# Optimization on Pareto sets:
# On a theory of multi-objective optimization

Abhishek Roy*

a2roy@ucsd.edu

Geelon So*

geelon@ucsd.edu

Yi-An Ma

yianma@ucsd.edu

University of California, San Diego

La Jolla, CA 92093

August 7, 2023

## Abstract

In multi-objective optimization, a single decision vector must balance the trade-offs between many objectives. Solutions achieving an optimal trade-off are said to be Pareto optimal—these are decision vectors for which improving any one objective must come at a cost to another. But as the set of Pareto optimal vectors can be very large, we further consider a more practically significant *Pareto-constrained optimization problem*, where the goal is to optimize a preference function constrained to the Pareto set.

We investigate local methods for solving this constrained optimization problem, which poses significant challenges because the constraint set is (i) implicitly defined, and (ii) generally non-convex and non-smooth, even when the objectives are. We define notions of optimality and stationarity, and provide an algorithm with a last-iterate convergence rate of $O(K^{-1/2})$ to stationarity when the objectives are strongly convex and Lipschitz smooth.

## 1   Introduction

The theory of optimization has provided the foundations for analyzing large-scale machine learning, giving us a language for understanding not only training accuracy, but also generalization (Hardt and Recht, 2022) and adaptive decision making (Puterman, 1994). However, in practice, we often need to further account for additional desiderata: resource constraints, fairness, fine-tunability, and so on. As a result, *multi-objective optimization* (MOO) has increasingly drawn interest from the machine learning community, since it naturally generalizes the single objective paradigm of classical learning while also being able to attend to these additional requirements.

Examples of machine learning settings formulated as MOO problems include those with multiple tasks (Sener and Koltun, 2018; Doersch and Zisserman, 2017), different data distributions (Dong et al., 2015; Huang et al., 2015), fairness requirements (Martinez et al., 2020; La Cava, 2023; Kamani et al., 2021), inverse reinforcement learning (Pirotta and Restelli, 2016), and the need to balance compute and power consumption among multiple algorithmic modules (Ghosh et al., 2013).

The solutions to MOO problems are those that achieve optimal trade-offs, or *Pareto optimality*; together, they form the *Pareto set*. But because the Pareto set generally does not contain a single solution, there is a need to make a further selection from the Pareto optimal solutions. Currently, there are two main approaches to making this selection. The first is to find a representative subsample

of the Pareto set: this reduces the number of solutions that need to be inspected before making a final decision (Lin et al., 2019; Liu et al., 2021; Kobayashi et al., 2019; Guerreiro et al., 2021). The other approach is to scalarize the multiple objectives into a single objective, say, by taking a linear combinations of the objectives (Mahapatra and Rajan, 2020).

However, as the number of objectives and dimensions increase, the Pareto set can become extremely large, forcing the size of a representative subsample to also become untenably large. Furthermore, the geometry of the Pareto set can be quite complicated, with "needle-like extensions" and "knees" (Kulkarni et al., 2022), which poses difficulties for sampling. Even with quadratic objectives in two dimensions, we can observe singularities in the Pareto set, see Sheftel et al. (2013) or Figure 1. The other scalarization approach is also not without difficulties. As the objective functions can be incomparable, it can be challenging to find a meaningful weighting of the objectives.

For a more principled selection, we assume that we have an additional *preference function* $f_0$, which we aim to optimize constrained to the Pareto set. In supervised learning tasks, this preference function is oftentimes the loss function of a generic dataset. In economic and decision making problems, it is usually taken to be the social welfare of the entire community of users. This approach has been considered in various contexts such as portfolio management (Thach et al., 1996) and manufacturing planning (Yamamoto, 2002), in addition to machine learning and optimization (Ye and Liu, 2022). While heuristics have been proposed, little is known about the convergence properties of these algorithms. This prompts us to ask:

*Given a set of objective functions $(f_1, \ldots, f_n)$ and a preference function $f_0$, what is a suitable approximate solution concept and what are efficient algorithms to achieve it?*

## 1.1 Main results

In MOO, we are given a set of $n$ objective functions $F \equiv (f_1, \ldots, f_n) : \mathbb{R}^d \to \mathbb{R}^n$ that are jointly optimized over a shared decision space $\mathbb{R}^d$:

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad F(x). \tag{1}$$

The solution concept for (1) is typically defined as the set of Pareto optimal solutions, $\text{Pareto}(F)$, which consists of decision vectors $x \in \mathbb{R}^d$ that make an optimal trade-off between objectives. And to further decide which trade-off to make, we consider the *Pareto-constrained optimization problem*, in which the aim is to optimize a preference function $f_0 : \mathbb{R}^d \to \mathbb{R}$ constrained to the Pareto set of $F$:

$$\underset{x \in \text{Pareto}(F)}{\text{minimize}} \quad f_0(x). \tag{2}$$

This problem is challenging not only because the constraint set is defined implicitly as the solution to the MOO problem from Equation (1), but because it is also non-convex and non-smooth. Even in the case of linear preference functions, the problem is known to be NP-hard (Fülöp, 1993). In fact, it is not obvious how to even define an appropriate relaxation of the problem such as stationarity that can be attained through optimization, given the challenges of non-convex non-smooth optimization (Zhang et al., 2020; Kornowski and Shamir, 2021; Li et al., 2020; Jordan et al., 2023). However, we show in this work that the Pareto set has additional geometry when the objectives are strongly convex that allows us to relax the Pareto-constrained optimization problem to a strong notion of stationarity that is necessary for optimality and that can be efficiently attained:

1. We show that the Pareto-constrained optimization problem has an equivalent reformulation as a smooth optimization problem over a linear constraint set (Proposition 1). This allows us to
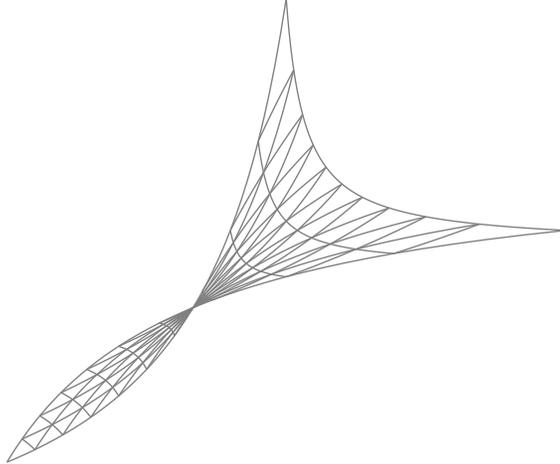
2

Figure 1: A Pareto set for three positive-definite quadratic objective functions in $\mathbb{R}^2$. The grid lines show the coordinate maps for $x^* : \Delta^2 \to \mathrm{Pareto}(f_1, f_2, f_3)$, where $\Delta^2$ is the 3-simplex. Even in this well-structured setting, the Pareto set is not convex or smooth.

   introduce solution concepts such as (approximate) *preference stationarity* in the standard way. Furthermore, we show that the solution concepts are geometrically meaningful (Proposition 4).

2. While the reformulation solves the issue of non-convexity and non-smoothness, the reformulated objective function remains implicit, which can make it hard to design optimization methods and provide simple analyses. If the objectives and preference are sufficiently smooth (Assumptions A–C), we construct a family of upper bounds for the reformulated objective function (Proposition 5), providing a general framework to analyze iterative gradient-based methods.

3. We provide the *Pareto majorization-minimization* algorithm, which iteratively (i) computes these upper bounds and (ii) minimizes them. In our setting, this amounts to solving a sequence of (i) unconstrained strongly-convex optimization problems and (ii) quadratic programs. We show that it suffices to solve the strongly convex programs up to $O(\varepsilon_0^2)$-optimality and the quadratic programs up to $O(\varepsilon_0)$-stationarity. Then, no more than $O(\varepsilon_0^{-2})$ rounds of optimization are needed to attain an $\varepsilon_0$-approximate preference stationary solution (Theorem 1).

## 1.2 Related work

Selecting a single decision out of all Pareto optimal decisions is a fundamental problem of MOO that does not appear in the classical single-objective setting; in MOO, there is no canonical total ordering of the solutions (Miettinen, 1999). Broadly, the approaches to making such a selection can be categorized as *a priori* and *a posteriori* (Hwang and Masud, 2012).[1]

   In the *a priori* setting, the preferences of the decision maker is known beforehand. While in the *a posteriori* approach, the goal is to present a decision maker with a representative spread of Pareto optimal options, from which the decision maker will make a final decision. But because the Pareto set can become very high-dimensional, the *a posteriori* approach becomes less viable (or needs to become more interactive) as the number of objectives and dimensions increase.

   Instead, we work in the *a priori* setting and consider optimization constrained to the Pareto set, also sometimes called *semivectorial bilevel optimization* or *optimization on efficient sets*, which can

---

[1]They also include two other categories: the *no-preference* and *interactive* approach. In the former, any Pareto optimal decision will do, while in the latter, candidates are presented adaptively to an interactive decision maker.

be considered an instantiation of bi-level optimization (Yamamoto, 2002; Bonnel and Morgan, 2006; Dempe, 2018). This problem is known to be NP-hard (Fülöp, 1993) and algorithms for solving this problem tend to focus on settings with: (i) linear preference functions (Philip, 1972; Benson, 1984; Liu and Ehrgott, 2018), (ii) linear objectives (Dauer, 1991; Bolintineanu, 1993; Tao et al., 1996; Yamamoto, 2002), or (iii) specific choices of preference functions such as the Tchebycheff weighting function (Steuer, 1989). To our knowledge, the only other algorithmic work that considers the general problem with nonlinear objectives is Ye and Liu (2022).

However, the stationary condition introduced by Ye and Liu (2022), defined as stationarity with respect to the proposed optimization dynamics, does not have a clear connection to preference optimality. In fact, as it is a non-trivial first-order stationary condition, the stationarity notion defined therein is not a necessary condition (Proposition 3); there are settings where such dynamics avoid optimal points (see Example 1).

We are able to introduce a simple and necessary condition for preference optimality by making use of the manifold structure of the Pareto set. While its smooth structure has previously been recognized (Hillermeier, 2001; Hamada et al., 2020), the prior focus has been on the extrinsic Pareto manifold within an ambient space, from which it inherits its smoothness. We take a different approach and work with the Pareto manifold intrinsically. Since it is diffeomorphic to the simplex, conceptually, this greatly simplifies optimization constrained to the Pareto set. And in order to overcome the implicit nature of the Pareto manifold, we use ideas from majorization-minimization and trust-region approaches to optimization, where approximate gradient information can be used to make provable improvements (Lange et al., 2000; Marumo et al., 2023).

## 2    Preliminaries

Let $f_1, \ldots, f_n : \mathbb{R}^d \to \mathbb{R}$ be objective functions, $f_0 : \mathbb{R}^d \to \mathbb{R}$ be a preference function. We assume:

(A) The objectives are strongly convex and twice-differentiable with Lipschitz-continuous gradients.

(B) The objectives have Lipschitz-continuous Hessians.

(C) The preference has Lipschitz-continuous gradients.

Let $[n] := \{1, \ldots, n\}$. We denote the $(n-1)$-simplex by $\Delta^{n-1}$, which is the set of convex weights:

$$\Delta^{n-1} := \left\{ \beta \in \mathbb{R}^n : \sum_{i \in [n]} \beta_i = 1 \text{ and } \forall i \in [n], \ \beta_i \geq 0 \right\}.$$

And given a convex weight $\beta \in \Delta^{n-1}$, we let $f_\beta$ denote the *scalarization*:

$$f_\beta(x) := \sum_{i \in [n]} \beta_i f_i(x). \tag{3}$$

For a detailed glossary, see Section 9.

Let us recall the definition of a Pareto optimal decision vector.

**Definition 1** (Pareto optimality). *Given objectives $f_1, \ldots, f_n$, we say that a decision vector $x \in \mathbb{R}^d$ is* Pareto optimal *if for all $x' \in \mathbb{R}^d$:*

$$f_i(x') < f_i(x) \qquad \Longrightarrow \qquad \exists j \quad \text{s.t.} \quad f_j(x') > f_j(x).$$

*We call the set of Pareto optimal decision vectors the* Pareto set *of $f_1, \ldots, f_n$, denoted* Pareto(F).

In words, the above condition states that there is no way to improve $f_i$ without also worsening some other $f_j$. When the objectives are smooth, a related local condition is Pareto stationarity:

**Definition 2** (Pareto stationarity). *Given objectives $f_1, \ldots, f_n$, we say that a decision vector $x \in \mathbb{R}^d$ is* Pareto stationary *if zero is a convex combination of the gradients:*

$$\nabla f_\beta(x) = 0, \qquad \textit{for some } \beta \in \Delta^{n-1},$$

*where $f_\beta$ is defined by Equation* (3).

In particular, Pareto stationarity is a necessary condition for Pareto optimality (Maruşciac, 1982). Furthermore, it is sufficient when the objectives are twice-differentiable and are strictly convex (Fliege et al., 2009). As we have assumed this, we have:

$$x \in \text{Pareto}(F) \quad \Longleftrightarrow \quad x \text{ is Pareto stationary.}$$

# 3  The Pareto manifold

It is not immediately evident from the definition of Pareto stationarity that $\text{Pareto}(F)$ is amenable to the Pareto-constrained optimization problem defined in Equation (2). In general, the Pareto set is non-smooth and non-convex. Even when the objectives are positive-definite quadratics, the Pareto set can have singularities (Sheftel et al., 2013). For example, see the Pareto set in Figure 1.

