  shows such a boundedness of the second term. Moreover, the locally Lipschitz continuity of  $f_i$  and  $F_i$  derived by the continuous differentiability of  $f_i$  and convexity  $F_i$  lead to the similar property for the third and fourth terms. Hence,  $p_\ell^{\text{acc}}$  is Hölder continuous with exponent  $1/2$  on  $W$ .



On the other hand, the definition (18) of  $p_\ell^{\text{acc}}$  and  $\theta_\ell^{\text{acc}}$  gives

$$\begin{aligned}
& \theta_\ell^{\text{acc}}(\hat{x}, \hat{y}) - \theta_\ell^{\text{acc}}(\check{x}, \check{y}) \leq \varphi_\ell^{\text{acc}}(p_\ell^{\text{acc}}(\check{x}, \check{y}); \hat{x}, \hat{y}) - \varphi_\ell^{\text{acc}}(p_\ell^{\text{acc}}(\check{x}, \check{y}); \check{x}, \check{y}) \\
& = \max_{i=1, \dots, m} [\langle \nabla f_i(\hat{y}), p_\ell^{\text{acc}}(\check{x}, \check{y}) - \hat{y} \rangle + g_i(p_\ell^{\text{acc}}(\check{x}, \check{y})) + f_i(\hat{y}) - F_i(\hat{x})] \\
& \quad - \max_{i=1, \dots, m} [\langle \nabla f_i(\check{y}), p_\ell^{\text{acc}}(\check{x}, \check{y}) - \check{y} \rangle + g_i(p_\ell^{\text{acc}}(\check{x}, \check{y})) + f_i(\check{y}) - F_i(\check{x})] \\
& \quad + \frac{\ell}{2} \left[ \|p_\ell^{\text{acc}}(\check{x}, \check{y}) - \hat{y}\|^2 - \|p_\ell^{\text{acc}}(\check{x}, \check{y}) - \check{y}\|^2 \right] \\
& \leq \max_{i=1, \dots, m} [\langle \nabla f_i(\check{y}), \check{y} - \hat{y} \rangle + \langle \nabla f_i(\hat{y}) - \nabla f_i(\check{y}), p_\ell^{\text{acc}}(\check{x}, \check{y}) - \hat{y} \rangle \\
& \quad + f_i(\hat{y}) - f_i(\check{y}) - F_i(\hat{x}) + F_i(\check{x})] \\
& \quad + \frac{\ell}{2} \langle 2p_\ell^{\text{acc}}(\check{x}, \check{y}) - \hat{y} - \check{y}, \check{y} - \hat{y} \rangle \\
& \leq \max_{i=1, \dots, m} \|\nabla f_i(\hat{y})\| \|\hat{y} - \check{y}\| + \|\hat{y} - p_\ell^{\text{acc}}(\check{x}, \check{y})\| \max_{i=1, \dots, m} \|\nabla f_i(\hat{y}) - \nabla f_i(\check{y})\| \\
& \quad + \max_{i=1, \dots, m} |f_i(\hat{y}) - f_i(\check{y})| + \max_{i=1, \dots, m} |F_i(\hat{x}) - F_i(\check{x})| \\
& \quad + \frac{\ell}{2} \|2p_\ell^{\text{acc}}(\check{x}, \check{y}) - \hat{y} - \check{y}\| \|\hat{y} - \check{y}\|,
\end{aligned}$$

where the second inequality follows from the relation  $\max_{i=1, \dots, m} a_i - \max_{i=1, \dots, m} b_i \leq \max_{i=1, \dots, m} (a_i - b_i)$  for all  $a, b \in \mathbf{R}^m$ , and the third inequality comes from (19b) and Cauchy-Schwarz inequalities. Since the above inequality holds even if we interchange  $(\hat{x}, \hat{y})$  and  $(\check{x}, \check{y})$ , we can show the Lipschitz continuity of  $\theta_\ell^{\text{acc}}$  on  $W$  in the same way as in the previous paragraph.  $\square$

Note that the Hölder exponent  $1/2$  mentioned in Proposition 4.1 (ii) is optimal, i.e., for some  $F_i$ ,  $p_\ell^{\text{acc}}$  is not Hölder continuous with exponent  $\alpha > 1/2$ . In fact, this result was also proved for multiobjective steepest direction in [39].

Proposition 4.1 suggests that we can use  $\|p_\ell^{\text{acc}}(x, y) - y\|_\infty < \varepsilon$  for some  $\varepsilon > 0$  as a stopping criteria. Now, we state below the proposed algorithm.

---

**Algorithm 2** Accelerated proximal gradient method for multiobjective optimization

---

**Input:** Set  $x^0 = y^1 \in \text{dom } F$ ,  $\ell \geq L$ ,  $\varepsilon > 0$ .

**Output:**  $x^*$ : A weakly Pareto optimal point

```

1:  $k \leftarrow 1$ 
2:  $t_1 \leftarrow 1$ 
3: loop
4:    $x^k \leftarrow p_\ell^{\text{acc}}(x^{k-1}, y^k)$ , where  $p_\ell^{\text{acc}}$  is defined by (18)
5:   if  $\|x^k - y^k\|_\infty < \varepsilon$  then
6:     return  $x^k$ 
7:   end if
8:    $t_{k+1} \leftarrow \sqrt{t_k^2 + 1/4} + 1/2$ 
9:    $\gamma_k \leftarrow (t_k - 1)/t_{k+1}$ 
10:   $y^{k+1} \leftarrow x^k + \gamma_k(x^k - x^{k-1})$ 
11:   $k \leftarrow k + 1$ 
12: end loop

```

---

We show below some properties of  $\{t_k\}$  and  $\{\gamma_k\}$ , related to stepsizes.

**Lemma 4.2.** *Let  $\{t_k\}$  and  $\{\gamma_k\}$  be defined by lines 2, 8 and 9 in Algorithm 2. Then, the following inequalities hold for all  $k \geq 1$ :*

(i)  $t_{k+1} \geq t_k + 1/2$  and  $t_k \geq (k+1)/2$ ; (ii)  $t_k^2 - t_{k+1}^2 + t_{k+1} = 0$ ; (iii)  $1 - \gamma_k^2 \geq \frac{1}{t_k}$ .

*Proof.* Part (i) : From the definition of  $\{t_k\}$ , we have

$$t_{k+1} = \sqrt{t_k^2 + \frac{1}{4}} + \frac{1}{2} \geq t_k + \frac{1}{2}.$$

Applying the above inequality recursively, we obtain

$$t_k \geq t_1 + \frac{k-1}{2} = \frac{k+1}{2}.$$

Part (ii) : An easy computation shows that

$$t_k^2 - t_{k+1}^2 + t_{k+1} = t_k^2 - \left[ \sqrt{t_k^2 + \frac{1}{4}} + \frac{1}{2} \right]^2 + \sqrt{t_k^2 + \frac{1}{4}} + \frac{1}{2} = 0.$$

Part (iii) : Part (i) of this lemma implies that  $t_{k+1} > t_k \geq 1$ . Thus, the definition of  $\gamma_k$  leads to

$$1 - \gamma_k^2 = 1 - \left( \frac{t_k - 1}{t_{k+1}} \right)^2 \geq 1 - \left( \frac{t_k - 1}{t_k} \right)^2 = \frac{2t_k - 1}{t_k^2} \geq \frac{2t_k - t_k}{t_k^2} = \frac{1}{t_k}.$$

□

We end this section by noting some remarks about the proposed algorithm.

- Remark 4.1.** (i) When  $m = 1$ , we can remove the term  $f_i(y) - F_i(x)$  from the subproblem (15), so Algorithm 2 corresponds to the Fast Iterative Shrinkage-Thresholding Algorithm (FISTA) [3] for single-objective optimization.
- (ii) Algorithm 2 produces two sequences  $x^k$  and  $y^k$ , in a similar way to the single-objective FISTA. In particular, the stopping condition (Step 5), the momentum update (Steps 8 and 9), and the update of the iterate (Step 10) are actually equivalent to the single-objective case.
- (iii) Since  $x \in \text{dom } F$  implies  $p_\ell^{\text{acc}}(x, y) \in \text{dom } F$ , every  $x^k$  computed by the above algorithm is in  $\text{dom } F$ . However,  $y^k$  is not necessarily in  $\text{dom } F$ .
- (iv) Since  $y^1 = x^0$ , it follows from (21) that  $\theta_\ell^{\text{acc}}(x^0, y^1) \leq 0$ , but the inequality  $\theta_\ell^{\text{acc}}(x^{k-1}, y^k) \leq 0$  does not necessarily hold for  $k \geq 2$ .
- (v) Like Remark 3.1 (ii), Algorithm 2 induces the accelerated versions of first-order algorithms such as the steepest descent [18], proximal point [7], and projected gradient methods [26].
- (vi) Like Remark 3.1 (iii), even if it is difficult to estimate  $L$ , we can update the constant  $\ell$  to satisfy  $F_i(p_\ell^{\text{acc}}(x^{k-1}, y^k)) - F_i(x^{k-1}) \leq \theta_\ell^{\text{acc}}(x^{k-1}, y^k)$  for all  $i = 1, \dots, m$  in each iteration by a finite number of backtracking steps. Moreover, we can restrict the assumption of  $\nabla f_i$ 's Lipschitz continuity on the level set  $\mathcal{L}_F(F(x^0))$  without affecting the analysis in the subsequent sections.

## 5 Convergence rate

This section shows that Algorithm 2 has a convergence rate of  $O(1/k^2)$  under the same assumptions used in the complexity analysis of Algorithm 1. As it is expected, this rate is better than the one obtained for Algorithm 1.

Let us first define some functions below, that will be useful for our complexity analysis. For  $k \geq 0$ , let  $\sigma_k: \mathbf{R}^n \rightarrow \mathbf{R} \cup \{-\infty\}$  and  $\rho_k: \mathbf{R}^n \rightarrow \mathbf{R}$  be defined by

$$\begin{aligned} \sigma_k(z) &:= \min_{i=1, \dots, m} [F_i(x^k) - F_i(z)], \\ \rho_k(z) &:= \|t_{k+1}x^{k+1} - (t_{k+1} - 1)x^k - z\|^2, \end{aligned} \tag{22}$$

respectively. We present a lemma on  $\sigma_k$  that will be helpful in the subsequent discussions.

**Lemma 5.1.** Let  $\{x^k\}$  and  $\{y^k\}$  be sequences generated by Algorithm 2. Then, the following inequalities hold for all  $z \in \mathbf{R}^n$  and  $k \geq 0$ :

$$\begin{aligned} \sigma_{k+1}(z) &\leq -\frac{\ell}{2} \left[ 2\langle x^{k+1} - y^{k+1}, y^{k+1} - z \rangle + \|x^{k+1} - y^{k+1}\|^2 \right] \\ &\quad - \frac{\ell - L}{2} \|x^{k+1} - y^{k+1}\|^2, \end{aligned} \tag{23}$$

$$\begin{aligned} \sigma_k(z) - \sigma_{k+1}(z) &\geq \frac{\ell}{2} \left[ 2\langle x^{k+1} - y^{k+1}, y^{k+1} - x^k \rangle + \|x^{k+1} - y^{k+1}\|^2 \right] \\ &\quad + \frac{\ell - L}{2} \|x^{k+1} - y^{k+1}\|^2. \end{aligned} \tag{24}$$

*Proof.* Suppose that  $z \in \mathbf{R}^n$  and  $k \geq 0$ . Recall that there exist  $\eta(x^k, y^{k+1}) \in \partial g(x^{k+1})$  and a Lagrange multiplier  $\lambda(x^k, y^{k+1}) \in \mathbf{R}^m$  that satisfy the KKT condition (19) for the subproblem (15). From the definition (22) of  $\sigma_{k+1}$ , we get

$$\sigma_{k+1}(z) = \min_{i=1, \dots, m} [F_i(x^{k+1}) - F_i(z)] \leq \sum_{i=1}^m \lambda_i(x^k, y^{k+1}) [F_i(x^{k+1}) - F_i(z)].$$

where the inequality follows from (19b). Taking  $p = x^{k+1}$ ,  $q = y^{k+1}$ , and  $r = z$  in (8), we have

$$\begin{aligned} \sigma_{k+1}(z) &\leq \sum_{i=1}^m \lambda_i(x^k, y^{k+1}) [\langle \nabla f_i(y^{k+1}), x^{k+1} - y^{k+1} \rangle + g_i(x^{k+1}) \\ &\quad + f_i(y^{k+1}) - F_i(z)] + \frac{L}{2} \|x^{k+1} - y^{k+1}\|^2. \end{aligned}$$

Hence, the convexity of  $f_i$  and  $g_i$  yields

$$\begin{aligned} \sigma_{k+1}(z) &\leq \sum_{i=1}^m \lambda_i(x^k, y^{k+1}) [\langle \nabla f_i(y^{k+1}), x^{k+1} - y^{k+1} \rangle + \langle \nabla f_i(y^{k+1}), y^{k+1} - z \rangle \\ &\quad + \langle \eta_i(x^k, y^{k+1}), x^{k+1} - z \rangle] + \frac{L}{2} \|x^{k+1} - y^{k+1}\|^2 \\ &= \sum_{i=1}^m \lambda_i(x^k, y^{k+1}) \langle \nabla f_i(y^{k+1}) + \eta_i(x^k, y^{k+1}), x^{k+1} - z \rangle + \frac{L}{2} \|x^{k+1} - y^{k+1}\|^2. \end{aligned}$$

Using (19a) with  $x = x^k$  and  $y = y^{k+1}$  and from the fact that  $x^{k+1} = p_\ell^{\text{acc}}(x^k, y^{k+1})$  (see line 4 of Algorithm 2), we obtain

$$\sigma_{k+1}(z) \leq -\ell \langle x^{k+1} - y^{k+1}, x^{k+1} - z \rangle + \frac{L}{2} \|x^{k+1} - y^{k+1}\|^2.$$

Moreover, simple calculations show that

$$\begin{aligned} \sigma_{k+1}(z) &\leq -\frac{\ell}{2} \left[ 2 \langle x^{k+1} - y^{k+1}, x^{k+1} - z \rangle - \|x^{k+1} - y^{k+1}\|^2 \right] - \frac{\ell - L}{2} \|x^{k+1} - y^{k+1}\|^2 \\ &= -\frac{\ell}{2} \left[ 2 \langle x^{k+1} - y^{k+1}, y^{k+1} - z \rangle + \|x^{k+1} - y^{k+1}\|^2 \right] - \frac{\ell - L}{2} \|x^{k+1} - y^{k+1}\|^2, \end{aligned}$$

which completes the proof of (23).

Now, let us show inequality (24). Again from the definition (22) of  $\sigma_k$ , we obtain

$$\begin{aligned} \sigma_k(z) - \sigma_{k+1}(z) &= \min_{i=1, \dots, m} [F_i(x^k) - F_i(z)] - \min_{i=1, \dots, m} [F_i(x^{k+1}) - F_i(z)] \\ &\geq -\max_{i=1, \dots, m} [F_i(x^{k+1}) - F_i(x^k)], \end{aligned} \tag{25}$$

where the inequality holds because

$$\min_{i=1, \dots, m} (u_i + v_i) - \min_{i=1, \dots, m} u_i \geq \min_{i=1, \dots, m} v_i \quad \text{for all } u, v \in \mathbf{R}^m.$$

Letting  $p = x^{k+1}$ ,  $q = y^{k+1}$ , and  $r = x^k$  in (8), we have

$$\begin{aligned}
\sigma_k(z) - \sigma_{k+1}(z) &\geq - \max_{i=1, \dots, m} [\langle \nabla f_i(y^{k+1}), x^{k+1} - y^{k+1} \rangle + g_i(x^{k+1}) + f_i(y^{k+1}) \\
&\quad - F_i(x^k)] - \frac{L}{2} \|x^{k+1} - y^{k+1}\|^2 \\
&= - \sum_{i=1}^m \lambda_i(x^k, y^{k+1}) [\langle \nabla f_i(y^{k+1}), x^{k+1} - y^{k+1} \rangle + g_i(x^{k+1}) \\
&\quad + f_i(y^{k+1}) - F_i(x^k)] - \frac{L}{2} \|x^{k+1} - y^{k+1}\|^2 \\
&= - \sum_{i=1}^m \lambda_i(x^k, y^{k+1}) [\langle \nabla f_i(y^{k+1}), x^k - y^{k+1} \rangle + f_i(y^{k+1}) - f_i(x^k)] \\
&\quad - \sum_{i=1}^m \lambda_i(x^k, y^{k+1}) [\langle \nabla f_i(y^{k+1}), x^{k+1} - x^k \rangle + g_i(x^{k+1}) - g_i(x^k)] \\
&\quad - \frac{L}{2} \|x^{k+1} - y^{k+1}\|^2,
\end{aligned}$$

where the first equality comes from (19b), and the second one follows by taking  $x^{k+1} - y^{k+1} = (x^k - y^{k+1}) + (x^{k+1} - x^k)$ . From the convexity of  $f_i$ , the first term of the above expression is nonnegative. Moreover, the convexity of  $g_i$  shows that

$$\begin{aligned}
\sigma_k(z) - \sigma_{k+1}(z) &\geq - \sum_{i=1}^m \lambda_i(x^k, y^{k+1}) \langle \nabla f_i(y^{k+1}) + \eta_i(x^k, y^{k+1}), x^{k+1} - x^k \rangle - \frac{L}{2} \|x^{k+1} - y^{k+1}\|^2.
\end{aligned}$$

Thus, (19a) with  $(x, y) = (x^k, y^{k+1})$  and direct calculations prove that

$$\begin{aligned}
\sigma_k(z) - \sigma_{k+1}(z) &\geq \ell \langle x^{k+1} - y^{k+1}, x^{k+1} - x^k \rangle - \frac{L}{2} \|x^{k+1} - y^{k+1}\|^2 \\
&= \frac{\ell}{2} \left[ 2 \langle x^{k+1} - y^{k+1}, x^{k+1} - x^k \rangle - \|x^{k+1} - y^{k+1}\|^2 \right] + \frac{\ell - L}{2} \|x^{k+1} - y^{k+1}\|^2 \\
&= \frac{\ell}{2} \left[ 2 \langle x^{k+1} - y^{k+1}, y^{k+1} - x^k \rangle + \|x^{k+1} - y^{k+1}\|^2 \right] + \frac{\ell - L}{2} \|x^{k+1} - y^{k+1}\|^2.
\end{aligned}$$

□

We can also show the following corollary of Lemma 5.1 (24).