This issue of non-smoothness arises because the set of Pareto stationary points naturally lives in a higher-dimensional space $\mathbb{R}^d \times \Delta^{n-1}$, in which it is a smooth $(n-1)$-dimensional submanifold. But when it is projected back down into $\mathbb{R}^d$, it can cross itself to create singularities. Formally, we define:

**Definition 3** (Pareto manifold). *The* Pareto manifold $\mathcal{P}(F) \subset \mathbb{R}^d \times \Delta^{n-1}$ *is the zero set:*

$$\mathcal{P}(F) = \big\{(x, \beta) : \nabla f_\beta(x) = 0\big\}.$$

The Pareto manifold consists of all $(x, \beta)$ such that $x$ is Pareto stationary and $\beta$ bears witness to the stationarity condition $\nabla f_\beta(x) = 0$. And of course, we can recover the Pareto set from the Pareto manifold simply by projecting down to its first component in $\mathbb{R}^d$:

$$x \in \text{Pareto}(F) \quad \Longleftrightarrow \quad (x, \beta) \in \mathcal{P}(F) \text{ for some } \beta \in \Delta^{n-1}.$$

But this projection can also collapse any smoothness structure that $\mathcal{P}(F)$ has. And indeed, it is a smooth submanifold of $\mathbb{R}^d \times \Delta^{n-1}$. To see this, notice that $\mathcal{P}(F)$ is the zero set of the map:

$$(x, \beta) \mapsto \nabla f_\beta(x),$$

whose partial derivative with respect to $x$ is invertible—the partial derivative is $\nabla^2 f_\beta$, which is positive-definite by strong convexity. The implicit function theorem then yields its manifold structure:

**Proposition 1** (Characterization of the Pareto manifold). *Define the map $x^* : \Delta^{n-1} \to \text{Pareto}(F)$:*

$$x^*(\beta) \equiv x_\beta := \underset{x \in \mathbb{R}^d}{\arg\min} \; f_\beta(x). \tag{4}$$

*Let $\nabla F(x) \in \mathbb{R}^{n \times d}$ be the Jacobian. Then, the map $x^*$ has derivative:*

$$\nabla x^*(\beta) = -\nabla^2 f_\beta(x_\beta)^{-1} \nabla F(x_\beta)^\top, \tag{5}$$

*so that the map $\beta \mapsto (x_\beta, \beta)$ is a diffeomorphism of $\Delta^{n-1}$ with the Pareto manifold $\mathcal{P}(F)$.*

Thus, one natural set of coordinates for the Pareto manifold is its parametrization by the simplex. This allows us to define an equivalent but smooth formulation of the Pareto-constrained optimization problem obtained by pulling $f_0$ back onto $\Delta^{n-1}$, which we shall now do.

5

# 4  The Pareto-constrained optimization problem

The Pareto-constrained optimization problem defined in Equation (2) has another formulation:

$$\underset{(x,\beta)\in\mathcal{P}(F)}{\text{minimize}}\ f_0(x), \tag{6}$$

where the constraint has been replaced with the Pareto manifold. The two are equivalent because $\text{Pareto}(F)$ is precisely the projection of $\mathcal{P}(F)$ onto $\mathbb{R}^d$. But the reformulation allows us to apply Proposition 1 to pullback the optimization problem onto $\Delta^{n-1}$:

$$\underset{\beta\in\Delta^{n-1}}{\text{minimize}}\ (f_0\circ x^*)(\beta), \tag{7}$$

which is now a smooth optimization problem over the simplex. We say that $x$ is *preference optimal* if it is a solution to (6); if $\beta$ solves (7), then correspondingly, $x^*(\beta)$ is preference optimal.

As $f_0$ and $x^*$ are smooth, so too is their composition $(f_0\circ x^*)$; we can define a stationarity condition in the standard way for smooth objectives on convex sets (Nesterov, 2003). We say that $x$ is *weakly preference stationary* if there is some $\beta$ such that $(x,\beta)\in\mathcal{P}(F)$ and $\beta$ is stationary in the usual sense for (7). For any given $x$, many $\beta$'s could satisfy the condition $(x,\beta)\in\mathcal{P}(F)$,

$$\Delta^{n-1}(x) := \big\{\beta\in\Delta^{n-1} : \nabla f_\beta(x)=0\big\}. \tag{8}$$

We say that $x$ is *preference stationary* if the stationary condition holds for all such $(x,\beta)$'s.

**Definition 4** (Preference stationarity). *We say that a point $x\in\text{Pareto}(F)$ is* weakly preference stationary *if there exists some $\beta\in\Delta^{n-1}(x)$ such that:*[2]

$$-\nabla(f_0\circ x^*)(\beta)^\top(\beta'-\beta)\le 0, \qquad \forall\beta'\in\Delta^{n-1}, \tag{9}$$

*where Equation (5) gives $\nabla x^*$. If (9) holds for all $\beta\in\Delta^{n-1}(x)$, then $x$ is* preference stationary.

From optimization on convex sets (Nesterov, 2003, Lemma 3.1.19), we immediately have:

**Proposition 2** (Necessary condition). *Preference optimality implies (weak) preference stationarity.*

While this definition of preference stationarity is appealing because it is necessary for preference optimality and because it is well-founded in standard optimization theory, it is not necessarily the only reasonable relaxation of preference optimality. For example, our notion of preference stationarity requires second-order information in $F$ for the term $\nabla^2 f_\beta(x)^{-1}$. It is natural to ask whether we could define stationarity with reference to only first-order information. It turns out that this is impossible, if we require the stationarity condition to be (i) non-trivial, (ii) necessary for preference optimality and (iii) decidable from local information at a single point $x$.

The reason is that the local behavior of the Pareto set about a point $x$ cannot be determined from $\nabla F(x)$ alone. Figure 2 shows two different Pareto sets that share the same first-order information at a point $x$. But the preference stationarity of $x$ with respect to $f_0$ also depends on its neighboring Pareto points. So to attain a non-trivial and necessary condition, we would either need to look at higher-order information or more than a one point. To formalize this, first define:

**Definition 5** (Preference genericity). *Let $\{v_0,v_1,\ldots,v_n\}\subset\mathbb{R}^d$ where $1<n\le d$. We say that this set is* preference generic *if there is a unique $\beta\in\Delta^{n-1}$ such that $\beta_1 v_1+\cdots\beta_n v_n=0$, and:*

$$v_0\notin\text{span}(v_1,\ldots,v_n).$$

---

[2]As $\nabla F(x_\beta)^\top\beta=\nabla f_\beta(x_\beta)=0$, Equation (9) can be simplified to $-\nabla(f_0\circ x^*)(\beta)^\top\beta'\le 0$, for all $\beta'\in\Delta^{n-1}$.
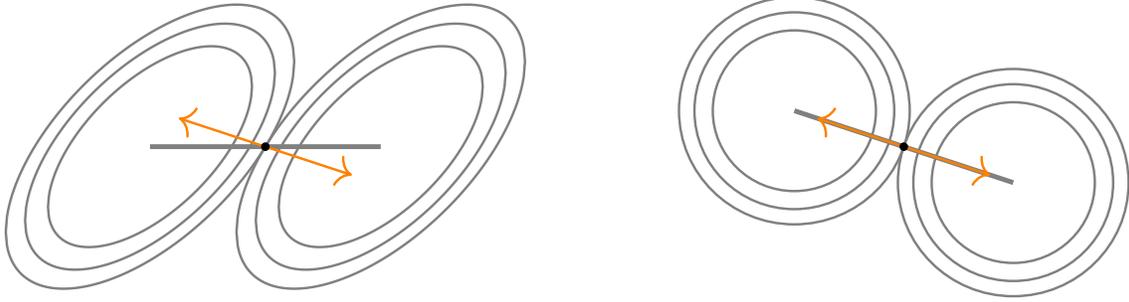
Figure 2: Two instances of Pareto($f_1, f_2$) are shown (thick gray lines), where $f_1$ and $f_2$ are positive-definite quadratic objectives in $\mathbb{R}^2$ (visualized by contour lines). At $x$ (the black dot), the two instances share the same local information $-\nabla f_1(x)$ and $-\nabla f_2(x)$ (orange arrows); they cross the contour lines at right angles. When $n = 2$, the Pareto set contains all $z$ such that $\nabla f_1(z) = -\lambda \nabla f_2(z)$ for $\lambda \geq 0$. Notice that if $f_0$ is strictly convex and $x$ does not minimize $f_0$ over $\mathbb{R}^2$, then $x$ cannot be stationary for both instances.

We also formalize stationarity conditions as *decision functions*, which are functions mapping continuous inputs to Boolean variables taking values of `true` or `false`.

**Definition 6** (Stationarity function). *A first-order stationary condition is a decision function:*

$$\text{Stationary} : \mathbb{R}^d \times \overset{n+1 \ times}{\cdots} \times \mathbb{R}^d \to \{\mathsf{true}, \mathsf{false}\}.$$

*Let $f_0$ be a smooth preference function and $f_1, \ldots, f_n$ be smooth, strongly convex objectives. We say that a first-order condition is* necessary *if the following holds:*

$$x \text{ is preference optimal} \quad \implies \quad \text{Stationary}\big(\nabla f_0(x), \ldots, \nabla f_n(x)\big) = \mathsf{true}.$$

**Proposition 3** (Necessary first-order conditions are trivial). *Suppose that* Stationary *is necessary. Then, it is trivial in the following sense: for any preference generic set of $v_0, \ldots, v_n \in \mathbb{R}^d$,*

$$\text{Stationary}(v_0, \ldots, v_n) = \mathsf{true}.$$

# 5 Pareto majorization-minimization

Let us now consider how to solve the Pareto-constrained optimization problem:

$$\underset{\beta \in \Delta^{n-1}}{\text{minimize}} \ (f_0 \circ x^*)(\beta). \tag{7}$$

As the problem has been reformulated as a smooth optimization problem on the simplex, this seems to open up local methods like gradient descent. But for this, there is a remaining issue that $x^*$ is defined implicitly as the solution of another optimization problem:

$$x^*(\beta) \equiv x_\beta := \underset{x \in \mathbb{R}^d}{\arg\min} \ f_\beta(x). \tag{4}$$

Because $x_\beta$ does not generally have a closed form, we also cannot explicitly compute $\nabla x^*(\beta)$, which is required if we wish to compute $\nabla (f_0 \circ x^*)(\beta)$ by the chain rule.

## 5.1 Approximating the gradient

We can, however, approximate the gradient. Define the following estimator, which uses local information $\nabla^2 f_\beta(x)$ and $\nabla F(x)$ at $x$ as a proxy for the corresponding local information at $x_\beta$:

$$\widehat{\nabla} x^*(x, \beta) := -\nabla^2 f_\beta(x)^{-1} \nabla F(x)^\top. \tag{10}$$

If $F$ has continuous second derivative, then $\widehat{\nabla} x^*(x, \beta)$ approaches $\nabla x^*(\beta)$ as $x$ goes to $x_\beta$; strong convexity implies that $\nabla^2 f_\beta$ has a continuous inverse. And so, there are many reasonable approaches to this problem: it is a smooth optimization problem on a convex set with approximate gradients. For example, we could use the gradient estimate to perform projected gradient descent on the simplex.

Then, the questions at hand: (a) how valid is the approximation $\widehat{\nabla} x^*(x, \beta)$, and (b) how can an optimization procedure make use of that information? It is certainly not the case that the approximation computed at $(x, \beta)$ for some distant $x$ should be as equally valid as one computed near $(x_\beta, \beta)$. One way we can capture the validity of the estimator $\widehat{\nabla} x^*(x, \beta)$ is by using it to construct a majorizing surrogate function, which is a function that upper bounds $f_0 \circ x^*$:

**Definition 7** (Majorizing surrogate). *A function $g : \Delta^{n-1} \to \mathbb{R}$ majorizes $f_0 \circ x^*$ if:*

$$f_0(x_{\beta'}) \leq g(\beta'), \tag{11}$$

*for all $\beta' \in \Delta^{n-1}$. We say that $g$ is a* surrogate *of $f$.*

Intuitively, the better the approximation is, the tighter an upper bound we could provably attain. And as an example, suppose that we have $\widehat{\nabla} x^*(x_\beta, \beta)$, which in this case is exactly $\nabla x^*(\beta)$. And suppose that we knew that $f_0 \circ x^*$ were 1-Lipschitz smooth. Then, the standard quadratic upper bound for Lipschitz smooth functions (Nesterov, 2003) yields a family of majorizing surrogates:

$$g(\beta'; x_\beta, \beta) = f_0(x_\beta) + \nabla f_0(x_\beta)^\top \widehat{\nabla} x^*(x_\beta, \beta)(\beta' - \beta) + \frac{1}{2} \|\beta' - \beta\|^2.$$

This means that we could use $g$ to bound how much improvement in $f_0$ is made by any iterative optimization scheme that takes a step from $\beta$ to $\beta'$: we can think of $g(\beta'; x, \beta)$ as extracting information from $\widehat{\nabla} x^*(x, \beta)$ to certify when an update $f_0(x_{\beta'})$ will improve upon $f_0(x_\beta)$.

## 5.2 Algorithms from upper bounds

Assuming we can obtain such bounds, we can use them not only to analyze optimization procedures, but we can also define a broad class of iterative methods that directly optimize the upper bounds. Suppose that we can compute a family of majorizing surrogates indexed over $\mathbb{R}^d \times \Delta^{n-1}$. Then, the idealized *Pareto majorization-minimization* (PMM) algorithm proceeds in rounds:

1. majorization: query $\widehat{\nabla} x^*(x_k, \beta_k)$ to construct a majorizing surrogate $g_k(\beta) \equiv g(\beta; x_k, \beta_k)$,

2. minimization: make updates $\beta_{k+1} \leftarrow \underset{\Delta^{n-1}}{\arg\min} \; g_k(\beta)$ and $x_{k+1} \leftarrow \underset{x \in \mathbb{R}^d}{\arg\min} \; f_{\beta_{k+1}}$.

The majorizing property of $g_k$ ensures that the iterates $f_0(x_{\beta_{k+1}})$ improve as $\beta_{k+1}$ optimizes $g_k$. We also operationalize the intuition that $g(\cdot; x, \beta)$ becomes more informative as $x$ approaches $x_\beta$ by optimizing $f_{\beta_{k+1}}$. Algorithm 1 is obtained by relaxing step 2, for we do not need to fully optimize $g_k$ and $f_{\beta_{k+1}}$, and we allow for any black-box optimizer. In theory, any iterative optimization method could be interpreted as an approximate PMM; this yields one framework for convergence analysis.

8

---

**Algorithm 1** Pareto majorization-minimization (PMM)

---

**Input:** objectives $F \equiv (f_1, \ldots, f_n)$, preference function $f_0$, and black-box optimizer $\widehat{\arg\min}$
**Initialize:** $(\beta_0, x_0) \in \Delta^{n-1} \times \mathbb{R}^d$

1: **for** $k = 1, \ldots, K$ **do**
2:     Compute a majorizing surrogate $g_k(\beta) \equiv g(\beta; x_k, \beta_k)$ satisfying Equation (11)
3:     Compute approximate minimizers

$$\beta_{k+1} \leftarrow \underset{\beta \in \Delta^{n-1}}{\widehat{\arg\min}} \, g_k(\beta) \qquad \text{and} \qquad x_{k+1} \leftarrow \underset{x \in \mathbb{R}^d}{\widehat{\arg\min}} \, f_{\beta_{k+1}}(x).$$

4: **end for**
5: **return** $(\beta_{K+1}, x_{K+1})$

---

# 6   Approximability from smoothness

In this section, we quantify the smoothness assumptions presented in Section 2. From them, we can derive the following implications:

- Assumption A allows us to bound the size of the Pareto set (Lemma 1).

- Assumption B additionally bounds the curvature of the Pareto manifold: we show that $\nabla x^*$ is well-behaved (Lemma 2) and is well-approximated by $\widehat{\nabla} x^*$ (Lemma 3).

- Assumption C further leads to error bounds when approximating gradient of $f_0 \circ x^*$ (Lemma 4). It also allows us to define a notion of approximate preference stationarity that is geometrically meaningful (Proposition 4) and can be verified using approximate information (Lemma 5).

Formally, we have:

**Assumption A.** *Let the objectives $f_1, \ldots, f_n : \mathbb{R}^d \to \mathbb{R}$ be twice differentiable, $\mu$-strongly convex, and have $L$-Lipschitz continuous gradient. That is, for all $i = 1, \ldots, n$,*

$$\mu\mathbf{I} \preceq \nabla^2 f_i(x) \preceq L\mathbf{I}.$$

*Thus, the condition number of $\nabla^2 f_i$ is upper bounded by $\kappa := L/\mu$. We also let $r$ be a scale parameter, defined by the maximum distance between any of the minimizers of the objectives:*

$$r := \max_{i,j \in [n]} \, \big\| \arg\min f_i(x) - \arg\min f_j(x) \big\|_2.$$

**Lemma 1** (Size of Pareto set). *Suppose $F$ satisfies Assumption A. Then $R \leq \sqrt{\kappa}\, r$, where:*

$$R := \mathrm{diam}\big(\mathrm{Pareto}(F)\big) \equiv \sup\big\{\|x - x'\|_2 : x, x' \in \mathrm{Pareto}(F)\big\}.$$

**Assumption B.** *Let the objectives $f_1, \ldots, f_n : \mathbb{R}^d \to \mathbb{R}$ have $L_H$-Lipschitz continuous Hessian. That is, for all $x, y \in \mathbb{R}^d$ and $i = 1, \ldots, n$, we have $\|\nabla^2 f_i(x) - \nabla^2 f_i(y)\|_2 \leq L_H \|x - y\|_2$.*

**Lemma 2** (Smoothness of $x^*$). *Suppose $F$ satisfies Assumptions A,B. Then, $x^* : \Delta^{n-1} \to \mathbb{R}^d$ is $M_0$-Lipschitz continuous and has $M_1$-Lipschitz continuous gradients, where:*

$$M_0 := \kappa R \qquad \text{and} \qquad M_1 := 2\kappa^2 R \left(1 + \frac{L_H R}{\mu}\right).$$

9

**Lemma 3** (Approximability of $\nabla x^*$). *If $F$ satisfies Assumptions A,B. Then:*

$$\|\nabla x^*(\beta) - \widehat{\nabla} x^*(x, \beta)\|_{1,2} \leq \frac{1}{\mu} \frac{M_1}{2M_0} \|\nabla f_\beta(x)\|_2.$$

**Assumption C.** *Let the preference function $f_0 : \mathbb{R}^d \to \mathbb{R}$ have $L_0$-Lipschitz continuous gradient. That is, for all $x, y \in \mathbb{R}^d$, we have $\|\nabla f_0(x) - \nabla f_0(y)\|_2 \leq L_0 \|x - y\|_2$.*

**Lemma 4** (Approximability of $\nabla(f_0 \circ x^*)$). *If $F$ and $f_0$ satisfy Assumptions A,B,C. Then:*

$$\left\| \nabla(f_0 \circ x^*)(\beta)^\top - \nabla f_0(x)^\top \widehat{\nabla} x^*(x, \beta) \right\|_{1,2} \leq \frac{1}{\mu} \left( \frac{M_1}{2M_0} \|\nabla f_0(x)\|_2 + L_0 M_0 \right) \|\nabla f_\beta(x)\|_2.$$

*We denote the right-hand side by $\mathrm{err}_{\nabla f_0}(x, \beta)$.*

## 6.1 An approximate solution concept

In practice, we generally can never exactly recover stationary points, so we further relax our target solution concept to an approximate version of preference stationarity in the standard way (Nesterov, 2013). To define our notion of approximation, we consider $\Delta^{n-1}$ as a metric space. While somewhat arbitrary, it is also fairly natural to endow $\Delta^{n-1}$ with the $\ell_1$-metric, so that it has unit diameter.

**Definition 8** (Approximate preference stationarity). *Let $\varepsilon_0, \varepsilon \geq 0$. A point $(x, \beta) \in \mathbb{R}^d \times \Delta^{n-1}$ is $(\varepsilon_0, \varepsilon)$-preference stationary if:*

$$-\nabla f_0(x_\beta)^\top \nabla x^*(\beta)(\beta' - \beta) \leq \varepsilon_0 \|\beta' - \beta\|_1, \qquad \forall \beta' \in \Delta^{n-1}. \tag{12a}$$

$$\|\nabla f_\beta(x)\|_2 \leq \varepsilon \tag{12b}$$

When the objectives and preference are sufficiently nice, then an approximate preference stationary solution $(\hat{x}, \hat{\beta})$ has an intuitive meaning: (a) there is a ball around $\hat{\beta}$ within which $f_0 \circ x^*$ decreases at most at an $O(\varepsilon_0)$-rate when moving away from $\hat{\beta}$, and (b) the point $\hat{x}$ is $O(\varepsilon)$-close to $x_{\hat{\beta}}$.

**Proposition 4** (Geometric meaning of approximate stationarity). *Let $F$ and $f_0$ satisfy Assumptions A,B,C and let $(\hat{x}, \hat{\beta})$ be $(\varepsilon_0, \varepsilon)$-preference stationary. The following hold:*

*a. if $\|\beta - \hat{\beta}\|_1 \leq s$, then $f_0(x_\beta) - f_0(x_{\hat{\beta}}) \geq -2\varepsilon_0 \|\beta - \hat{\beta}\|_1$, and*

*b. $\|\hat{x} - x_{\hat{\beta}}\|_2 \leq \varepsilon/\mu$,*

*where we let $R$ is defined in Lemma 1 and $s := \frac{2\mu^2 \varepsilon_0}{L_0 L^2 R^2}$.*

**Lemma 5** (Verifiability of approximate stationarity). *Let $F$ and $f_0$ satisfy Assumptions A,B,C. Then $(\hat{x}, \hat{\beta})$ is $(\varepsilon_0, \varepsilon)$-preference stationary if $\|\nabla f_{\hat{\beta}}(\hat{x})\|_2 \leq \varepsilon$, and for some $x \in \mathbb{R}^d$ and $\alpha \in (0, 1)$,*

*1. an $\alpha \cdot \varepsilon_0$-approximate stationary condition holds:*

$$-\nabla f_0(x)^\top \widehat{\nabla} x^*(x, \hat{\beta})(\beta' - \hat{\beta}) \leq \alpha \cdot \varepsilon_0 \|\beta' - \hat{\beta}\|_1, \tag{13}$$

*2. an error bound holds:*

$$\mathrm{err}_{\nabla f_0}(\hat{\beta}, x) \leq (1 - \alpha) \cdot \varepsilon_0.$$

# 7 Analysis of Pareto majorization-minimization

In this section, we give an explicit majorizing family of positive-definite quadratics surrogates, and we provide a condition for when Algorithm 1 converges.

## 7.1 A family of majorizing surrogates

Because $f_0$ and $x^*$ are respectively $L_0$- and $M_1$-Lipschitz smooth (Assumption C and Lemma 3), their composition is also Lipschitz smooth and admits the quadratic upper bound:

$$f_0(x_{\beta'}) \le f_0(x_\beta) + \nabla(f_0 \circ x^*)(\beta)^\top(\beta' - \beta) + \frac{1}{2}nL_0M_1\|\beta' - \beta\|_2^2,$$

where the dimension $n$ in the second term comes from $\|\beta' - \beta\|_1^2 \le n\|\beta' - \beta\|_2^2$. And even though the gradient $\nabla(f_0 \circ x^*)(\beta)^\top$ is implicit, we can approximate it using $\nabla f_0(x)^\top\widehat{\nabla}x^*(x, \beta)$ where the error is bounded by Lemma 4. This implies the following family of majorizing surrogates:

**Proposition 5** (A family of majorizing surrogates). *Suppose $F$ and $f_0$ satisfy Assumptions A,B,C. Let $\mathrm{err}_{\nabla f_0}(x, \beta)$ be as defined in Lemma 4. Define:*

$$g(\beta'; x, \beta) := f_0(x_\beta) + \nabla f_0(x)^\top\widehat{\nabla}x^*(x, \beta)(\beta' - \beta) + \frac{1}{2}\mu_g\|\beta' - \beta\|_2^2 + \mathrm{err}_{\nabla f_0}(x, \beta), \quad (14)$$

*where $\mu_g := nL_0M_1$. Then $g(\beta'; \beta, x)$ majorizes $f_0 \circ x^*$, satisfying Equation (11).*

Note that technically we cannot explicitly compute the value $g(\beta'; x, \beta)$ because it contains the term $f_0(x_\beta)$. However, we can compute the difference $g(\beta'; x, \beta) - g(\beta; x, \beta)$, which is enough to optimize $g$ and to prove descent for the iterates of any given optimization scheme.

## 7.2 Convergence analysis

We now give the convergence result for the Pareto majorization-minimization algorithm. We make use of the sufficient condition provided by Lemma 5, which can be determined using approximate information. And as Algorithm 1 can make use of any black-box optimizer, we state the result in terms of the convergence guarantees of the black-box optimizers.

In particular, the PMM algorithm uses two optimizers: one for the surrogate $g(\,\cdot\,; x, \beta)$ and another for the scalarized objective $f_\beta(\cdot)$. As we aim to achieve $\varepsilon_0$-preference stationarity, we also ask the optimizer for the surrogate $g$ to achieve $O(\varepsilon_0)$-approximate stationarity.

But approximate stationarity with respect to $g$ only transfers to $f_0$ when the surrogate is sufficiently tight, which depends on the performance of the optimizer for $f_\beta(\cdot)$. It turns out that we shall require that it achieves $O(\varepsilon_0^2)$-optimality. This is because when we optimize a positive-definite quadratic over a convex set, finding an $\varepsilon_0$-approximate stationary point $\hat\beta$ means finding an $O(\varepsilon_0^2)$-approximately optimal point (Lemmas 13 and 14):

$$g(\hat\beta; x, \beta) < g(\beta^*; x, \beta) + O(\varepsilon_0^2),$$

where $\beta^*$ minimizes the surrogate. But, the surrogate contains an approximation error $\mathrm{err}_{\nabla f_0}(\beta, x)$. If this error term is larger than $\Omega(\varepsilon_0^2)$, then it is possible for the surrogate to fail to either (i) decide that the current iterate $\beta$ is $\varepsilon_0$-preference stationary or (ii) make progress by finding some $\hat\beta$ that certifiably improves on $f_0$. We preclude this by requiring the optimizer for $f_\beta$ to achieve $O(\varepsilon_0^2)$-optimality.

**Theorem 1** (Convergence of PMM). *Let $F$ and $f_0$ satisfy Assumptions A,B,C. Fix $0 < \varepsilon^{1/2} \leq \varepsilon_0 \leq 1$. Let $\hat{x}_\beta$ and $\hat{\beta}$ be the approximate solutions that are returned by the black-box optimizer for $g(\,\cdot\,; x, \beta)$ and $f_\beta(\cdot)$, defined in Equation (14) and Equation (3), respectively:*

$$\hat{\beta} \leftarrow \widehat{\arg\min_{\beta' \in \Delta^{n-1}}} g(\beta'; x, \beta) \qquad and \qquad \hat{x}_\beta \leftarrow \widehat{\arg\min_{x \in \mathbb{R}^d}} f_\beta(x).$$

*Given constants $c_1, c_2 > 0$, suppose that the black-box optimizer achieves the following guarantees:*

  *1. the approximate minimizer $\hat{\beta}$ is $O(\varepsilon_0)$-approximately stationary:*

$$- \nabla g(\hat{\beta}; x, \beta) v \leq c_1 \cdot \varepsilon_0 \|v\|_2, \qquad \forall v \in T_{\Delta^{n-1}}(\beta).$$

  *2. the approximate minimizer $\hat{x}_\beta$ is an $O(\varepsilon_0^2)$-approximate solution:*

$$\|\nabla f_\beta(\hat{x}_\beta)\| \leq c_2 \cdot \varepsilon.$$

*Let $(x_k, \beta_k)_k$ be the iterates of Algorithm 1. Then, there exist $c_1(f_0, F)$ and $c_2(f_0, F)$ bounded away from zero and some $K$ such that $(f_0 \circ x^*)(\beta_k)$ is monotonically decreasing for $k \in [K]$ and $(x_K, \beta_K)$ is an $(\varepsilon_0, \varepsilon)$-preference stationary point. Furthermore, $K$ is no more than $O(\varepsilon_0^{-2})$:*

$$K \leq \frac{2\mu_g \cdot (f^* - f_*)}{c_1^2 \cdot \varepsilon_0^2},$$

*where $f^* := \max f_0(x)$ and $f_* = \min f_0(x)$ are optimized over the compact set $\mathrm{Pareto}(F)$.*

**Remark 1.** *Algorithm 1 makes calls to sub-routines at each iteration to solve two sub-problems. As the problems are strongly-convex and Lipschitz-smooth, they can be solved using (projected) gradient descent with iteration complexity $O(\log(1/\varepsilon_0))$. And so, taking the computational cost of the sub-problems into account only increases the rate obtained in Theorem 1 by logarithmic factors.*

## 8 Conclusion

In this work, we provide a principled and efficient way to select a decision vector from the Pareto set of a set of objectives $f_1, \ldots, f_n$ given an additional preference function $f_0$. A main contribution of this work is to provide a geometrically-meaningful notion of (approximate) preference stationarity. This is non-trivial due to the non-smoothness and non-convexity of the Pareto set. We also provide a simple algorithm that achieves $\varepsilon_0$-approximate stationarity with iteration complexity of $O(\varepsilon_0^{-2})$.