**Corollary 5.1.** *Let  $\{x^k\}$  and  $\{y^k\}$  be sequences generated by Algorithm 2. Then, we have*

$$\begin{aligned}
\sigma_{k_1}(z) - \sigma_{k_2}(z) &\geq \frac{\ell}{2} \left[ \|x^{k_2} - x^{k_2-1}\|^2 - \|x^{k_1} - x^{k_1-1}\|^2 + \sum_{k=k_1}^{k_2-1} \frac{1}{t_k} \|x^k - x^{k-1}\|^2 \right]
\end{aligned}$$

for any  $k_2 \geq k_1 \geq 1$ .

*Proof.* Let  $k \geq 1$ . Since  $\ell \geq L$ , Lemma 5.1 (24) implies

$$\begin{aligned}
\sigma_k(z) - \sigma_{k+1}(z) &\geq \frac{\ell}{2} \left[ 2 \langle x^{k+1} - y^{k+1}, y^{k+1} - x^k \rangle + \|x^{k+1} - y^{k+1}\|^2 \right] \\
&= \frac{\ell}{2} \left[ \|x^{k+1} - x^k\|^2 - \|y^{k+1} - x^k\|^2 \right],
\end{aligned}$$

where the equality holds from (2) with  $(a, b, c) = (y^{k+1}, x^{k+1}, x^k)$ . Hence, the definition of  $y^{k+1}$  given in line 10 of Algorithm 2 yields

$$\sigma_k(z) - \sigma_{k+1}(z) \geq \frac{\ell}{2} \left[ \|x^{k+1} - x^k\|^2 - \gamma_k^2 \|x^k - x^{k-1}\|^2 \right].$$

Applying this inequality repeatedly, we have

$$\begin{aligned} \sigma_{k_1}(z) - \sigma_{k_2}(z) & \\ & \geq \frac{\ell}{2} \left[ \|x^{k_2} - x^{k_2-1}\|^2 - \|x^{k_1} - x^{k_1-1}\|^2 + \sum_{k=k_1}^{k_2-1} (1 - \gamma_k^2) \|x^k - x^{k-1}\|^2 \right]. \end{aligned}$$

Using Lemma 4.2 (iii), we get the desired inequality.  $\square$

Before analyzing the convergence rate of Algorithm 2, we show that the objective function values at  $x^k$  for any  $k \geq 0$  never exceed the ones at the initial point, that is,  $\{x^k\}$  belongs to the level set  $\mathcal{L}_F(F(x^0))$  (see (10) for the definition of  $\mathcal{L}_F$ ). However, note that Algorithm 2 does not guarantee the monotonically decreasing property  $F(x^{k+1}) \leq F(x^k)$ .

**Theorem 5.1.** *Algorithm 2 generates a sequence  $\{x^k\}$  such that*

$$F_i(x^k) \leq F_i(x^0) \quad \text{for all } i = 1, \dots, m, k \geq 0.$$

*Proof.* Let  $i = 1, \dots, m$  and  $p \geq 1$ . Then, we have

$$F_i(x^p) - F_i(x^{p+1}) \geq - \max_{i=1, \dots, m} [F_i(x^{p+1}) - F_i(x^p)].$$

Noting that the right-hand side of the above inequality is the same as (25), with similar arguments used in the proof of (24) in Lemma 5.1, we obtain

$$\begin{aligned} F_i(x^p) - F_i(x^{p+1}) & \geq \frac{\ell}{2} \left[ 2 \langle x^{p+1} - y^{p+1}, y^{p+1} - x^p \rangle + \|x^{p+1} - y^{p+1}\|^2 \right] \\ & \quad + \frac{\ell - L}{2} \|x^{p+1} - y^{p+1}\|^2. \end{aligned} \quad (26)$$

Note that this inequality also holds for  $p = 0$ . Again, in the same way as in the proof of Corollary 5.1, we get

$$F_i(x^1) - F_i(x^k) \geq \frac{\ell}{2} \left[ \|x^k - x^{k-1}\|^2 - \|x^1 - x^0\|^2 + \sum_{p=1}^{k-1} \frac{1}{t_p} \|x^p - x^{p-1}\|^2 \right].$$

Since  $t_1 = 1$ , the above inequality reduces to

$$F_i(x^1) - F_i(x^k) \geq \frac{\ell}{2} \left[ \|x^k - x^{k-1}\|^2 + \sum_{p=2}^{k-1} \frac{1}{t_p} \|x^p - x^{p-1}\|^2 \right] \geq 0.$$

Moreover, (26) with  $p = 0$  and the fact that  $y^1 = x^0$  imply  $F_i(x^1) \leq F_i(x^0)$ , so we can conclude that  $F_i(x^k) \leq F_i(x^0)$ .  $\square$

The following result provides the fundamental relation for our convergence rate analysis.

**Lemma 5.2.** *Let  $\{x^k\}$  and  $\{y^k\}$  be sequences generated by Algorithm 2. Also, let  $\sigma_k$  and  $\rho_k$  be defined by (22). Then, we have*

$$t_{k+1}^2 \sigma_{k+1}(z) + \frac{\ell}{2} \rho_k(z) + \frac{\ell - L}{2} \sum_{p=1}^k t_{p+1}^2 \|x^{p+1} - y^{p+1}\|^2 \leq \frac{\ell}{2} \|x^0 - z\|^2$$

for all  $k \geq 0$  and  $z \in \mathbf{R}^n$ .

*Proof.* Let  $p \geq 1$  and  $z \in \mathbf{R}^n$ . Recall from Lemma 5.1 that

$$\begin{aligned} -\sigma_{p+1}(z) &\geq \frac{\ell}{2} \left[ 2\langle x^{p+1} - y^{p+1}, y^{p+1} - z \rangle + \|x^{p+1} - y^{p+1}\|^2 \right] \\ &\quad + \frac{\ell - L}{2} \|x^{p+1} - y^{p+1}\|^2, \\ \sigma_p(z) - \sigma_{p+1}(z) &\geq \frac{\ell}{2} \left[ 2\langle x^{p+1} - y^{p+1}, y^{p+1} - x^p \rangle + \|x^{p+1} - y^{p+1}\|^2 \right] \\ &\quad + \frac{\ell - L}{2} \|x^{p+1} - y^{p+1}\|^2. \end{aligned}$$

To get a relation between  $\sigma_p(z)$  and  $\sigma_{p+1}(z)$ , we multiply the second inequality above by  $(t_{p+1} - 1)$  and add it to the first one:

$$\begin{aligned} (t_{p+1} - 1)\sigma_p(z) - t_{p+1}\sigma_{p+1}(z) &\geq \frac{\ell}{2} \left[ t_{p+1} \|x^{p+1} - y^{p+1}\|^2 + 2\langle x^{p+1} - y^{p+1}, t_{p+1}y^{p+1} - (t_{p+1} - 1)x^p - z \rangle \right] \\ &\quad + \frac{\ell - L}{2} t_{p+1} \|x^{p+1} - y^{p+1}\|^2. \end{aligned}$$

Multiplying this inequality by  $t_{p+1}$  and using the relation  $t_p^2 = t_{p+1}^2 - t_{p+1}$  (cf. Lemma 4.2 (ii)), we get

$$\begin{aligned} t_p^2\sigma_p(z) - t_{p+1}^2\sigma_{p+1}(z) &\geq \frac{\ell}{2} \left[ \|t_{p+1}(x^{p+1} - y^{p+1})\|^2 + 2t_{p+1}\langle x^{p+1} - y^{p+1}, t_{p+1}y^{p+1} - (t_{p+1} - 1)x^p - z \rangle \right] \\ &\quad + \frac{\ell - L}{2} t_{p+1}^2 \|x^{p+1} - y^{p+1}\|^2. \end{aligned}$$

Applying (2) with  $(a, b, c) = (t_{p+1}y^{p+1}, t_{p+1}x^{p+1}, (t_{p+1} - 1)x^p + z)$  to the right-hand side of the last inequality, we get

$$\begin{aligned} t_p^2\sigma_{p+1}(z) - t_{p+1}^2\sigma_p(z) &\geq \frac{\ell}{2} \left[ \|t_{p+1}x^{p+1} - (t_{p+1} - 1)x^p - z\|^2 - \|t_{p+1}y^{p+1} - (t_{p+1} - 1)x^p - z\|^2 \right] \\ &\quad + \frac{\ell - L}{2} t_{p+1}^2 \|x^{p+1} - y^{p+1}\|^2. \end{aligned}$$

Recall that  $\rho_p(z) := \|t_{p+1}x^{p+1} - (t_{p+1} - 1)x^p - z\|^2$ . Then, from the definition of  $y^p$  defined in line 10 of Algorithm 2, we get

$$t_p^2\sigma_p(z) - t_{p+1}^2\sigma_{p+1}(z) \geq \frac{\ell}{2} [\rho_p(z) - \rho_{p-1}(z)] + \frac{\ell - L}{2} t_{p+1}^2 \|x^{p+1} - y^{p+1}\|^2.$$

Now, let  $k \geq 0$ . Adding the above inequality from  $p = 0$  to  $p = k$  and using  $t_1 = 1$  and  $\rho_0(z) = \|x^1 - z\|^2$ , we have

$$\sigma_1(z) - t_{k+1}^2\sigma_{k+1}(z) \geq \frac{\ell}{2} \left[ \rho_k(z) - \|x^1 - z\|^2 \right] + \frac{\ell - L}{2} \sum_{p=1}^k t_{k+1}^2 \|x^{k+1} - y^{k+1}\|^2.$$

Lemma 4.2 (23) with  $k = 0$  and  $y^1 = x^0$  lead to

$$\begin{aligned} \sigma_1(z) &\leq -\frac{\ell}{2} \left[ \|x^1 - z\|^2 - \|x^0 - z\|^2 \right] - \frac{\ell - L}{2} \|x^1 - y^1\|^2 \\ &\leq -\frac{\ell}{2} \left[ \|x^1 - z\|^2 - \|x^0 - z\|^2 \right], \end{aligned}$$

where the second inequality follows since  $\ell \geq L$ . From the above two inequalities, we can derive the desired inequality.  $\square$

Finally, using Lemma 5.2, we can evaluate the convergence rate of Algorithm 2 with the following theorem.

**Theorem 5.2.** *Under Assumption 3.1, Algorithm 2 generates a sequence  $\{x^k\}$  such that*

$$u_0(x^k) \leq \frac{2\ell R}{(k+1)^2},$$

where  $R \geq 0$  is given in (14), and  $u_0$  is a merit function defined by (11).

*Proof.* Let  $k \geq 0$ . Since  $\rho_k(z) \geq 0$ , Lemma 5.2 gives

$$t_{k+1}^2 \sigma_{k+1}(z) \leq \frac{\ell}{2} \|x^0 - z\|^2.$$

It follows from Lemma 4.2 (i) that

$$\sigma_{k+1}(z) \leq \frac{2\ell}{(k+2)^2} \|x^0 - z\|^2.$$

With similar arguments used in the proof of Theorem 3.1 (see [41, Theorem 5.2]), we get the desired inequality.  $\square$

We end this section by showing that the global convergence of Algorithm 2, in terms of weak Pareto optimality, is also guaranteed by using the above complexity result.

**Corollary 5.2.** *Suppose that Assumption 3.1 holds. Then, every accumulation point of the sequence  $\{x^k\}$  generated by Algorithm 2 is weakly Pareto optimal for (7). In particular, if the level set  $\mathcal{L}_F(F(x^0))$  is bounded, then  $\{x^k\}$  has accumulation points, and they are all weakly Pareto optimal.*

*Moreover, if each  $F_i$  is strictly convex, then the accumulation points are Pareto optimum, i.e., there does not exist any points with the same or smaller objective function values and with at least one objective function value being strictly smaller.*

*Proof.* The first claim is clear from the lower-semicontinuity of  $F_i$  for all  $i = 1, \dots, m$  as well as Theorems 2.1 and 5.2, and the second one is easy since Theorem 5.1 holds. The third is also obvious from the relationship between weak Pareto and Pareto optimalities [40, Lemma 2.2].  $\square$

## 6 Efficient computation of the subproblem via its dual

In the previous section, we proved global convergence and complexity results of Algorithm 2. Now, we want to show how practical is the proposed method. In particular, we now discuss a way of computing the subproblem (15). First, define

$$\psi_i(z; x, y) := \langle \nabla f_i(y), z - y \rangle + g_i(z) + f_i(y) - F_i(x) + \frac{\ell}{2} \|z - y\|^2 \quad (27)$$

for all  $i = 1, \dots, m$ . Then, fixing some  $\ell \geq L$ , we can rewrite the objective function  $\varphi_\ell^{\text{acc}}(z; x, y)$  of (15) as

$$\varphi_\ell^{\text{acc}}(z; x, y) = \max_{i=1, \dots, m} \psi_i(z; x, y).$$

Recall that  $\Delta^m \subseteq \mathbf{R}^m$  represents the standard simplex (1). Since  $\max_{i=1, \dots, m} q_i = \max_{\lambda \in \Delta^m} \sum_{i=1}^m \lambda_i q_i$  for any  $q \in \mathbf{R}^m$ , we get

$$\varphi_\ell^{\text{acc}}(z; x, y) = \max_{\lambda \in \Delta^m} \sum_{i=1}^m \lambda_i \psi_i(z; x, y).$$

Then, the subproblem (15) reduces to the following minimax problem:

$$\min_{z \in \mathbf{R}^n} \max_{\lambda \in \Delta^m} \sum_{i=1}^m \lambda_i \psi_i(z; x, y). \quad (28)$$

We can see that  $\mathbf{R}^n$  is convex,  $\Delta^m$  is compact and convex, and  $\sum_{i=1}^m \lambda_i \psi_i(z; x, y)$  is convex for  $z$  and concave for  $\lambda$ . Therefore, Sion's minimax theorem [37] shows that the above problem is equivalent to

$$\max_{\lambda \in \Delta^m} \min_{z \in \mathbf{R}^n} \sum_{i=1}^m \lambda_i \psi_i(z; x, y). \quad (29)$$

The definition (27) of  $\psi_i$  yields

$$\begin{aligned} \min_{z \in \mathbf{R}^n} \sum_{i=1}^m \lambda_i \psi_i(z; x, y) &= \min_{z \in \mathbf{R}^n} \left[ \sum_{i=1}^m \lambda_i g_i(z) + \frac{\ell}{2} \left\| z - y + \frac{1}{\ell} \sum_{i=1}^m \lambda_i \nabla f_i(y) \right\|^2 \right] \\ &\quad - \frac{1}{2\ell} \left\| \sum_{i=1}^m \lambda_i \nabla f_i(y) \right\|^2 + \sum_{i=1}^m \lambda_i \{f_i(y) - F_i(x)\} \\ &= \ell \mathcal{M}_{\frac{1}{\ell} \sum_{i=1}^m \lambda_i g_i} \left( y - \frac{1}{\ell} \sum_{i=1}^m \lambda_i \nabla f_i(y) \right) \\ &\quad - \frac{1}{2\ell} \left\| \sum_{i=1}^m \lambda_i \nabla f_i(y) \right\|^2 + \sum_{i=1}^m \lambda_i \{f_i(y) - F_i(x)\}, \end{aligned}$$

where  $\mathcal{M}$  is the Moreau envelope (4). Based on the discussion above, we obtain the dual problem of (15) as follows:

$$\begin{aligned} \max_{\lambda \in \mathbf{R}^m} \quad & \omega(\lambda) \\ \text{s.t.} \quad & \lambda \geq 0 \quad \text{and} \quad \sum_{i=1}^m \lambda_i = 1, \end{aligned} \tag{30}$$

where

$$\begin{aligned} \omega(\lambda) &:= \ell \mathcal{M}_{\frac{1}{\ell} \sum_{i=1}^m \lambda_i g_i} \left( y - \frac{1}{\ell} \sum_{i=1}^m \lambda_i \nabla f_i(y) \right) \\ &\quad - \frac{1}{2\ell} \left\| \sum_{i=1}^m \lambda_i \nabla f_i(y) \right\|^2 + \sum_{i=1}^m \lambda_i [f_i(y) - F_i(x)]. \end{aligned} \tag{31}$$

If we can find the global optimal solution  $\lambda^*$  of this dual problem (30), we can construct the optimal solution  $z^*$  of the original subproblem (15) as

$$z^* = \mathbf{prox}_{\frac{1}{\ell} \sum_{i=1}^m \lambda_i^* g_i} \left( y - \frac{1}{\ell} \sum_{i=1}^m \lambda_i^* \nabla f_i(y) \right),$$

where  $\mathbf{prox}$  denotes the proximal operator (5). This is because the equivalence between (28) and (29) induces

$$\sum_{i=1}^m \lambda_i^* \psi_i(z^*; x, y) = \max_{\lambda \in \Delta^n} \min_{z \in \mathbf{R}^n} \sum_{i=1}^m \lambda_i \psi_i(z; x, y) = \min_{z \in \mathbf{R}^n} \max_{\lambda \in \Delta^n} \sum_{i=1}^m \lambda_i \psi_i(z; x, y),$$

which means that  $z^*$  attains the minimum in (28). Since  $\sum_{i=1}^m \lambda_i \psi_i(z; x, y)$  is concave for  $\lambda$ , it is clear that  $\omega(\lambda) = \min_{z \in \mathbf{R}^n} \sum_{i=1}^m \lambda_i \psi_i(z; x, y)$  is concave. Furthermore,  $\omega$  is differentiable, as the following theorem shows.