# 9 Proofs and derivations

| Symbol | Usage |
|---|---|
| $\Delta^{n-1}$ | the $(n-1)$-simplex equipped with the $\ell_1$-metric, see Definition 2 |
| $\Delta^{n-1}(x)$ | the set of $\beta$ satisfying $\nabla f_\beta(x) = 0$, see Equation (8) |
| $\nabla x^*, \widehat{\nabla} x^*$ | derivative of the map $x^*$ and its approximation, see Equations (5) and (10) |
| $\mathrm{err}_{\nabla f_0}(x, \beta)$ | bound on the approximation error of $\nabla(f_0 \circ x^*)$, see Lemma 4 |
| $F, (f_1, \ldots, f_n)$ | the set of objective functions |
| $f_0$ | the preference function |
| $f_\beta(x)$ | the scalarized objective $\sum_i \beta_i f_i(x)$, see Equation (3) |
| $g(\beta'; x, \beta)$ | majorizing surrogate for $f(x_{\beta'})$, see Equation (11) |
| $\kappa$ | condition number $\kappa := L/\mu$ for $\nabla^2 f_i$, see Assumption A |
| $L, L_H, L_0$ | Lipschitz parameters for $\nabla f_i$, $\nabla^2 f_i$, and $\nabla f_0$, see Assumptions A, B, C |
| $M_0, M_1$ | Lipschitz parameters for $x^*$ and $\nabla x^*$, see Lemma 2 |
| $\mu$ | strong convexity parameter for $f_i$, see Assumption A |
| $\mu_g$ | strong convexity parameter $nL_0 M_1$ for the surrogate $g$, see Equation (11) |
| $\mathrm{Pareto}(F)$ | the set of Pareto optimal solutions of $F$, see Definition 1 |
| $r$ | distance between the minimizers of $f_1, \ldots, f_n$, see Assumption A |
| $\mathcal{P}(F)$ | the Pareto manifold, see Definition 3 |
| $R$ | $\mathrm{diam}\big(\mathrm{Pareto}(F)\big) := \sup\big\{ \|x - x'\|_2 : x, x' \in \mathrm{Pareto}(F) \big\}$, see Lemma 1 |
| $x^*(\beta), x_\beta$ | stationary point for $f_\beta$, see Equation (4) |

## 9.1 The Pareto manifold

**Proposition 1** (Characterization of the Pareto manifold)**.** *Define the map* $x^* : \Delta^{n-1} \to \mathrm{Pareto}(F)$*:*

$$x^*(\beta) \equiv x_\beta := \underset{x \in \mathbb{R}^d}{\arg\min}\, f_\beta(x). \tag{4}$$

*Let* $\nabla F(x) \in \mathbb{R}^{n \times d}$ *be the Jacobian. Then, the map* $x^*$ *has derivative:*

$$\nabla x^*(\beta) = -\nabla^2 f_\beta(x_\beta)^{-1} \nabla F(x_\beta)^\top, \tag{5}$$

*so that the map* $\beta \mapsto (x_\beta, \beta)$ *is a diffeomorphism of* $\Delta^{n-1}$ *with the Pareto manifold* $\mathcal{P}(F)$*.*

**Proof of Proposition 1**    The map $x^*$ is well-defined because $f_\beta$ is strictly convex—it is the convex combination of strictly convex objectives, so it has a unique minimizer. Furthermore, because the objectives are smooth, the stationarity condition $\nabla f_\beta(x) = 0$ uniquely holds at $x^*(\beta)$:

$$\nabla f_\beta(x_\beta) = 0.$$

Define the map $\zeta(x, \beta) = \nabla f_\beta(x)$. Then, the Pareto manifold is precisely the zero set $\mathcal{P}(F) = \zeta^{-1}(0)$, and which can be parametrized by simplex $\Delta^{n-1}$ via the map $\beta \mapsto (x_\beta, \beta)$.

In fact, it is a smooth parametrization. To see this, we apply the implicit function theorem (Theorem 2), which states that the map $x^*$ is smooth at $\beta$ when $\nabla_x \zeta(x_\beta, \beta)$ is invertible. Indeed, we have that $\zeta$ is continuously differentiable, with:

$$\nabla_x \zeta(x, \beta) = \sum_{i \in [n]} \beta_i \nabla^2 f_i(x) = \nabla^2 f_\beta(x),$$

$$\nabla_\beta \zeta(x, \beta) = \nabla_\beta \left( \sum_{i \in [n]} \beta_i \nabla f_i(x) \right) = \nabla F(x)^\top.$$

Because $f_\beta$ is strictly convex, it has positive definite Hessian, implying invertibility $\det \nabla_x \zeta(x_\beta, \beta) \neq 0$. Furthermore, Theorem 2 also implies that the derivative of $\nabla x^*$ is given by Equation (5). It follows that the map $\beta \mapsto (x_\beta, \beta)$ is smooth. It also has a smooth inverse. Namely, the projection onto the second component $(x_\beta, \beta) \mapsto \beta$. Thus, $\mathcal{P}(F)$ is diffeomorphic with $\Delta^{n-1}$. ∎

**Theorem 2** (Implicit function theorem, Spivak (2018)). *Let $f : \mathbb{R}^d \times \mathbb{R}^n \to \mathbb{R}^d$ be continuously differentiable on an open set containing $(a, b)$ and let $f(a, b) = 0$. Let $\nabla_u f(u, v)$ be the $d \times d$ matrix:*

$$\left[ \nabla_u f(u, v) \right]_{ij} = \nabla_{u_j} f_i(u, v).$$

*If $\det \nabla_u f(a, b) \neq 0$, there are open sets $U \subset \mathbb{R}^d$ and $V \subset \mathbb{R}^n$ containing $a$ and $b$ respectively with the following property: for each $v \in V$ there is a unique $g(v) \in U$ such that $f(g(v), v) = 0$. Furthermore, the map $g$ is differentiable with derivative given by:*

$$\nabla g(v) = -\left[ \nabla_u f(g(v), v) \right]^{-1} \nabla_v f(g(v), v).$$

## 9.2   Solution concepts to Pareto-constrained optimization

In this section, we elaborate on how the different solution concepts (optimality, stationarity, approximate stationarity) relate to each other for the Pareto-constrained optimization problem:

$$\underset{\beta \in \Delta^{n-1}}{\text{minimize}} \ (f_0 \circ x^*)(\beta). \tag{7}$$

We can call any optimal solution *preference optimal*:

**Definition 9** (Preference optimality). *A decision vector $x \in \text{Pareto}(F)$ is* preference optimal *if:*

$$f_0(x) \leq f_0(x'), \qquad \text{for all } x' \in \text{Pareto}(F).$$

Because preference optimality is a global condition, it is generally computationally infeasible to verify. By considering Equation (7) as a smooth optimization problem over the simplex, we relax the solution concept in the standard way to the first-order stationarity condition in terms of $\beta$:

$$-\nabla(f_0 \circ x^*)(\beta)(\beta' - \beta) \leq 0, \qquad \forall \beta' \in \Delta^{n-1}. \tag{9}$$

Given a stationary point $\beta$, we can push forward this stationary condition to $x^*(\beta)$, which we say is *weakly preference stationary*. We reproduce the definition from before:

**Definition 4** (Preference stationarity). *We say that a point $x \in \mathrm{Pareto}(F)$ is weakly preference stationary if there exists some $\beta \in \Delta^{n-1}(x)$ such that:*[3]

$$- \nabla(f_0 \circ x^*)(\beta)^\top (\beta' - \beta) \leq 0, \qquad \forall \beta' \in \Delta^{n-1}, \tag{9}$$

*where Equation (5) gives $\nabla x^*$. If (9) holds for all $\beta \in \Delta^{n-1}(x)$, then $x$ is preference stationary.*

Finally, to relax the exact stationarity condition to an approximate one, we appeal to the standard notion of an approximate stationary point (Nesterov, 2013). In our setting, we can make use of:

**Definition 10** (Approximate stationary point, Marumo et al. (2023)). *Let $\mathcal{C}$ be a closed and convex set, and let $f : \mathcal{C} \to \mathbb{R}$ be a smooth objective function. A point $\beta \in \mathcal{C}$ is an $\varepsilon$-approximate stationary point of $f$ if for all $\beta' \in \mathcal{C}$, the following holds:*

$$-\nabla f(\beta)^\top (\beta' - \beta) \leq \varepsilon \|\beta' - \beta\|.$$

Specializing Definition 10 to the optimization of $f_0 \circ x^*$ over $\Delta^{n-1}$ yields an approximate stationary condition for $\beta$. And because we are ultimately interested in $x^*(\beta)$, which is the solution of to optimizing $f_\beta$ over $\mathbb{R}^d$, we can also make use of Definition 10 to also define the appropriate approximate stationary condition on $x$. This leads us to Definition 8, which we reproduce here:

**Definition 8** (Approximate preference stationarity). *Let $\varepsilon_0, \varepsilon \geq 0$. A point $(x, \beta) \in \mathbb{R}^d \times \Delta^{n-1}$ is $(\varepsilon_0, \varepsilon)$-preference stationary if:*

$$-\nabla f_0(x_\beta)^\top \nabla x^*(\beta)(\beta' - \beta) \leq \varepsilon_0 \|\beta' - \beta\|_1, \qquad \forall \beta' \in \Delta^{n-1}. \tag{12a}$$
$$\|\nabla f_\beta(x)\|_2 \leq \varepsilon \tag{12b}$$

**Proposition 4** (Geometric meaning of approximate stationarity). *Let $F$ and $f_0$ satisfy Assumptions $A$,$B$,$C$ and let $(\hat{x}, \hat{\beta})$ be $(\varepsilon_0, \varepsilon)$-preference stationary. The following hold:*

a. *if $\|\beta - \hat{\beta}\|_1 \leq s$, then $f_0(x_\beta) - f_0(x_{\hat{\beta}}) \geq -2\varepsilon_0 \|\beta - \hat{\beta}\|_1$, and*

b. *$\|\hat{x} - x_{\hat{\beta}}\|_2 \leq \varepsilon/\mu$,*

*where we let $R$ is defined in Lemma 1 and $s := \frac{2\mu^2 \varepsilon_0}{L_0 L^2 R^2}$.*

**Proof of Proposition 4**

(a) Recall that $x_\beta$ is the minimizer of $f_\beta$, by definition. Because $f_\beta$ is $\mu$-strongly convex, we can bound the distance between $x$ and $x_\beta$ by:

$$\|x - x_\beta\| \leq \frac{1}{\mu} \|\nabla f_\beta(x)\| \leq \frac{\varepsilon}{\mu},$$

where the second inequality follows from condition (12a).

(b) Let $\beta_s := (1 - s)\beta + s\beta'$ parametrize the line connecting $\beta$ and $\beta'$. Let $\gamma : [0, 1] \to \mathrm{Pareto}(F)$ be the path $\gamma(s) := x^*(\beta_s)$, so that:

$$d\gamma(s) = \nabla x^*(\beta_s)(\beta' - \beta) \, ds.$$

---

[3]As $\nabla F(x_\beta)^\top \beta = \nabla f_\beta(x_\beta) = 0$, Equation (9) can be simplified to $-\nabla(f_0 \circ x^*)(\beta)^\top \beta' \leq 0$, for all $\beta' \in \Delta^{n-1}$.

We can now upper bound the difference:

$$f_0(x_\beta) - f_0(x_{\beta'}) = -\int_\gamma \nabla f_0(x_{\beta_s})^\top d\gamma(s)$$

$$= -\int_\gamma \left[\nabla f_0(x_{\beta_s}) - \nabla f_0(x_\beta) + \nabla f_0(x_\beta)\right]^\top d\gamma(s)$$

$$\leq \int_\gamma L_0 \|x_{\beta_s} - x_\beta\| \, |d\gamma(s)| + \int_\gamma \left(-\nabla f_0(x_\beta)^\top d\gamma(s)\right).$$

Let's bound the integrals separately. Since $x_{\beta_s} = \int_0^s d\gamma(s)(\beta' - \beta)$, we have by Lemma 8:

$$\|x_{\beta_s} - x_\beta\| \leq \frac{LR}{\mu} \|\beta - \beta'\|_1 \cdot s.$$

We also have $|d\gamma(s)| \leq \mu^{-1} LR \|\beta - \beta'\|_1$, by Lemma 8. The first integral is bounded by:

$$\int_\gamma L_0 \|x_{\beta_s} - x_\beta\| \, |d\gamma(s)| \leq \int_0^1 \frac{L_0 L^2 R^2}{\mu^2} \|\beta - \beta'\|_1^2 \cdot s \, ds = \frac{1}{2} \frac{L_0 L^2 R^2}{\mu^2} \|\beta - \beta'\|_1^2.$$

For the second integral, first note that condition (12b) implies:

$$-\nabla f_0(x_\beta)^\top d\gamma(s) = -\nabla f_0(x_\beta)^\top \nabla x^*(x_\beta)(\beta' - \beta) \leq \varepsilon_0 \|\beta - \beta'\|_1,$$

yielding the other bound:

$$\int_\gamma \left(-\nabla f_0(x_\beta)^\top d\gamma(s)\right) \leq \int_0^1 \varepsilon_0 \|\beta - \beta'\|_1 \, ds = \varepsilon_0 \|\beta - \beta'\|_1.$$

Putting these two together, we obtain:

$$f_0(x_\beta) - f_0(x_{\beta'}) \leq \frac{1}{2} \frac{L_0 L^2 R^2}{\mu^2} \|\beta - \beta'\|_1^2 + \varepsilon_0 \|\beta - \beta'\|_1.$$

It follows that if we restrict $\|\beta - \beta'\|_1 \leq \frac{2\mu^2 \varepsilon_0}{L_0 L^2 R^2}$, one of the factors of $\|\beta - \beta'\|_1$ in the first term can be absorbed into the constant, proving the result:

$$f_0(x_\beta) \leq f_0(x_{\beta'}) + 2\varepsilon_0 \|\beta - \beta'\|_1.$$

∎

### 9.2.1 Weak preference stationarity and degeneracy

The solution concepts are related by:

$$\begin{array}{ccccccc} \text{preference} & & \text{preference} & & \text{weak preference} & & \text{approximate preference} \\ \text{optimality} & \subset & \text{stationarity} & \subset & \text{stationarity} & \subset & \text{stationarity} \end{array}$$

It is fairly clear that the first and last inequalities are strict. Here, we discuss the inner inequality.