**Theorem 6.1.** *The function  $\omega: \mathbf{R}^m \rightarrow \mathbf{R}$  defined by (31) is continuously differentiable at every  $\lambda \in \mathbf{R}^m$  and*

$$\begin{aligned} \nabla \omega(\lambda) &= g \left( \mathbf{prox}_{\frac{1}{\ell} \sum_{i=1}^m \lambda_i g_i} \left( y - \frac{1}{\ell} \sum_{i=1}^m \lambda_i \nabla f_i(y) \right) \right) \\ &\quad + J_f(y) \left( \mathbf{prox}_{\frac{1}{\ell} \sum_{i=1}^m \lambda_i g_i} \left( y - \frac{1}{\ell} \sum_{i=1}^m \lambda_i \nabla f_i(y) \right) - y \right) + f(y) - F(x), \end{aligned}$$

where  $\mathbf{prox}$  is the proximal operator (5), and  $J_f(y)$  is the Jacobian matrix at  $y$  given by

$$J_f(y) := (\nabla f_1(y), \dots, \nabla f_m(y))^\top.$$

*Proof.* Define

$$h(z, \lambda) := \sum_{i=1}^m \lambda_i g_i(z) + \frac{\ell}{2} \left\| z - y + \frac{1}{\ell} \sum_{i=1}^m \lambda_i \nabla f_i(y) \right\|^2.$$



Clearly,  $h$  is continuous on  $\mathbf{R}^n \times \mathbf{R}^m$ . Moreover,  $h_z(\cdot) := h(z, \cdot)$  is continuously differentiable and

$$\nabla_{\lambda} h_z(\lambda) = g(z) + J_f(y) \left( z - y + \frac{1}{\ell} \sum_{i=1}^m \lambda_i \nabla f_i(y) \right).$$

Furthermore,

$$\mathbf{prox}_{\frac{1}{\ell} \sum_{i=1}^m \lambda_i g_i} \left( y - \frac{1}{\ell} \sum_{i=1}^m \lambda_i \nabla f_i(y) \right) = \underset{z \in \mathbf{R}^n}{\operatorname{argmin}} h(z, \lambda)$$

is also continuous at every  $\lambda \in \mathbf{R}^m$  (cf. [36, Theorem 2.26 and Exercise 7.38]). Therefore, the well-known result in first order differentiability analysis of the optimal value function [6, Theorem 4.13] gives

$$\begin{aligned} & \nabla_{\lambda} \left[ \ell \mathcal{M}_{\frac{1}{\ell} \sum_{i=1}^m \lambda_i g_i} \left( y - \frac{1}{\ell} \sum_{i=1}^m \lambda_i \nabla f_i(y) \right) \right] \\ &= g \left( \mathbf{prox}_{\frac{1}{\ell} \sum_{i=1}^m \lambda_i g_i} \left( y - \frac{1}{\ell} \sum_{i=1}^m \lambda_i \nabla f_i(y) \right) \right) \\ & \quad + J_f(y) \left( \mathbf{prox}_{\frac{1}{\ell} \sum_{i=1}^m \lambda_i g_i} \left( y - \frac{1}{\ell} \sum_{i=1}^m \lambda_i \nabla f_i(y) \right) - y + \frac{1}{\ell} \sum_{i=1}^m \lambda_i \nabla f_i(y) \right). \end{aligned}$$

On the other hand, we have

$$\begin{aligned} & \nabla_{\lambda} \left[ -\frac{1}{2\ell} \left\| \sum_{i=1}^m \lambda_i \nabla f_i(y) \right\|^2 + \sum_{i=1}^m \lambda_i \{f_i(y) - F_i(x)\} \right] \\ &= -\frac{1}{\ell} J_f(y) \sum_{i=1}^m \lambda_i \nabla f_i(y) + f(y) - F(x). \end{aligned}$$

Adding the above two equalities, we get the desired result.  $\square$

This theorem shows that the dual problem (30) is an  $m$ -dimensional differentiable convex optimization problem. Hence, if we can compute the proximal operator of  $\sum_{i=1}^m \lambda_i g_i$  quickly, then we can solve (30) using convex optimization techniques such as the interior point method [9]. In addition, for cases where  $n \gg m$ , the computational cost is much lower than solving the subproblem (15) directly. In particular, when  $m = 2$ , eliminating a variable with  $\lambda_2 = 1 - \lambda_1$  reduces (30) to a one-dimensional optimization that can be solved quickly using, for example, Brent's method [10]. Note, for example, that if  $g_i(x) = g_1(x)$  for all  $i = 1, \dots, m$ , or if  $g_i(x) = g_i(x_{I_i})$  and the index sets  $I_i$  do not overlap each other, then we can evaluate the proximal operator of  $\sum_{i=1}^m \lambda_i g_i$  from the proximal operator of each  $g_i$ . Furthermore, even if there is an overlap, we can compute such a proximal operator immediately for special functions, for example,  $m = 2, g_1(x) = \|x\|_1, g_2(x) = \|x\|_2^2$  ( $\lambda_1 g_1(x) + \lambda_2 g_2(x)$  is the elastic net [50] when  $\lambda_1 > 0$  and  $\lambda_2 > 0$ . The elastic net has a proximal operator in closed-form [35, Section 6.5.3]).

## 7 Numerical experiments

This section illustrates the proposed method's performance compared to the proximal gradient method without acceleration (Algorithm 1), and the algorithm below. Unlike the proposed Algorithm 2, Algorithm 3 does not include the term  $f_i(y) - F_i(x)$ , which was the key to the proof of Theorem 5.2, in the subproblem solved in Step 4. Therefore, the convergence rate of Algorithm 3 is still theoretically unknown. However, since it is the easiest algorithm to conceive from the scalar optimization FISTA, and Algorithm 3 is consistent with [16] when  $g = 0$ , we use it as a comparison in the numerical experiments.

### 7.1 Test problems

We generate a new list of convex multiobjective optimization test problems by processing the problem list of [31] based on the following three criteria:

---

**Algorithm 3** Accelerated proximal gradient method for multiobjective optimization (without  $f_i(y) - F_i(x)$ )

---

**Input:** Set  $x^0 = y^1 \in \text{dom } F, \ell \geq L, \varepsilon > 0$ .

**Output:**  $x^*$ : A weakly Pareto optimal point

```

1:  $k \leftarrow 1$ 
2:  $t_1 \leftarrow 1$ 
3: loop
4:    $x^k \leftarrow$  the optimal solution of (17)
5:   if  $\|x^k - y^k\|_\infty < \varepsilon$  then
6:     return  $x^k$ 
7:   end if
8:    $t_{k+1} \leftarrow \sqrt{t_k^2 + 1/4} + 1/2$ 
9:    $\gamma_k \leftarrow (t_k - 1)/t_{k+1}$ 
10:   $y^{k+1} \leftarrow x^k + \gamma_k(x^k - x^{k-1})$ 
11:   $k \leftarrow k + 1$ 
12: end loop

```

---

- Extracting convex problems: Since our proposed method is designed for convex problems, we selected only the convex problems from the original problem list.
- Dealing with various dimensions of  $n$ : For some problems, we enhance the variety by using different values of  $n$ .
- Including  $g$ : The original test problems include both constrained and unconstrained problems. For constrained problems, we set  $g_i$  as the indicator function (6) corresponding to the constraint for every  $i = 1, \dots, m$ . For unconstrained problems, we consider two types:  $g_i = 0$  and  $g_i = \|x - i + 1\|_1 / [(i - 1)n]$  (i.e. using  $\ell_1$ -norm) for each  $i = 1, \dots, m$ .

Based on this, the new list is given in Table 1.

Table 1: List of test problems

Problem name	$m$	$n$	$g_i$
JOS1 [27]	2	5, 10, 20, 50, 100, 200, 500, 1000	0 or $\ell_1$
ZDT1 [48]	2	5, 10, 20, 50, 100, 200, 500, 1000	indicator function
SD [38]	2	4	indicator function
TOI4 [43]	2	4	0 or $\ell_1$
TRIDIA [43]	3	3	0 or $\ell_1$
FDS [17]	3	5, 10, 20, 50, 100	0 or $\ell_1$
LF1 [32]	4	30, 100, 1000	0 or $\ell_1$

## 7.2 Experimental settings

The experiments are carried out on a machine with 2.4 GHz Intel Xeon Silver 4210R CPU and 64 GB memory, implementing all codes in Python 3.9.5. In all algorithms, we convert the subproblem into its dual as discussed in Section 6 and solve it using the trust-region interior point method [11] with the scientific library SciPy. The stopping tolerance for solving the subproblem is  $10^{-11}$ , except for the difficult problem TRIDIA where we use  $10^{-6}$ . Also, we use backtracking procedure to determine a parameter  $\ell$ , where the initial value of  $\ell$  is 1 and the constant multiplied to  $\ell$  is 2. We set the general stopping criteria as  $\varepsilon = 10^{-5}$  for each experiment. Moreover, we choose 100 initial points, commonly for both algorithms, uniformly, and randomly between the bounds given in [31]. The source code used here is available at <https://github.com/zalgo3/zfista>.

## 7.3 Evaluation metrics

We use the following metrics to assess the algorithms' performance:

- **The number of iterations:** The number of iterations required to satisfy the stopping criteria.

- **Time:** The time needed to meet the stopping criteria.
- **Purity [1]:** The ratio of the solutions obtained by a given solver within the approximated Pareto frontier. Let  $PF_{p,s}$  be the set of function values of the solutions obtained by solver  $s \in \mathcal{S}$  for problem  $p \in \mathcal{P}$  that are not dominated by other solutions, and let  $PF_p$  be the set of  $\bigcup_{s \in \mathcal{S}} PF_{p,s}$  that are not dominated by other solutions. The purity is defined by  $|PF_{p,s} \cap PF_p|/|PF_p|$ .
- **Hypervolume [49]:** The sum of the volumes of the hyperrectangles where the line segment connecting the reference point and each point of  $PF_{p,s}$  forms a diagonal. We set as the reference point the maximum value of each objective function in  $PF_p$ .
- **Spread metrics ( $\Gamma$  and  $\Delta$ ) [14]:** The metric representing how well-distributed the obtained Pareto frontier is. Let  $PF_{p,s} \cap PF_p$  be formed by  $F^1, \dots, F^N$ . Assume that  $F_j^{a_j^i} \leq F_j^{a_j^{i+1}}$  for some  $\{a_j^i\} \subseteq \{1, \dots, m\}$  and for each  $i = 1, \dots, N, j = 1, \dots, m$ . Moreover, set  $F_j^{a_j^0}$  and  $F_j^{a_j^{N+1}}$  as the points in  $PF_p$  where  $F_j$  is largest and smallest, respectively. When  $N \geq 2$ , the spread metrics  $\Gamma$  and  $\Delta$  are defined by

$$\Gamma_{p,s} := \max_{j=1, \dots, m} \max_{i=0, \dots, N} \delta_{i,j}$$

and

$$\Delta_{p,s} := \max_{j=1, \dots, m} \frac{\delta_{0,j} + \delta_{N,j} + \sum_{i=1}^N |\delta_{i,j} - \bar{\delta}_j|}{\delta_{0,j} + \delta_{N,j} + (N-1)\bar{\delta}_j},$$

where  $\delta_{i,j} := F_j^{a_j^{i+1}} - F_j^{a_j^i}$  and  $\bar{\delta}_j := \sum_{i=1}^N \delta_{i,j}/N$ . On the other hand when  $N \leq 1$ , we define  $\Gamma_{p,s} = \Delta_{p,s} = \infty$ .

We also obtained performance profiles [15] for each of the evaluation metrics to provide a comprehensive comparison of the algorithms. Suppose that a metric  $t_{p,s}$  is defined for a solver  $s \in \mathcal{S}$  and a problem  $p \in \mathcal{P}$ . We assume that the smaller  $t_{p,s}$  is, the better. The performance profile  $R_s(\tau)$  of a solver  $s \in \mathcal{S}$  is defined as

$$R_s(\tau) := \frac{1}{|\mathcal{P}|} |p \in \mathcal{P} \mid r_{p,s} \leq \tau|,$$

where  $r_{p,s}$  is the performance ratio given by  $r_{p,s} := t_{p,s}/\min_{s \in \mathcal{S}} t_{p,s}$ . Note that for hypervolume and spread metrics, we took the reciprocal when calculating the performance ratio, as larger metric values correspond to better performance for them.

## 7.4 Results of the experiments

Let us first illustrate the behaviour of the algorithms. For this, we take the problem JOS1 [27] with  $n = 50$  and  $g_i$  as the  $\ell_1$ -norm. In Figure 1, we plot the objective function values for  $k = 0$  (i.e., at the initial points),  $k = 10$ , and the terminal points of each algorithm, respectively. The set of terminal points are in fact the Pareto solutions obtained. Here, “Normal”, “Accelerated”, and “Accelerated (without  $f_i(y) - F_i(x)$ )” means, respectively, Algorithm 1, Algorithm 2 and Algorithm 3. As we can see, all the algorithms were able to find a wide range of Pareto solutions in this case. However, the objective function values at  $k = 10$  are smaller when using the accelerated Algorithm 2 and Algorithm 3. Moreover, from Figure 2, we see that Algorithm 2 and Algorithm 3 converge faster than the non-accelerated Algorithm 1, despite oscillations. In this example, we can also see that Algorithm 2 were faster and obtained a more uniform Pareto frontier than Algorithm 3.

We now check the performance of the algorithms. As we explained in Section 7.2, for each problem of Table 1, we run the algorithms with 100 different initial points. Table 2 shows the average of the computational time and iteration counts for each problem. For problems with different values of  $n$ , we just show the smallest and the biggest  $n$  for convenience. From the table, it is possible to see that acceleration is in general more efficient in terms of time. In fact, by checking the performance profiles given in Figure 3 and Figure 4, we observe that our proposed Algorithm 2 performs better in terms of time and iteration counts. It is interesting to see from Table 2, however, that there are cases where Algorithm 2 does not perform well.

Besides the performance, it is usually important to see how good the Pareto frontier is. Thus, once again we show performance profiles, this time for purity (Figure 5), hypervolume (Figure 6), spread metric  $\Gamma$  (Figure 7) and spread metric  $\Delta$  (Figure 8). Clearly, our proposed Algorithm 2 outperforms the other two algorithms, obtaining better Pareto frontiers. We can thus conclude that at least among the test problems considered, Algorithm 2 seem promising both in terms of performance and uniform Pareto frontiers.

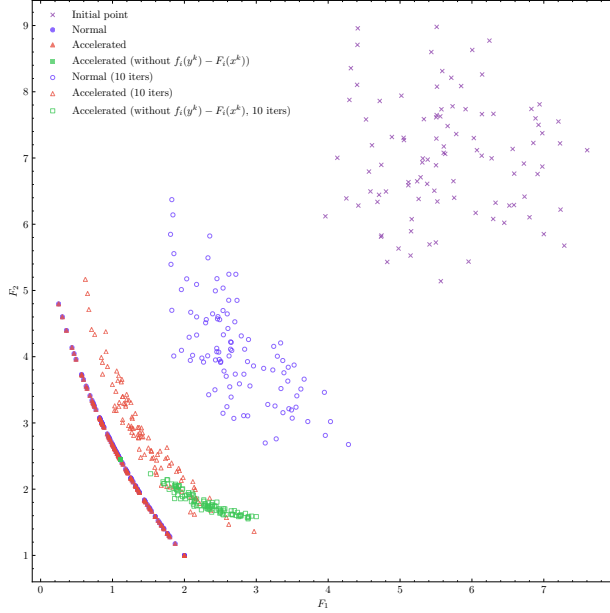


Figure 1: Objective function values for problem JOS1 with  $n = 50$ , and  $\ell_1$  norm for  $g_i$

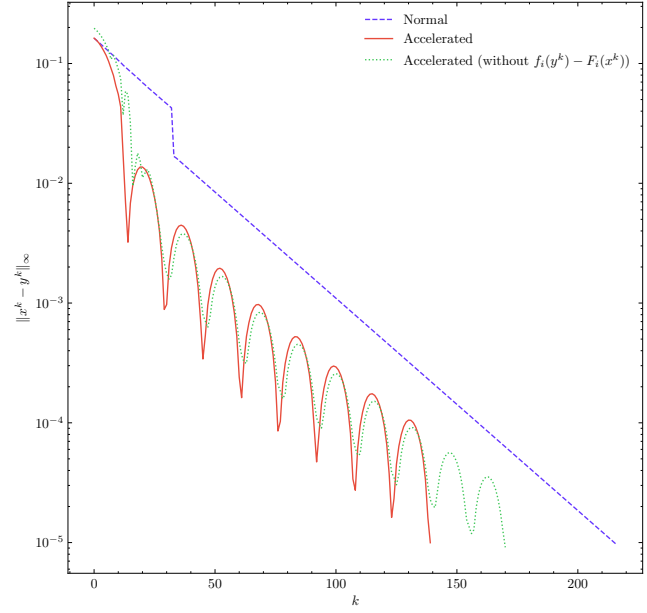


Figure 2: An example of  $\{\|x^k - y^k\|_\infty\}$  for problem JOS1 with  $n = 50$ , and  $\ell_1$  norm for  $g_i$

Table 2: Average computational costs

Problem	$n$	$g_i$	Total time (s)			Iteration counts		
			Alg. 1	Alg. 2	Alg. 3	Alg. 1	Alg. 2	Alg. 3
JOS1	5	0	0.032	0.029	0.027	23.82	27.89	27.89
JOS1	5	$\ell_1$	0.783	0.355	0.336	22.20	21.26	28.09
JOS1	1000	0	3.674	0.260	0.207	3203.05	155.00	155.00
JOS1	1000	$\ell_1$	183.957	47.913	46.197	2901.50	732.72	644.11
ZDT1	5	ind.	0.743	0.279	0.234	38.81	11.03	8.87
ZDT1	1000	ind.	1.738	0.840	0.634	32.68	14.30	9.986
SD	4	ind.	1.063	0.806	1.026	36.58	33.02	32.94
TOI	4	0	0.013	0.018	0.015	3.95	4.57	5.18
TOI	4	$\ell_1$	1.109	0.841	1.035	20.95	18.41	22.90
TRIDIA	3	0	94.447	0.981	5.842	3177.21	6.35	172.89
TRIDIA	3	$\ell_1$	79.562	3.892	78.616	1348.16	25.80	860.90
FDS	5	0	22.897	12.719	14.433	286.4	152.35	170.83
FDS	5	$\ell_1$	16.611	12.150	1330.424	127.48	91.39	13178.25
FDS	100	0	3805.058	4007.926	3607.842	644.45	117.27	158.72
FDS	100	$\ell_1$	4773.474	5412.802	5880.479	767.81	177.37	474.55
LFR1	30	0	4.904	8.335	4.362	9.18	11.67	6.69
LFR1	30	$\ell_1$	10.337	10.184	160.399	8.91	11.4	1224.37
LFR1	1000	0	10.928	13.804	10.916	8.54	10.07	8.36
LFR1	1000	$\ell_1$	26.669	31.038	26.566	8.55	10.31	8.56