It turns out that a point $x$ can be weakly preference stationary without being preference stationary. However, this can only happen if $x$ is also a point of singularity in Pareto($F$). Geometrically, if we consider Pareto($F$) as the projection of $\mathcal{P}(F)$ onto its first component in $\mathbb{R}^d$, the this means

that multiple points are collapsed onto $x$. Algebraically, this means that the set of gradients $\nabla f_1(x), \ldots, \nabla f_n(x)$ fails to have full (Pareto) rank (Smale, 1973; Hamada et al., 2020).

To elaborate, recall the set:

$$\Delta^{n-1}(x) := \{\beta \in \Delta^{n-1} : \nabla f_\beta(x) = 0\}. \tag{8}$$

Then, $x$ is Pareto stationary if there is some $\beta$ in $\Delta^{n-1}(x)$, so that:

$$\sum_{i \in [n]} \beta_i \nabla f_i(x) = 0,$$

and the rank of this set of gradients is at most $n - 1$. Since $\Delta^{n-1}$ does not contain any collinear vectors, if $\Delta^{n-1}(x)$ contains more than a single point, then the rank of the set of gradients must be strictly less than $n - 1$. This leads us to the definition:

**Definition 11** (Pareto genericity). *Let $\{v_1, \ldots, v_n\} \subset \mathbb{R}^d$. This set is Pareto generic if:*

$$\beta_1 v_1 + \cdots + \beta_n v_n = 0, \qquad \text{for some } \beta \in \Delta^{n-1},$$

*and the non-degeneracy condition holds:* $\mathrm{rank}(v_1, \ldots, v_n) = n - 1$.

If $\nabla F(x)$ is Pareto generic, then $\Delta^{n-1}(x)$ contains a unique $\beta$, so we immediately have:

**Proposition 6** (Generic and weak implies strong preference stationarity). *If $\nabla F(x)$ is Pareto generic and $x$ is weakly preference stationary, then $x$ is preference stationary.*

However, when the gradients $\nabla F(x)$ are not Pareto generic, then weak preference stationarity can be strictly weaker. Let $(x, \beta)$ where $\beta \in \Delta^{n-1}(x)$ be weakly preference stationary, so that:

$$-\nabla f_0(x)^\top \underbrace{\left(-\nabla^2 f_\beta(x)^{-1} \nabla F(x)^\top (\beta' - \beta)\right)}_{\nabla x^*(\beta)(\beta' - \beta)} \leq 0, \qquad \forall \beta' \in \Delta^{n-1}.$$

We can simplify this by using the fact that $\nabla f_\beta(x) = \nabla F(x)^\top \beta = 0$. Then, one way for the stationary condition to be fulfilled is for the underlined term to be normal to $\nabla f_0(x)$:

$$-\nabla^2 f_\beta(x)^{-1} \nabla F(x)^\top \beta' \in \mathrm{span}(\nabla f_0(x))^\perp, \qquad \forall \beta' \in \Delta^{n-1}.$$

This statement has the following geometric interpretation. These vectors are contained in the Clarke tangent cone of $\mathrm{Pareto}(F)$ at $x$. If these are the only vectors in the tangent cone, then this above condition states that $-\nabla f_0(x)$ is contained in the normal cone of $\mathrm{Pareto}(F)$ at $x$.

But, in general, the tangent cone contains the union of subspaces:

$$\bigcup_{\beta \in \Delta^{n-1}(x)} \left\{ -\nabla^2 f_\beta(x)^{-1} \nabla F(x)^\top \beta' : \beta' \in \Delta^{n-1} \right\}.$$

And so, when $\Delta^{n-1}(x)$ does not contain a unique vector, the tangent cone can contain more vectors. By selecting different $\beta$'s, we recover different slices of the tangent cone. This also means that even if the above normality condition holds for one $\beta$, it may fail to hold for a different $\tilde{\beta} \in \Delta^{n-1}(x)$. In this case, $(x, \beta)$ is weakly preference stationary while $(x, \tilde{\beta})$ may not be.

17

### 9.2.2 Insufficiency of first-order information

**Proposition 3** (Necessary first-order conditions are trivial). *Suppose that* Stationary *is necessary. Then, it is trivial in the following sense: for any preference generic set of* $v_0, \ldots, v_n \in \mathbb{R}^d$,

$$\text{Stationary}(v_0, \ldots, v_n) = \text{true}.$$

**Proof of Proposition 3** It suffices to show that there exist $f_0$, $F$, and $x^\star$ such that $x^\star$ is preference optimal and for $i = 0, \ldots, n$:

$$v_i = \nabla f_i(x^\star). \tag{15}$$

And since $x^\star$ is preference optimal, any necessary stationary condition must accept:

$$\text{Stationary}(v_0, \ldots, v_n) = \text{true}.$$

Without loss of generality, let $x^\star = 0$ by an affine transformation. To construct $f_0$ and $F$, we can simply consider a family of positive-definite quadratics:

- Let the preference function $f_0$ be:

$$f_0(x) = \frac{1}{2}\|x + v_0\|^2.$$

  Notice that $\nabla f_0(x^*) = v_0$.

- Let the objectives $f_1, \ldots, f_n$ share the same Hessian:

$$f_i(x) = \frac{1}{2}\|A(x - z_i)\|^2,$$

  where $A \in \mathbb{R}^{d \times d}$ is full rank and $z_i \in \mathbb{R}^d$. Let $H = A^\top A$ for short.

We show that we can set $A$ and the $z_i$'s so that $x^\star$ is preference optimal while Equation (15) holds.

By Lemma 6, the Pareto set is the convex hull $\mathcal{C} := \text{conv}(z_1, \ldots, z_n)$. Notice that the choice of $H$ and $v_i$'s determines the $z_i$'s, since we require $\nabla f_i(x^*) = v_i$, which expands to:

$$z_i = -H^{-1}v_i, \qquad \forall i \in [n].$$

From convex optimization, $x^\star = 0$ is preference optimal if (i) $x^\star \in \mathcal{C}$ and (ii) $\mathcal{C}$ is normal to $\nabla f_0$. Indeed, these two conditions can be fulfilled:

(i) Because $v_1, \ldots, v_n$ is assumed to be Pareto generic, zero is a convex combination of the $v_i$'s. As the $z_i$'s are related to the $v_i$'s by a linear transformation, this also implies that zero is a convex combination of the $z_i$'s (with the same set of convex weights).

(ii) We need to show that the subspace $\text{span}(v_1, \ldots, v_n)$ can be mapped into $\text{span}(v_0)^\perp$ by the map $v \mapsto -H^{-1}v$ where $H$ is positive definite. Lemma 7 shows that such a map $H$ exists as long as $v_0 \notin \text{span}(v_1, \ldots, v_n)$, which is assumed from preference genericity.

Thus, there exists $f_0$ and $F$ that is preference optimal at $x^\star$ with matching first-order information. A necessary stationary condition must therefore be accepted. ∎

**Remark 2.** *Suppose that* Stationary *is not necessary, but that we can design some optimization method that provably converges to a stationary point in* $\{x : \text{Stationary}(x) = \text{true}\}$. *Then, this also means that there are settings in which the method provably avoids preference optimal points.*

In the remainder of this section, we prove Lemma 6 and Lemma 7 used above.

**Lemma 6.** *Let $f_1, \ldots, f_n : \mathbb{R}^d \to \mathbb{R}$ be positive-definite quadratics with a shared Hessian:*

$$f_i(x) = \frac{1}{2}\|A(x - z_i)\|^2,$$

*where $A \in \mathbb{R}^{d \times d}$ is full rank and $z_i \in \mathbb{R}^d$. Then, the Pareto set is the convex hull:*

$$\text{Pareto}(f_1, \ldots, f_n) = \text{conv}(z_1, \ldots, z_n).$$

*Proof.* As the objectives $f_1, \ldots, f_n$ are strongly convex, optimality is equivalent to stationarity. Thus, $x \in \text{Pareto}(f_1, \ldots, f_n)$ if and only if there exists some $\beta \in \Delta^{n-1}$ such that:

$$0 = \sum_{i \in [n]} \beta_i \nabla f_i(x),$$

which, when expanded, states that:

$$(A^\top A)x = (A^\top A)\sum_{i \in [n]} \beta_i z_i.$$

But as $A$ is invertible, this is equivalent to:

$$x = \sum_{i \in [n]} \beta_i z_i,$$

which is to say that $x \in \text{conv}(z_1, \ldots, z_n)$. $\square$

**Lemma 7.** *Let $U$ and $V$ be linear subspaces of $\mathbb{R}^d$ such that $U \cap V^\perp = \{0\}$. Then, there exists some positive definite map $H : \mathbb{R}^d \to \mathbb{R}^d$ such that $H(U) \subset V$.*

*Proof.* If $S \subset \mathbb{R}^d$ is a subspace, let $\Pi_S : \mathbb{R}^d \to \mathbb{R}^d$ be the projection onto $S$. Define the map:

$$H := \Pi_V + \Pi_{V^\perp}\Pi_{U^\perp}.$$

Then $H$ satisfies the following:

- $H$ is positive definite. To see this, let $0 \neq x \in \mathbb{R}^d$ have decomposition $x = x_1 + x_2$, where $x_1 \in U$ and $x_2 \in U^\perp$. Then:

$$x^\top Hx = \underbrace{x_1^\top \Pi_V x_1 + 2x_1 \Pi_V x_2 + x_2 \Pi_V x_2}_{x^\top \Pi_V x} + \underbrace{x_1^\top \Pi_{V^\perp} x_2 + x_2^\top \Pi_{V^\perp} x_2}_{x^\top \Pi_{V^\perp}\Pi_{U^\perp} x}$$

$$= \|\Pi_V x_1\|^2 + \underbrace{x_1^\top x_2}_{0} + x_1 \Pi_V x_2 + \|x_2\|^2 \geq \frac{1}{2}\|\Pi_V x_1 + x_2\|^2 > 0,$$

  where the last inequality is strict because $x \neq 0$ and $U \cap V^\perp = \{0\}$.

- $H(U) \subset V$. If $x \in U$, then by definition $\Pi_{U^\perp} x = 0$ so that $Hx = \Pi_V x \in V$.

$\square$

### 9.2.3 An example of a first-order stationarity condition avoiding optimality

In this section, we discuss the first-order stationarity condition of Ye and Liu (2022), defined to as stationarity with respect to their optimization dynamics, *Pareto navigating gradient descent* (PNG). We show that it fails to be a necessary condition for preference optimality.

Despite that, their condition and dynamics have appealing properties since (i) they do not require second-order information, which is computationally more expensive, and (ii) their dynamics largely satisfies what they call the *Pareto improvement property*, which ensures that each objective enjoys monotonic improvement during optimization:

$$\frac{d}{dt} f_i(x_t) \leq 0, \qquad \text{for all } i \in [n].$$

As the goal of Pareto improvement can be at odds with preference optimality, this leads to an open question: when and how should we balance Pareto improvement with preference optimality?

**Definition 12** (PNG stationarity, Ye and Liu (2022)). *Let $c > 0$. Define the PNG vector $v_c(x)$:*

$$v_c(x) := \arg\min_{v \in \mathbb{R}^d} \frac{1}{2} \|\nabla f_0(x) - v\|^2$$
$$\text{s.t.} \quad \nabla f_i(x)^\top v \geq c, \qquad \text{for all } i \in [n].$$

*Let $\varepsilon > 0$. A vector $x \in \mathbb{R}^d$ is $(c, \varepsilon)$-PNG stationary if $v_c(x) = \lambda \nabla f_0(x)$ for some $\lambda \leq 0$ and:*

$$\min_{\beta \in \Delta^{n-1}} \|\nabla f_\beta(x)\| = \varepsilon.$$

In the following example, we consider a two-dimensional example with two objectives. Let the standard basis be denoted $\mathbf{e}_1, \mathbf{e}_2 \in \mathbb{R}^2$, and let the objective functions $f_1, f_2 : \mathbb{R}^2 \to \mathbb{R}$ be defined:

$$f_1(x) = \frac{1}{2} \|A(x + \mathbf{e}_1)\|^2 \qquad \text{and} \qquad f_2(x) = \frac{1}{2} \|A(x - \mathbf{e}_1)\|^2, \tag{16}$$

where $A \in \mathbb{R}^{2 \times 2}$ is full-rank. Lemma 6 shows that the Pareto set of the objectives $\text{Pareto}(f_1, f_2)$ is the line segment from $-\mathbf{e}_1$ to $\mathbf{e}_2$. That is, the Pareto set is invariant under changes of $A$. However, the PNG stationarity condition is not, since the constraint set changes with $A$:

$$\{v : (x + \mathbf{e}_1)^\top Hv \geq c\} \cap \{v : (x - \mathbf{e}_1)^\top Hv \geq c\},$$

where $H = A^\top A$. Due to this discrepancy, PNG stationary points can fail to be preference optimal.

**Example 1.** *Let the preference function be: $f_0(x) = \frac{1}{2}\|x - \mathbf{e}_2\|^2$, and let the objectives $f_1, f_2$ be defined as in the above Equation (16) with:*

$$H = A^\top A = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}. \tag{17}$$

*Then, the unique preference optimal point is the origin 0. However, the $(c, \varepsilon)$-PNG stationary point is bounded away from 0. It even converges to $\mathbf{e}_1$ as the error tolerance $\varepsilon$ goes to zero.*

*Proof.* Consider the PNG vector $v_c(x)$ when $x$ is in the region:

$$\mathcal{C} = \left\{ x \in \mathbb{R}^2 : \nabla f_0(x)^\top \nabla f_i(x) < 0, \text{ for } i = 1, 2 \right\} \cap \left\{ \mathbf{e}_2^\top x > 0 \right\}.$$

Here, both constraints $\nabla f_i(x)^\top v \geq c$ are active in the constrained optimization problem that defines the PNG vector; and so, $v_c(x)$ is the vertex point of the constraint set, satisfying:

$$\nabla f_1(x)^\top v_c(x) = \nabla f_2(x)^\top v_c(x) = c.$$

Expanding out the gradients, we obtain:

$$(x + \mathbf{e}_1)^\top H v_c(x) = c \qquad \text{and} \qquad (x - \mathbf{e}_1)^\top H v_c(x) = c.$$

This implies that $\mathbf{e}_1^\top H v_c(x) = 0$. Now suppose that $x_{\mathrm{PNG}} \in \mathcal{C}$ is PNG stationary. Then, by definition, it must satisfy $\nabla f_0(x_{\mathrm{PNG}}) \in \mathrm{span}(v_c(x_{\mathrm{PNG}}))$, so it has the form:

$$x_{\mathrm{PNG}} = \mathbf{e}_2 + \lambda u, \qquad \text{where } \mathbf{e}_1^\top H u = 0.$$

Whenever the standard basis vectors are not eigenvectors of $H$, the line $\mathbf{e}_2 + \lambda u$ intersects $\mathrm{Pareto}(f_1, f_2)$ away from 0. In this example, we let $A$ satisfy $H = A^\top A$ where $H$ is given by Equation (17).