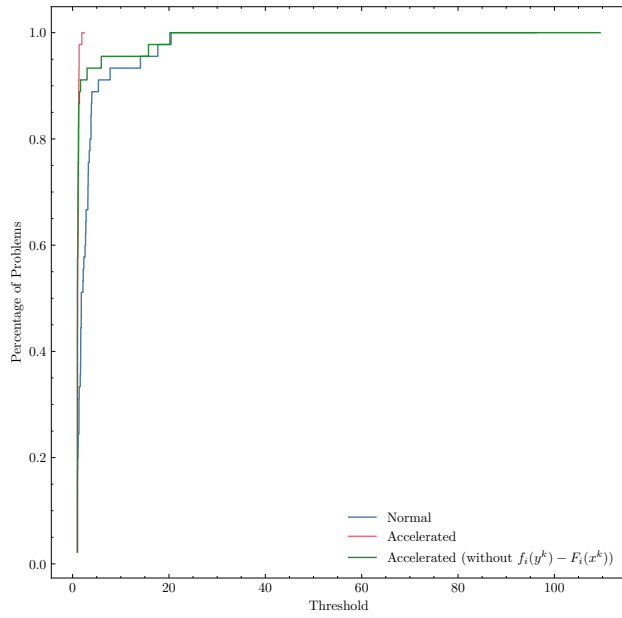


Figure 3: Performance profile:computational time

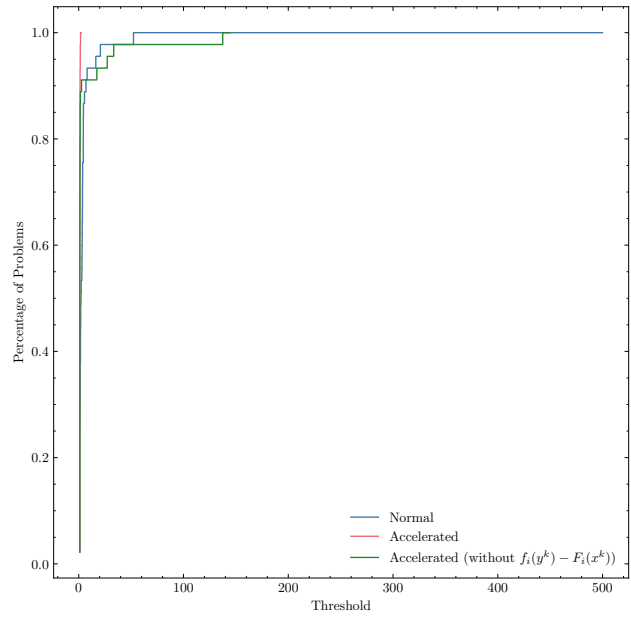


Figure 4: Performance profile:iteration counts

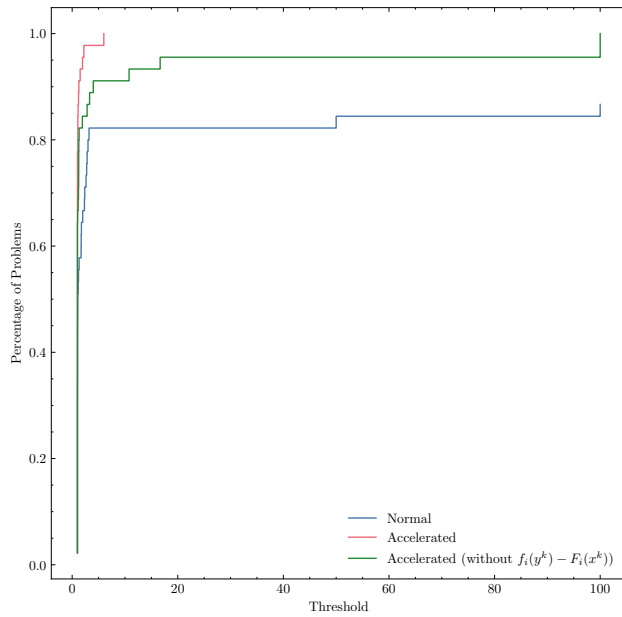


Figure 5: Performance profile:purity

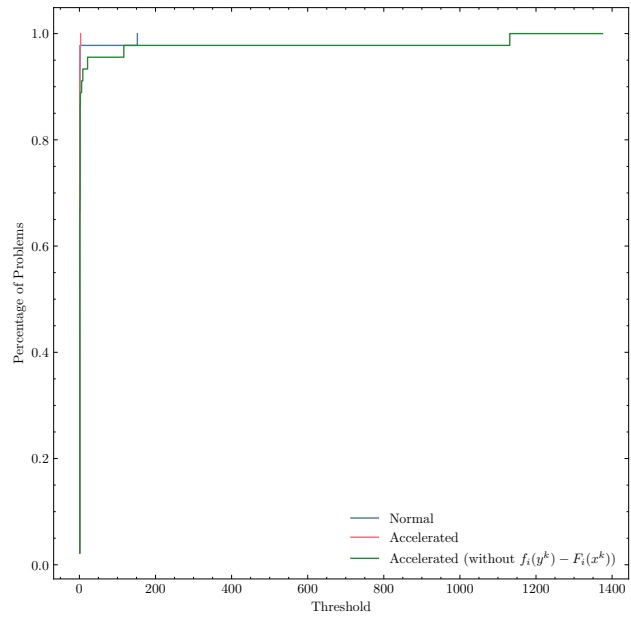


Figure 6: Performance profile: hypervolume

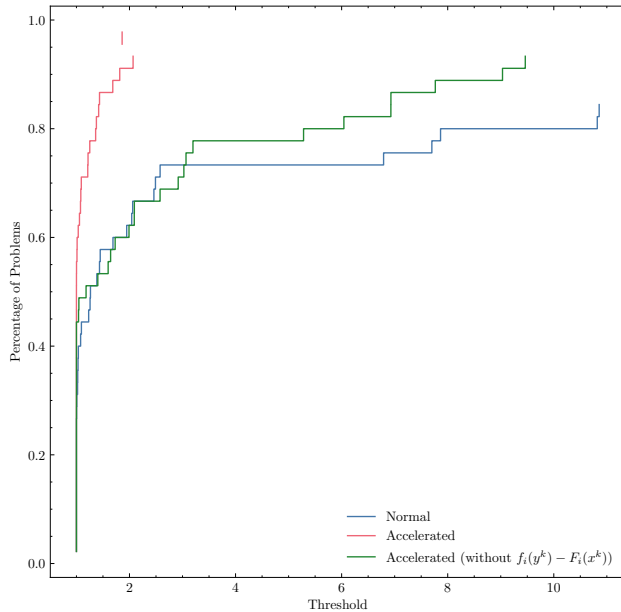


Figure 7: Performance profile:spread metrics ( $\Gamma$ )

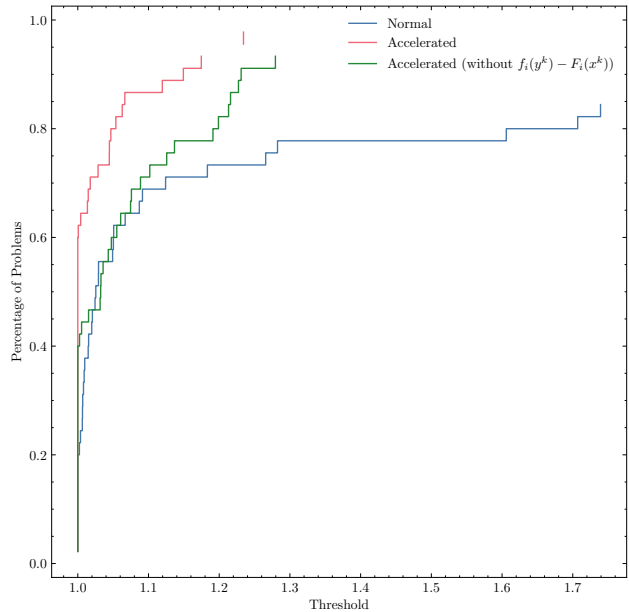


Figure 8: Performance profile:spread metrics ( $\Delta$ )

## 8 Conclusion

By putting information of the previous points into the subproblem, we have successfully accelerated the proximal gradient method for multiobjective optimization and proved its convergence rate under natural assumptions, which was an open problem. Moreover, we showed an efficient way of computing the subproblem via its dual. As the experiments suggested, the proposed methods are also effective from the numerical point of view.

This paper shows the convergence rate for the sequence of the merit function values and the classical global convergence concerning accumulation points but does not provide the global convergence of the sequence of iterates itself. For single-objective optimization, by changing the update rule for the parameter  $t_k$ , Chambolle and Dossal have proposed a variant with the iterates' global convergence [13]. In the multi-objective optimization problem, it may also be possible to modify the algorithm similarly and obtain the global convergence of iterates. Moreover, since many schemes for single-objective optimization had been developed, following the idea of Nesterov's acceleration technique, this paper may also contribute to the development of various multiobjective optimization methods. Extensions to vector optimization and its generalization, the vector optimization problem with variable domination structure [8, 28], may also be worth considering. Such extensions will be subjects of future works.

## Acknowledgments

This work was supported by the Grant-in-Aid for Scientific Research (C) (21K11769 and 19K11840) and Grant-in-Aid for JSPS Fellows (20J21961) from the Japan Society for the Promotion of Science.

## References

- [1] Bandyopadhyay, S., Pal, S. and Aruna, B.: Multiobjective GAs, quantitative indices, and pattern classification, *IEEE Transactions on Systems, Man and Cybernetics, Part B (Cybernetics)*, Vol. 34 (2004), 2088–2099.
- [2] Beck, A.: *First-Order Methods in Optimization*, Society for Industrial and Applied Mathematics, Philadelphia, Pennsylvania, USA, 2017.
- [3] Beck, A. and Teboulle, M.: A fast iterative shrinkage-thresholding algorithm for linear inverse problems, *SIAM Journal on Imaging Sciences*, Vol. 2 (2009), 183–202.