Then, the line $\mathbf{e}_2 + \lambda u$ runs through $\mathbf{e}_1$ and $\mathbf{e}_2$. We can verify that $\mathcal{C}$ contains all points on this line between its two endpoints. When $x = \mathbf{e}_2 + \lambda(\mathbf{e}_1 - \mathbf{e}_2)$ and $\lambda \in (0, 1)$, we have:

$$\begin{aligned}
\nabla f_0(x)^\top \nabla f_1(x) &= (x - \mathbf{e}_2)^\top H(x + \mathbf{e}_1) \\
&= \lambda(\mathbf{e}_1 - \mathbf{e}_2)^\top H\big((1 + \lambda)\mathbf{e}_1 + (1 - \lambda)\mathbf{e}_2\big) = -\lambda(1 - \lambda),
\end{aligned}$$

and similarly, we have:

$$\begin{aligned}
\nabla f_0(x)^\top \nabla f_2(x) &= (x - \mathbf{e}_2)^\top H(x - \mathbf{e}_1) \\
&= \lambda(\mathbf{e}_1 - \mathbf{e}_2)^\top H\big((\lambda - 1)(\mathbf{e}_1 - \mathbf{e}_2)\big) = -\lambda(1 - \lambda).
\end{aligned}$$

This implies that for all $c > 0$ and $\varepsilon > 0$, the $(c, \varepsilon)$-PNG stationary point is bounded away from 0, converging to $\mathbf{e}_1$ as $\varepsilon$ goes to zero. $\qquad\square$

## 9.3   Implications of smoothness assumptions

**Lemma 1** (Size of Pareto set). *Suppose $F$ satisfies Assumption $A$. Then $R \leq \sqrt{\kappa}r$, where:*

$$R := \mathrm{diam}\big(\mathrm{Pareto}(F)\big) \equiv \sup\left\{ \|x - x'\|_2 : x, x' \in \mathrm{Pareto}(F) \right\}.$$

**Proof of Lemma 1**   Because each $f_i$ is $\mu$-strongly convex and $L$-Lipschitz smooth, so too is the convex combination $f_\beta$. This implies the upper and lower bounds:

$$\frac{1}{2}\mu \sum_{i \in [n]} \beta_i \|x - x_i\|_2^2 \leq f_\beta(x) \leq \frac{1}{2}L \sum_{i \in [n]} \beta_i \|x - x_i\|_2^2.$$

It follows that the minimizer of $f_\beta$ is bounded:

$$f_\beta(x_\beta) \leq \frac{1}{2}Lr^2.$$

On the other hand, if a point $\|x - x_i\| > 2s$ for some $i \in [n]$, then by reverse triangle inequality, $\|x - x_j\| > s$ for all $j \in [n]$. This implies that:

$$\|x - x_i\| > 2s \qquad \implies \qquad f_\beta(x) > \frac{1}{2}\mu s^2.$$

It follows that if $\|x - x_i\| > 2\sqrt{L/\mu}$ for some $i$, then $x$ is not a Pareto optimal point. ∎

**Lemma 2** (Smoothness of $x^*$). *Suppose $F$ satisfies Assumptions A,B. Then, $x^* : \Delta^{n-1} \to \mathbb{R}^d$ is $M_0$-Lipschitz continuous and has $M_1$-Lipschitz continuous gradients, where:*

$$M_0 := \kappa R \qquad and \qquad M_1 := 2\kappa^2 R \left(1 + \frac{L_H R}{\mu}\right).$$

**Proof of Lemma 2** That $x^*$ is Lipschitz continuous with Lipschitz continuous gradients follows from the following two lemmas:

**Lemma 8.** *Let $F \equiv (f_1, \ldots, f_n)$ be a set of twice-differentiable objective functions and let $f_0$ be a smooth preference function. Suppose the objectives are $L$-Lipschitz smooth and $\mu$-strongly convex:*

$$\mu\mathbf{I} \preceq \nabla^2 f_i \preceq L\mathbf{I}.$$

*Let $R := \mathrm{diam}\big(\mathrm{Pareto}(F)\big)$. Then, the map $x^* : (\Delta^{n-1}, \ell_1) \to (\mathbb{R}^d, \ell_2)$ is $LR/\mu$-Lipschitz.*

*Proof.* Recall from Equation (5) that $\nabla x^*(\beta) = -\nabla^2 f_\beta(x_\beta)^{-1}\nabla F(x_\beta)^\top$. The following holds:

$$\begin{aligned}
\|\nabla x^*(\beta)\|_{1,2} &\overset{(i)}{\leq} \left\|\nabla^2 f_\beta(x_\beta)^{-1}\right\|_2 \cdot \left\|\nabla F(x_\beta)^\top\right\|_{1,2} \\
&\overset{(ii)}{\leq} \frac{1}{\mu} \cdot LR,
\end{aligned}$$

where (i) is a property of the $\|\cdot\|_{1,2}$-norm, (ii) uses $\mu\mathbf{I} \preceq \nabla^2 f_\beta(x_\beta)$ and Lemma 10. □

**Lemma 9.** *Let $\beta, \beta' \in \Delta^{n-1}$. Then,*

$$\left\|\nabla x^*(\beta) - \nabla x^*(\beta')\right\|_{1,2} \leq \frac{2L^2 R}{\mu^2}\left(1 + \frac{L_H R}{\mu}\right) \cdot \|\beta - \beta'\|_1.$$

*Proof.* By definition, we have:

$$\left\|\nabla x^*(\beta) - \nabla x^*(\beta')\right\|_{1,2} = \left\| -\nabla^2 f_\beta(x_\beta)^{-1}\nabla F(x_\beta)^\top + \nabla^2 f_{\beta'}(x_{\beta'})^{-1}\nabla F(x_{\beta'})^\top\right\|_{1,2}.$$

We can add and subtract $\nabla^2 f_\beta(x_\beta)^{-1}\nabla F(x_{\beta'})^\top$ inside the norm on the right-hand side (RHS):

$$(\text{RHS}) = \left\| -\nabla^2 f_\beta(x_\beta)^{-1} \cdot \left[\nabla F(x_\beta) - \nabla F(x_{\beta'})\right]^\top + \left[\nabla^2 f_\beta(x_\beta)^{-1} - \nabla^2 f_{\beta'}(x_{\beta'})^{-1}\right] \cdot \nabla F(x_{\beta'})^\top\right\|_{1,2}.$$

We can bound the two terms in the norm separately. For the first:

$$\left\| -\nabla^2 f_\beta(x_\beta)^{-1} \cdot \left[\nabla F(x_\beta) - \nabla F(x_{\beta'})\right]^\top\right\|_{1,2} \overset{(i)}{\leq} \frac{L}{\mu} \cdot \|x_\beta - x_{\beta'}\| \overset{(ii)}{\leq} \frac{L^2 R}{\mu^2}\|\beta - \beta'\|_1,$$

where (i) follows the same argument as Lemma 3, and (ii) applies Lemma 8. For the second term, we can add and subtract $\nabla^2 f_{\beta'}(x_\beta)^{-1}\nabla F(x_{\beta'})^\top$ to obtain:

$$\left\|\left[\nabla^2 f_\beta(x_\beta)^{-1} - \nabla^2 f_{\beta'}(x_{\beta'})^{-1}\right]\cdot\nabla F(x_{\beta'})^\top\right\|_{1,2}$$
$$= \left\|\left[\nabla^2 f_\beta(x_\beta)^{-1} - \nabla^2 f_{\beta'}(x_\beta)^{-1} + \nabla^2 f_{\beta'}(x_\beta)^{-1} - \nabla^2 f_{\beta'}(x_{\beta'})^{-1}\right]\cdot\nabla F(x_{\beta'})^\top\right\|_{1,2}$$
$$\leq \left(\frac{L}{\mu^2}\|\beta - \beta'\|_1 + \frac{L_H}{\mu^2}\frac{LR}{\mu}\|\beta - \beta'\|_1\right)\cdot LR.$$

where $\nabla^2 f_\beta(x)^{-1} - \nabla^2 f_{\beta'}(x)^{-1}$ is bounded by Lemma 12; $\nabla^2 f_\beta(x)^{-1} - \nabla^2 f_\beta(x')^{-1}$ is bounded by Lemma 11 and Lemma 8; and $\|\nabla F(x_{\beta'})^\top\|_{1,2}$ is bounded by Lemma 10. $\square$

The result follows by substituting in the definitions of $M_0$ and $M_1$. ∎

**Lemma 3** (Approximability of $\nabla x^*$). *If $F$ satisfies Assumptions A,B. Then:*

$$\|\nabla x^*(\beta) - \widehat{\nabla}x^*(x,\beta)\|_{1,2} \leq \frac{1}{\mu}\frac{M_1}{2M_0}\|\nabla f_\beta(x)\|_2.$$

*Proof.* Recall that $x_\beta := x^*(\beta)$. Then, by definition, we have:

$$\left\|\widehat{\nabla}x^*(x,\beta) - \nabla x^*(\beta)\right\|_{1,2} = \left\|-\nabla^2 f_\beta(x)^{-1}\nabla F(x)^\top + \nabla^2 f_\beta(x_\beta)^{-1}\nabla F(x_\beta)^\top\right\|_{1,2}.$$

We can add and subtract $\nabla^2 f_\beta(x)^{-1}\nabla F(x_\beta)^\top$ inside the norm on the right-hand side (RHS) to get:

$$(\text{RHS}) = \left\|-\nabla^2 f_\beta(x)^{-1}\cdot\left[\nabla F(x) - \nabla F(x_\beta)\right]^\top + \left[\nabla^2 f_\beta(x)^{-1} - \nabla^2 f_\beta(x_\beta)^{-1}\right]\cdot\nabla F(x_\beta)^\top\right\|_{1,2}$$
$$\overset{(i)}{\leq} \frac{L}{\mu}\cdot\|x - x_\beta\| + \frac{L_H}{\mu^2}\|x - x_\beta\|\cdot LR$$
$$\overset{(ii)}{\leq} \frac{L}{\mu^2}\left(1 + \frac{L_H R}{\mu}\right)\cdot\|\nabla f_\beta(x)\|,$$

where (i) the first blue term uses $\mu\mathbf{I} \preceq \nabla^2 f_\beta$ and the $L$-Lipschitz smoothness of the objectives, while the bracket orange term follows from Lemma 11 and the final purple term follows from Lemma 10, and (ii) uses the $\mu$-strong convexity of $f_\beta$. $\square$

**Lemma 4** (Approximability of $\nabla(f_0 \circ x^*)$). *If $F$ and $f_0$ satisfy Assumptions A,B,C. Then:*

$$\left\|\nabla(f_0 \circ x^*)(\beta)^\top - \nabla f_0(x)^\top\widehat{\nabla}x^*(x,\beta)\right\|_{1,2} \leq \frac{1}{\mu}\left(\frac{M_1}{2M_0}\|\nabla f_0(x)\|_2 + L_0 M_0\right)\|\nabla f_\beta(x)\|_2.$$

*We denote the right-hand side by* $\text{err}_{\nabla f_0}(x,\beta)$.

**Proof of Lemma 4**   Add and subtract $\nabla f_0(x)^\top\nabla x^*(\beta)$ within the norm on the right-hand side:

$$(\text{RHS}) = \left\|\left(\nabla f_0(x_\beta)^\top - \nabla f_0(x)\right)^\top\nabla x^*(\beta) + \nabla f_0(x)^\top\left(\nabla x^*(\beta) - \widehat{\nabla}x^*(x,\beta)\right)\right\|_{1,2}$$
$$\leq L_0 M_0\|x_\beta - x\| + \|\nabla f_0(x)\|\cdot\frac{1}{\mu}\frac{M_1}{2M_0}\|\nabla f_\beta(x)\|_2,$$

where we use the fact that $f_0$ is $L_0$-Lipschitz smooth by Assumption C, that $x^*$ is $M_0$-Lipschitz continuous by Lemma 2, and that $\|\nabla x^*(\beta) - \widehat{\nabla} x^*(\beta)\|_{1,2}$ is bounded by Lemma 3. The result follows from upper bounding $\|x_\beta - x\|$ by $\mu$-strong convexity of $f_\beta$:

$$\|x_\beta - x\| \leq \frac{1}{\mu}\|\nabla f_\beta(x)\|.$$

∎

**Lemma 5** (Verifiability of approximate stationarity). *Let $F$ and $f_0$ satisfy Assumptions A,B,C. Then $(\hat{x}, \hat{\beta})$ is $(\varepsilon_0, \varepsilon)$-preference stationary if $\|\nabla f_{\hat{\beta}}(\hat{x})\|_2 \leq \varepsilon$, and for some $x \in \mathbb{R}^d$ and $\alpha \in (0, 1)$,*

1. *an $\alpha \cdot \varepsilon_0$-approximate stationary condition holds:*
$$-\nabla f_0(x)^\top \widehat{\nabla} x^*(x, \hat{\beta})(\beta' - \hat{\beta}) \leq \alpha \cdot \varepsilon_0 \|\beta' - \hat{\beta}\|_1, \tag{13}$$

2. *an error bound holds:*
$$\mathrm{err}_{\nabla f_0}(\hat{\beta}, x) \leq (1 - \alpha) \cdot \varepsilon_0.$$