- [4] Bertsekas, D. P.: *Nonlinear Programming*, Athena Scientific, Belmont, Massachusetts, second edition, 1999.
- [5] Boţ, R. I. and Grad, S. M.: Inertial forward-backward methods for solving vector optimization problems, *Optimization*, Vol. 67 (2018), 959–974.
- [6] Bonnans, J. F. and Shapiro, A.: *Perturbation Analysis of Optimization Problems*, Springer New York, New York, NY, USA, 2000.
- [7] Bonnel, H., Iusem, A. N. and Svaiter, B. F.: Proximal methods in vector optimization, *SIAM Journal on Optimization*, Vol. 15 (2005), 953–970.
- [8] Bouza, G. and Tammer, C.: A steepest descent-like method for vector optimization problems with variable domination structure, *Journal of Nonlinear and Variational Analysis*, Vol. 6 (2022), 605–618.
- [9] Boyd, S. and Vandenberghe, L.: *Convex Optimization*, Cambridge University Press, Cambridge, England, 2004.
- [10] Brent, R. P.: *Algorithms for Minimization without Derivatives*, Prentice-Hall, New Jersey, 1973.
- [11] Byrd, R. H., Hribar, M. E. and Nocedal, J.: An interior point algorithm for large-scale nonlinear programming, *SIAM Journal on Optimization*, Vol. 9 (1999), 877–900.
- [12] Carrizo, G. A., Lotito, P. A. and Maciel, M. C.: Trust region globalization strategy for the nonconvex unconstrained multiobjective optimization problem, *Mathematical Programming*, Vol. 159 (2016), 339–369.
- [13] Chambolle, A. and Dossal, C.: On the convergence of the iterates of the “Fast Iterative Shrinkage/Thresholding Algorithm”, *Journal of Optimization Theory and Applications*, Vol. 166 (2015), 968–982.
- [14] Custódio, A. L., Madeira, J. F., Vaz, A. I. and Vicente, L. N.: Direct multisearch for multiobjective optimization, *SIAM Journal on Optimization*, Vol. 21 (2011), 1109–1140.
- [15] Dolan, E. D. and Moré, J. J.: Benchmarking optimization software with performance profiles, *Mathematical Programming, Series B*, Vol. 91 (2002), 201–213.
- [16] El Moudden, M. and El Mouatasim, A.: Accelerated diagonal steepest descent method for unconstrained multiobjective optimization, *Journal of Optimization Theory and Applications*, Vol. 188 (2021), 220–242.
- [17] Fliege, J., Graña Drummond, L. M. and Svaiter, B. F.: Newton’s method for multiobjective optimization, *SIAM Journal on Optimization*, Vol. 20 (2009), 602–626.
- [18] Fliege, J. and Svaiter, B. F.: Steepest descent methods for multicriteria optimization, *Mathematical Methods of Operations Research*, Vol. 51 (2000), 479–494.
- [19] Fliege, J., Vaz, A. I. F. and Vicente, L. N.: Complexity of gradient descent for multiobjective optimization, *Optimization Methods and Software*, Vol. 34 (2019), 949–959.
- [20] Fukuda, E. H. and Graña Drummond, L. M.: Inexact projected gradient method for vector optimization, *Computational Optimization and Applications*, Vol. 54 (2013), 473–493.
- [21] Fukuda, E. H. and Graña Drummond, L. M.: A survey on multiobjective descent methods, *Pesquisa Operacional*, Vol. 34 (2014), 585–620.
- [22] Gandibleux, X., Sevaux, M., Sörensen, K. and T’kindt, V.: *Metaheuristics for Multiobjective Optimisation*, Vol. 535 of *Lecture Notes in Economics and Mathematical Systems*, Springer Berlin Heidelberg, Berlin, Heidelberg, 2004.
- [23] Gass, S. and Saaty, T.: The computational algorithm for the parametric objective function, *Naval Research Logistics Quarterly*, Vol. 2 (1955), 39–45.
- [24] Geoffrion, A. M.: Proper efficiency and the theory of vector maximization, *Journal of Mathematical Analysis and Applications*, Vol. 22 (1968), 618–630.

- [25] Gonçalves, M. L. N., Lima, F. S. and Prudente, L. F.: Globally convergent Newton-type methods for multiobjective optimization, *Computational Optimization and Applications*, Vol. 83 (2022), 403–434.
- [26] Graña Drummond, L. M. and Iusem, A. N.: A projected gradient method for vector optimization problems, *Computational Optimization and Applications*, Vol. 28 (2004), 5–29.
- [27] Jin, Y., Olhofer, M. and Sendhoff, B.: Dynamic weighted aggregation for evolutionary multi-objective optimization: Why does it work and how?, in *Proceedings of the 3rd Annual Conference on Genetic and Evolutionary Computation*, GECCO'01, San Francisco, CA, USA, 2001, Morgan Kaufmann Publishers Inc.
- [28] Köbis, E., Köbis, M. A. and Tammer, C.: A first bibliography on set and vector optimization problems with respect to variable domination structures, *Journal of Nonlinear and Variational Analysis*, Vol. 6 (2022), 725–735.
- [29] Lucambio Pérez, L. R. and Prudente, L. F.: Nonlinear conjugate gradient methods for vector optimization, *SIAM Journal on Optimization*, Vol. 28 (2018), 2690–2720.
- [30] Marler, R. T. and Arora, J. S.: The weighted sum method for multi-objective optimization: New insights, *Structural and Multidisciplinary Optimization*, Vol. 41 (2010), 853–862.
- [31] Mita, K., Fukuda, E. H. and Yamashita, N.: Nonmonotone line searches for unconstrained multiobjective optimization problems, *Journal of Global Optimization*, Vol. 75 (2019), 63–90.
- [32] Moré, J. J., Garbow, B. S. and Hillstom, K. E.: Testing unconstrained optimization software, *ACM Transactions on Mathematical Software*, Vol. 7 (1981), 17–41.
- [33] Moreau, J.-J.: Proximité et dualité dans un espace hilbertien, *Bulletin de la Société Mathématique de France*, Vol. 93 (1965), 273–299.
- [34] Nesterov, Y.: A method for solving the convex programming problem with convergence rate  $O(1/k^2)$ , *Doklady Akademii Nauk SSSR*, Vol. 269 (1983), 543–547.
- [35] Parikh, N. and Boyd, S.: *Proximal Algorithms*, Vol. 1, Now Publishers, Inc., Boston - Delft, 2014.
- [36] Rockafellar, R. T. and Wets, R. J. B.: *Variational Analysis*, Vol. 317 of *Grundlehren der mathematischen Wissenschaften*, Springer Berlin Heidelberg, Berlin, Heidelberg, 1998.
- [37] Sion, M.: On general minimax theorems, *Pacific Journal of Mathematics*, Vol. 8 (1958), 171–176.
- [38] Stadler, W. and Dauer, J.: Multicriteria optimization in engineering: a tutorial and survey, in Kamat, M. P. ed., *Progress in Aeronautics and Astronautics: Structural Optimization: Status and Promise*, Vol. 150, Washington DC, 1992, American Institute of Aeronautics and Astronautics.
- [39] Svaiter, B. F.: The multiobjective steepest descent direction is not Lipschitz continuous, but is Hölder continuous, *Operations Research Letters*, Vol. 46 (2018), 430–433.
- [40] Tanabe, H., Fukuda, E. H. and Yamashita, N.: Proximal gradient methods for multiobjective optimization and their applications, *Computational Optimization and Applications*, Vol. 72 (2019), 339–361.
- [41] Tanabe, H., Fukuda, E. H. and Yamashita, N.: Convergence rates analysis of a multiobjective proximal gradient method, *Optimization Letters*, Vol. 17 (2023), 333–350.
- [42] Tanabe, H., Fukuda, E. H. and Yamashita, N.: New merit functions for multiobjective optimization and their properties, arXiv:2010.09333, 2023.
- [43] Toint, P. L.: Test problems for partially separable optimization and results for the routine PSPMIN, Namur Report, 1983.
- [44] Wang, X., Wang, Y. and Wang, G.: An accelerated augmented lagrangian method for multi-criteria optimization problem, *Journal of Industrial and Management Optimization*, Vol. 16 (2020), 1–9.
- [45] Yosida, K.: *Functional Analysis*, Vol. 123 of *Classics in Mathematics*, Springer Berlin Heidelberg, Berlin, Heidelberg, sixth edition, 1995.



- [46] Zadeh, L. A.: Optimality and non-scalar-valued performance criteria, *IEEE Transactions on Automatic Control*, Vol. 8 (1963), 59–60.
- [47] Zhao, X., Jolaoso, L. O., Shehu, Y. and Yao, J.-C.: Convergence of a nonmonotone projected gradient method for nonconvex multiobjective optimization, *Journal of Nonlinear and Variational Analysis*, Vol. 5 (2021), 441–457.
- [48] Zitzler, E., Deb, K. and Thiele, L.: Comparison of multiobjective evolutionary algorithms: empirical results, *Evolutionary Computation*, Vol. 8 (2000), 173–195.
- [49] Zitzler, E., Technische, E. and Zürich, H.: *Evolutionary Algorithms for Multiobjective Optimization: Methods and Applications*, Phd thesis, Swiss Federal Institute of Technology Zurich, 1999.
- [50] Zou, H. and Hastie, T.: Regularization and variable selection via the elastic net, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, Vol. 67 (2005), 301–320.