**Proof of Lemma 5** For $(\varepsilon, \varepsilon_0)$-preference stationarity, we require that $\|\nabla f_{\hat{\beta}}(\hat{x})\|_2 \leq \varepsilon$ and:

$$\nabla f_0(x_{\hat{\beta}})^\top \nabla x^*(\hat{\beta})(\beta' - \hat{\beta}) + \varepsilon_0 \cdot \|\beta' - \hat{\beta}\|_1 \geq 0.$$

Then by Lemma 4, the left-hand side is lower bounded:

$$\nabla f_0(x)^\top \widehat{\nabla} x^*(x, \hat{\beta})(\beta' - \hat{\beta}) - \mathrm{err}_{\nabla f_0}(\hat{\beta}, x) \cdot \|\beta' - \hat{\beta}\|_1 + \varepsilon_0 \|\beta' - \hat{\beta}\|_1,$$

$$= \underbrace{\nabla f_0(x)^\top \widehat{\nabla} x^*(x, \hat{\beta})(\beta' - \hat{\beta}) + \alpha \cdot \varepsilon_0 \|\beta' - \hat{\beta}\|_1}_{\geq 0} + \underbrace{(1 - \alpha) \cdot \varepsilon_0 \|\beta' - \hat{\beta}\|_1 - \mathrm{err}_{\nabla f_0}(\hat{\beta}, x) \cdot \|\beta' - \hat{\beta}\|_1}_{\geq 0},$$

for $\alpha \in (0, 1)$. The two terms are lower bounded by zero by conditions (1) and (2), respectively. ∎

## 9.4 Convergence for Pareto majorization-minimization

**Theorem 1** (Convergence of PMM). *Let $F$ and $f_0$ satisfy Assumptions A,B,C. Fix $0 < \varepsilon^{1/2} \leq \varepsilon_0 \leq 1$. Let $\hat{x}_\beta$ and $\hat{\beta}$ be the approximate solutions that are returned by the black-box optimizer for $g(\cdot; x, \beta)$ and $f_\beta(\cdot)$, defined in Equation (14) and Equation (3), respectively:*

$$\hat{\beta} \leftarrow \widehat{\arg\min_{\beta' \in \Delta^{n-1}}} g(\beta'; x, \beta) \qquad and \qquad \hat{x}_\beta \leftarrow \widehat{\arg\min_{x \in \mathbb{R}^d}} f_\beta(x).$$

*Given constants $c_1, c_2 > 0$, suppose that the black-box optimizer achieves the following guarantees:*

1. *the approximate minimizer $\hat{\beta}$ is $O(\varepsilon_0)$-approximately stationary:*
$$-\nabla g(\hat{\beta}; x, \beta)v \leq c_1 \cdot \varepsilon_0 \|v\|_2, \qquad \forall v \in T_{\Delta^{n-1}}(\beta).$$

2. *the approximate minimizer $\hat{x}_\beta$ is an $O(\varepsilon_0^2)$-approximate solution:*
$$\|\nabla f_\beta(\hat{x}_\beta)\| \leq c_2 \cdot \varepsilon.$$

*Let $(x_k, \beta_k)_k$ be the iterates of Algorithm 1. Then, there exist $c_1(f_0, F)$ and $c_2(f_0, F)$ bounded away from zero and some $K$ such that $(f_0 \circ x^*)(\beta_k)$ is monotonically decreasing for $k \in [K]$ and $(x_K, \beta_K)$ is an $(\varepsilon_0, \varepsilon)$-preference stationary point. Furthermore, $K$ is no more than $O(\varepsilon_0^{-2})$:*

$$K \leq \frac{2\mu_g \cdot \left(f^* - f_*\right)}{c_1^2 \cdot \varepsilon_0^2},$$

*where $f^* := \max f_0(x)$ and $f_* = \min f_0(x)$ are optimized over the compact set $\mathrm{Pareto}(F)$.*

**Proof of Theorem 1** Fix $k > 1$. For short, we let:

$$(x, \beta) \equiv (x_{k-1}, \beta_{k-1}) \qquad \text{and} \qquad (\hat{x}, \hat{\beta}) \equiv (x_k, \beta_k).$$

*Claim.* At each iteration, either (i) the preference improves by at least a constant:

$$f_0(x_{\hat{\beta}}) - f_0(x_\beta) \leq -\frac{1}{2} \frac{c_1}{\mu_g} \cdot \varepsilon_0^2,$$

or (ii) the point $(\hat{x}, \hat{\beta})$ is $(\varepsilon_0, \varepsilon)$-preference stationary.

Assuming the claim holds, the theorem immediately follows: if the algorithm in $K$ steps has not found an $(\varepsilon_0, \varepsilon)$-preference stationary point, then the value $f_0(x_{\beta_k})$ must decrease every iteration by a constant. But because $f_0 \circ x^*$ is lower bounded over $\Delta^{n-1}$ by $f^*$, this can happen at most:

$$\frac{2\mu_g \cdot (f^* - f_*)}{c_1^2 \cdot \varepsilon_0^2} \quad \text{times.}$$

*Proof of the claim.* Let $\beta^* := \arg\min_{\beta' \in \Delta^{n-1}} g(\beta'; x, \beta)$. Lemma 13 shows that an approximate stationary point $\hat{\beta}$ of a strongly convex function is close to the exact stationary point $\beta^*$:

$$\|\hat{\beta} - \beta^*\|_2 \leq \frac{c_1 \varepsilon_0}{\mu_g} =: \delta, \tag{18}$$

where we let $\delta$ denote this constant for short.

We can analyze $\hat{\beta}$ through $\beta^*$. There are two cases, leading to either (1) $O(\varepsilon_0)$-preference stationarity or (2) $O(\varepsilon_0^2)$-constant descent. The two cases depend on the suboptimality of $\beta$.

Case 1: $\|\beta^* - \beta\|_2 < 2\delta$. Here, $\beta$ is fairly close to the optimum $\beta^*$ of the surrogate. We show that the approximate stationarity of $\hat{\beta}$ with respect to the surrogate implies approximate preference stationarity. We do so via Lemma 5, which states that $(\hat{x}, \hat{\beta})$ is $(\varepsilon_0, \varepsilon)$-preference stationary provided:

$$\|\nabla f_{\hat{\beta}}(\hat{x})\|_2 \leq \varepsilon \tag{19}$$

$$-\nabla f_0(x)^\top \widehat{\nabla} x^*(x, \hat{\beta})(\beta' - \hat{\beta}) \leq \frac{1}{2}\varepsilon_0 \|\beta' - \hat{\beta}\|_1, \qquad \forall \beta' \in \Delta^{n-1} \tag{20}$$

$$\text{err}_{\nabla f_0}(x, \hat{\beta}) \leq \frac{1}{2}\varepsilon_0 \tag{21}$$

While Equation (19) is immediate from our choice of $c_2$, defined in the last section of the proof, the others do not follow automatically from approximate stationarity with respect to the surrogate: the surrogate is derived from local information at $(x, \beta)$, while we would like guarantees at $(x, \hat{\beta})$. But because $\beta^*$ is close to both $\beta$ and $\hat{\beta}$, we can control all of these. By triangle inequality:

$$\|\beta - \hat{\beta}\|_2 \leq \|\beta - \beta^*\|_2 + \|\beta^* - \hat{\beta}\|_2 < 3\delta, \tag{22}$$

combining Equation (18) and the assumption that $\|\beta^* - \beta\|_2 < 2\delta$.

We now show Equation (20). We have for all $\beta' \in \Delta^{n-1}$,

$$-\nabla f_0(x)^\top \widehat{\nabla} x^*(x, \hat{\beta})(\beta' - \hat{\beta})$$

$$\overset{(i)}{\leq} -\nabla f_0(x)^\top \widehat{\nabla} x^*(x, \beta)(\beta' - \hat{\beta}) + \|\nabla f_0(x)^\top (\widehat{\nabla} x^*(x, \hat{\beta}) - \widehat{\nabla} x^*(x, \beta))\|_\infty \cdot \|\beta' - \hat{\beta}\|_1$$

$$\overset{(ii)}{\leq} c_1 \varepsilon_0 \cdot \|\beta' - \hat{\beta}\|_1 + \|\nabla f_0(x)\|_2 \cdot \frac{L}{\mu^2} \|\beta - \hat{\beta}\|_2 \cdot \|\beta' - \hat{\beta}\|_1$$

$$\overset{(iii)}{\leq} \frac{1}{2} \cdot 2 \left( c_1 \varepsilon_0 + \frac{L\|\nabla f_0(x)\|_2}{\mu^2} \cdot 3\delta \right) \cdot \|\beta' - \hat{\beta}\|_1 \tag{23}$$

$$\overset{(iv)}{\leq} \frac{1}{2} \varepsilon_0 \cdot \|\beta' - \hat{\beta}\|_1,$$

where (i) adds and subtracts $\nabla f_0(x)^\top \widehat{\nabla} x^*(x, \beta)(\beta' - \hat{\beta})$ and applies Hölder's inequality, (ii) substitutes in Condition 1 for the first term and bounds the second via Lemma 12, and (iii) bounds $\|\beta - \hat{\beta}\|_2$ using Equation (22), and (iv) applies the definition of $c_1$, set in the last section of the proof.

To show Equation (21), we have:

$$\text{err}_{\nabla f_0}(x, \hat{\beta}) \overset{(i)}{=} \text{err}_{\nabla f_0}(x, \beta) + \frac{1}{\mu} \left( \frac{M_1}{2M_0} \|\nabla f_0(x)\|_2 + L_0 M_0 \right) \left( \|\nabla f_{\hat{\beta}}(x)\|_2 - \|\nabla f_\beta(x)\|_2 \right)$$

$$\overset{(ii)}{\leq} \text{err}_{\nabla f_0}(x, \beta) + \frac{1}{\mu} \left( \frac{M_1}{2M_0} \|\nabla f_0(x)\|_2 + L_0 M_0 \right) \|\nabla F(x)^\top\|_2 \cdot \|\hat{\beta} - \beta\|_2$$

$$\overset{(iii)}{\leq} \frac{1}{2} \cdot \frac{2}{\mu} \left( \frac{M_1}{2M_0} \|\nabla f_0(x)\|_2 + L_0 M_0 \right) \left\{ c_2 \varepsilon + \|\nabla F(x)^\top\|_2 \cdot 3\delta \right\} \tag{24}$$

$$\overset{(iv)}{\leq} \frac{1}{2} \varepsilon_0.$$

where (i) expands out $\text{err}_{\nabla f_0}$, (ii) uses the fact that $\beta \mapsto \|\nabla F(x)^\top \beta\|_2$ is $\|\nabla F(x)^\top\|_2$-Lipschitz in $\beta$ with respect to the $\ell_2$-norm, (iii) applies the definition of $\text{err}_{\nabla f_0}$ and the inequality Equation (22), and (iv) follows by definition of $c_1$ and $c_2$, set in the last section of the proof.

As Equations (19) to (21) hold, Lemma 5 shows that $(\hat{x}, \hat{\beta})$ is $(\varepsilon_0, \varepsilon)$-preference stationary.

<u>Case 2</u>: $\|\beta^* - \beta\|_2 \geq 2\delta$. Here $\beta$ is suboptimal and $\beta^*$ achieves a large descent:

$$f_0(x_{\beta^*}) - f_0(x_\beta) \overset{(i)}{\leq} g(\beta^*; x, \beta) - f_0(x_\beta)$$

$$\overset{(ii)}{\leq} \text{err}_{\nabla f_0}(x, \beta) - \frac{1}{2} \mu_g \|\beta^* - \beta\|_2^2$$

$$\overset{(iii)}{\leq} \frac{1}{\mu} \left( \frac{M_1}{2M_0} \|\nabla f_0(x)\|_2 + L_0 M_0 \right) c_2 \varepsilon_0^2 - 2\mu_g \delta^2 \tag{25}$$

$$\overset{(iv)}{\leq} -\frac{3}{2} \mu_g \delta^2, \tag{26}$$

where (i) uses the majorizing property of $g$, (ii) follows from Lemma 14, (iii) applies the definition of $\text{err}_{\nabla f_0}(x, \beta)$ along with the assumption that $\varepsilon \leq \varepsilon_0^2$, and (iv) applies the definition of $c_2$.

The large descent also carries over to $\hat{\beta}$ because it is approximately stationary:

$$f_0(x_{\hat{\beta}}) - f_0(x_\beta) \overset{(i)}{\leq} g(\hat{\beta}; x, \beta) - f_0(x_\beta)$$

$$\overset{(ii)}{=} g(\beta^*; x, \beta) - f_0(x_\beta) + \big(g(\hat{\beta}; x, \beta) - g(\beta^*; x, \beta)\big)$$

$$\overset{(iii)}{\leq} -\frac{3}{2}\mu_g \delta^2 + c_1 \varepsilon_0 \cdot \delta = -\frac{1}{2}\frac{c_1}{\mu_g} \cdot \varepsilon_0^2,$$

where (i) uses the majorizing property of $g$, (ii) adds and subtracts $g(\beta^*; x, \beta)$ and (iii) applies Equation (26) and Lemma 13.

Thus, the preference improves by at least a constant. To finish proving the claim, we need to verify that it is indeed possible to set $c_1$ and $c_2$ appropriately.

Setting $c_1$ and $c_2$: we tabled a few inequalities above. Recall:

For Equation (19), we need:
$$c_2 \leq 1.$$

For Equation (23), we need:
$$2\left(c_1 \varepsilon_0 + \frac{3L \|\nabla f_0(x)\|_2}{\mu^2} \cdot \frac{c_1 \varepsilon_0}{\mu_g}\right) \leq \varepsilon_0.$$

For Equation (24), we need:
$$\frac{2}{\mu}\left(\frac{M_1}{2M_0}\|\nabla f_0(x)\|_2 + L_0 M_0\right)\left\{c_2 \varepsilon + 3\|\nabla F(x)^\top\|_2 \cdot \frac{c_1 \varepsilon_0}{\mu_g}\right\} \leq \varepsilon_0.$$

For Equation (25), we need:
$$\frac{1}{\mu}\left(\frac{M_1}{2M_0}\|\nabla f_0(x)\|_2 + L_0 M_0\right)c_2 \varepsilon_0^2 \leq \frac{1}{2}\mu_g\left(\frac{c_1 \varepsilon_0}{\mu_g}\right)^2.$$

It is unenlightening but straightforward to verify that it suffices to set:

$$c_1 \cdot \max\left\{2 + \frac{6L\|\nabla f_0(x)\|_2}{\mu^2 \cdot \mu_g}, \frac{12}{\mu \cdot \mu_g}\left(\frac{M_1}{2M_0}\|\nabla f_0(x)\|_2 + L_0 M_0\right) \cdot \|\nabla F(x)^\top\|_2\right\} \leq 1$$

$$c_2 \cdot \max\left\{1, \frac{2}{\mu}\left(\frac{M_1}{2M_0}\|\nabla f_0(x)\|_2 + L_0 M_0\right) \cdot \left(2 \vee \frac{\mu_g}{c_1^2}\right)\right\} \leq 1,$$

where $a \vee b := \max\{a, b\}$.

A concerned reader may wonder whether $c_1$ and $c_2$ may be bounded away from zero, as claimed in the theorem statement: we need to ensure that $\|\nabla f_0(x)\|_2$ and $\|\nabla F(x)^\top\|_2$ do not blow up. Indeed, this holds because the iterates $x_k$ remain within a constant distance of the Pareto set. In particular, since $c_2 \leq 1$, by Condition 2, we have that the $k$th iterate satisfies:

$$\|x_k - x_{\beta_k}\| \leq \frac{\varepsilon}{\mu},$$

which follows from $\mu$-strong convexity of $f_{\beta_k}$. Thus, all iterates of the algorithm are within $\varepsilon/\mu$ of the Pareto set and also satisfy for all $k, k' \in \mathbb{N}$:

$$\|x_k - x_{k'}\| \leq R + 2\varepsilon/\mu.$$

Then, by $L_0$-Lipschitz smoothness, we can bound:

$$\|\nabla f_0(x_k)\| \leq \|\nabla f_0(x_1)\| + \|\nabla f_0(x_k) - \nabla f_0(x_1)\|$$
$$\leq \|\nabla f_0(x_1)\| + L_0 \cdot (R + 2\varepsilon/\mu).$$

Similarly, by $L$-Lipschitz smoothness, we also have:

$$\|\nabla F(x_k)^\top\|_2 \leq \|\nabla F(x_1)^\top\|_2 + \|\nabla F(x_k)^\top - \nabla F(x_1)^\top\|_2$$
$$\leq \|\nabla F(x_1)^\top\|_2 + nL \cdot (R + 2\varepsilon/\mu).$$

■

### 9.4.1 Analytic lemma: gradient bound

**Lemma 10.** *Let* $R := \mathrm{diam}\big(\mathrm{Pareto}(F)\big)$. *Then for any* $x_\beta = x^*(\beta)$,

$$\big\|\nabla F(x_\beta)^\top\big\|_{1,2} \leq LR.$$

*Proof.* By definition, we have:

$$\big\|\nabla F(x_\beta)^\top\big\|_{1,2} = \sup_{\|z\|_1 = 1} \bigg\| \sum_{i \in [n]} z_i \nabla f_i(x_\beta) \bigg\|_2$$

$$\overset{(i)}{\leq} \sup_{\|z\|_1 = 1} \sum_{i \in [n]} |z_i| \cdot \|\mathrm{sign}(z_i) \cdot \nabla f_i(x_\beta)\|_2$$

$$\overset{(ii)}{\leq} \max_{i \in [n]} \|\nabla f_i(x_\beta)\|_2$$

$$\overset{(iii)}{\leq} \max_{i \in [n]} L\|x - x_i\|_2,$$

where (i) follows from Jensen's inequality, (ii) holds because the max is no smaller than the average, (iii) applies $L$-Lipschitz smoothness. In particular, let $x_i = \arg\min f_i(x)$, so that $\nabla f_i(x_i) = 0$. Then:

$$\|\nabla f_i(x_\beta) - \nabla f_i(x_i)\|_2 \leq L\|x_\beta - x_i\|_2.$$

The result holds because $x_\beta$ and all $x_i$'s are contained in $\mathrm{Pareto}(F)$. □

## 9.5 Analytic lemmas: matrix inverses

**Lemma 11.** *Let* $M : \mathbb{R}^d \to \mathbb{R}^{d \times d}$ *be* $L$-Lipschitz satisfying $M(x) \succeq \mu\mathbf{I}$ where $\mathbb{R}^d$ has the $\ell_2$-norm and $\mathbb{R}^{d \times d}$ the operator norm. Then, the map $x \mapsto M(x)^{-1}$ is $L/\mu^2$-Lipschitz.

*Proof.* For short, let us denote $M(x)$ by $M_x$. Note that $\mathbf{I} = (M_x' + M_x - M_x')M_x^{-1}$, so that:

$$M_x^{-1} - M_{x'}^{-1} = M_x^{-1} - M_{x'}^{-1}\big(M_{x'} + M_x - M_{x'}\big)M_x^{-1}$$
$$= M_x^{-1} - M_x^{-1} - M_{x'}^{-1}\big(M_x - M_{x'}\big)M_x^{-1} = -M_{x'}^{-1}\big(M_x - M_{x'}\big)M_x^{-1},$$

which is series of unenlightening algebraic manipulations. But now, we may apply $L$-Lipschitz continuity to obtain $\|M_x - M_{x'}\| \leq L\|x - x'\|$ and the $\mu$-lower bound to obtain $\|M_x^{-1}\|, \|M_{x'}^{-1}\| \leq \mu^{-1}$. Together, we obtain $L/\mu^2$-Lipschitz continuity:

$$\big\|M(x)^{-1} - M(x')^{-1}\big\| \leq \frac{L}{\mu^2}\|x - x'\|.$$

□

**Lemma 12.** *Let $M_1, \ldots, M_n$ be positive-definite matrices in $\mathbb{R}^{d \times d}$ equipped with the operator norm, and let $\Delta^{n-1}$ be equipped with the $\ell_1$ norm. Suppose the following holds:*

$$\mu \mathbf{I} \preceq M_1, \ldots, M_n \preceq L\mathbf{I}.$$

*Then, the map $\beta \mapsto M_\beta^{-1}$ where $M_\beta := \sum_{i \in [n]} \beta_i M_i$ has bounded derivative $\|\nabla_\beta M_\beta^{-1}\|_{1,2} \leq L/\mu^2$.*

*Proof.* We can compute the derivative of the above map:

$$\nabla_\beta M_\beta^{-1} = -M_\beta^{-1}(\nabla_\beta M_\beta) M_\beta^{-1},$$

where $\nabla_\beta M_\beta d\beta = M_{d\beta}$. The upper bound on the $M_i$'s implies that $\|\nabla_\beta M_\beta\|_{1,2} \leq L$. And on the other hand, the lower bound implies that $\|M_\beta^{-1}\|_2 \leq \mu^{-1}$. $\qquad\square$

### 9.5.1  Analytic lemmas: constrained optimization of strongly convex functions

**Lemma 13.** *Let $f : \mathbb{R}^n \to \mathbb{R}$ be smooth and convex and let $\mathcal{C} \subset \mathbb{R}^n$ be a convex constraint set. Suppose that $\beta^*, \hat{\beta} \in \mathcal{C}$ are stationary and $\varepsilon$-approximately stationary, respectively:*

$$-\nabla f(\beta^*)^\top(\beta - \beta^*) \leq 0 \qquad and \qquad -\nabla f(\hat{\beta})^\top(\beta - \hat{\beta}) \leq \varepsilon\|\beta - \hat{\beta}\|, \quad \forall \beta \in \mathcal{C}.$$

*Then, $f(\hat{\beta}) - f(\beta^*) \leq \varepsilon\|\hat{\beta} - \beta^*\|$. Furthermore, if $f$ is $\mu$-strongly convex, then $\|\hat{\beta} - \beta^*\| \leq \varepsilon/\mu$.*

*Proof.* For the first part, we apply the mean value theorem, which states that there exists some $\beta$ that is a convex combination of $\hat{\beta}$ and $\beta^*$ such that:

$$\begin{aligned}
f(\hat{\beta}) - f(\beta^*) &\overset{(i)}{=} \nabla f(\beta)^\top(\hat{\beta} - \beta^*) \\
&\overset{(ii)}{\leq} \nabla f(\hat{\beta})^\top(\hat{\beta} - \beta^*) \\
&\overset{(iii)}{\leq} \varepsilon\|\hat{\beta} - \beta^*\|,
\end{aligned}$$

where (i) applies the mean value theorem, (ii) uses the monotonicity of gradients of convex functions:

$$\left(\nabla f(\hat{\beta}) - \nabla f(\beta)\right)^\top(\hat{\beta} - \beta) \geq 0$$

and that $\hat{\beta} - \beta = \lambda(\hat{\beta} - \beta^*)$ for some $\lambda \in [0, 1]$, and (iii) applies the $\varepsilon$-stationarity condition.

For the second part, by strong convexity, we have on the one hand:

$$\left(\nabla f(\hat{\beta}) - \nabla f(\beta^*)\right)^\top(\hat{\beta} - \beta^*) \geq \mu\|\hat{\beta} - \beta^*\|^2.$$

And on the other, by stationarity and $\varepsilon$-stationarity, we have that:

$$\left(\nabla f(\hat{\beta}) - \nabla f(\beta^*)\right)^\top(\hat{\beta} - \beta^*) \geq \varepsilon\|\hat{\beta} - \beta^*\|.$$

Dividing through by $\|\hat{\beta} - \beta^*\|$ yields the result. $\qquad\square$

**Lemma 14.** *Let $\mathcal{C} \subset \mathbb{R}^n$ be a convex constraint set with $\beta \in \mathcal{C}$, and let $Q : \mathcal{C} \to \mathbb{R}$ be a quadratic:*

$$Q(\beta') = c + v^\top(\beta' - \beta) + \frac{1}{2}C\|\beta' - \beta\|^2, \tag{27}$$

*where $c \in \mathbb{R}$, $v \in \mathbb{R}^n$, and $C > 0$. Let $\beta^* \in \mathcal{C}$ minimize $Q$. If $\|\beta^* - \beta\| \geq \varepsilon > 0$, then:*

$$Q(\beta^*) - Q(\beta) \leq -\frac{1}{2}C\varepsilon^2.$$

29

*Proof.* Define the quadratic function $q : \mathbb{R} \to \mathbb{R}$ by:

$$q(\lambda) = c + \lambda v^\top (\beta^* - \beta) + \frac{1}{2} C \lambda^2 \|\beta^* - \beta\|^2$$

$$= c + \frac{1}{2} C \|\beta^* - \beta\|^2 \lambda (\lambda - 2\lambda^*) \tag{28}$$

where $\lambda^* = -\frac{v^\top (\beta^* - \beta)}{C \|\beta^* - \beta\|^2}$ minimizes $q$. Restricting $Q$ to the line between $\beta$ and $\beta^*$, we get:

$$Q\big(\beta + \lambda(\beta^* - \beta)\big) = q(\lambda),$$

for $\lambda \in [0, 1]$. This follows by expanding the definition of $Q$.

Notice that $q$ monotonically decreases on the interval $0 \le \lambda \le \lambda^*$, and also that $q$ monotonically increases for $\lambda > \lambda^*$. Because $Q(\beta^*) = q(1)$ minimizes $Q$ on the convex set $\mathcal{C}$, $q$ must be descending on $\lambda \in [0, 1]$. Thus, $1 \le \lambda^*$. It follows that $1 - 2\lambda^* \le -1$. Plugging in into Equation (28), we have:

$$Q(\beta^*) = q(1) \le c - \frac{1}{2} C \|\beta^* - \beta\|^2.$$

Applying $Q(\beta_0) = c$ and $\|\beta^* - \beta\| \ge \varepsilon$ yields the result. $\qquad\square$

# Acknowledgements

# References

Harold P Benson. Optimization over the efficient set. *Journal of Mathematical Analysis and Applications*, 98(2):562–580, 1984.

S Bolintineanu. Minimization of a quasi-concave function over an efficient set. *Mathematical Programming*, 61:89–110, 1993.

Henri Bonnel and Jacqueline Morgan. Semivectorial bilevel optimization problem: penalty approach. *Journal of Optimization Theory and Applications*, 131:365–382, 2006.

Jerald P Dauer. Optimization over the efficient set using an active constraint approach. *Zeitschrift für Operations Research*, 35:185–195, 1991.

Stephan Dempe. *Bilevel optimization: theory, algorithms and applications*, volume 3. TU Bergakademie Freiberg, Fakultät für Mathematik und Informatik, 2018.

Carl Doersch and Andrew Zisserman. Multi-task self-supervised visual learning. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 2051–2060, 2017.

Daxiang Dong, Hua Wu, Wei He, Dianhai Yu, and Haifeng Wang. Multi-task learning for multiple language translation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1723–1732, 2015.

Joerg Fliege, LM Grana Drummond, and Benar Fux Svaiter. Newton's method for multiobjective optimization. *SIAM Journal on Optimization*, 20(2):602–626, 2009.

János Fülöp. On the equivalence between a linear bilevel programming problem and linear optimization over the efficient set. *Techn. Rep. WP*, pages 93–1, 1993.

Shaona Ghosh, Chris Lovell, and Steve R Gunn. Towards Pareto descent directions in sampling experts for multiple tasks in an on-line learning paradigm. In *2013 AAAI Spring Symposium Series*, 2013.

Andreia P. Guerreiro, Carlos M. Fonseca, and Luís Paquete. The hypervolume indicator: Computational problems and algorithms. *ACM Comput. Surv.*, 54(6), 2021.

Naoki Hamada, Kenta Hayano, Shunsuke Ichiki, Yutaro Kabata, and Hiroshi Teramoto. Topology of Pareto sets of strongly convex problems. *SIAM Journal on Optimization*, 30(3):2659–2686, 2020.

Moritz Hardt and Benjamin Recht. *Patterns, predictions, and actions: Foundations of machine learning.* Princeton University Press, 2022.

Claus Hillermeier. Generalized homotopy approach to multiobjective optimization. *Journal of Optimization Theory and Applications*, 110(3):557–583, 2001.

Zhen Huang, Jinyu Li, Sabato Marco Siniscalchi, I-Fan Chen, Ji Wu, and Chin-Hui Lee. Rapid adaptation for deep neural networks through multi-task learning. In *Sixteenth Annual Conference of the International Speech Communication Association*, 2015.

C-L Hwang and Abu Syed Md Masud. *Multiple objective decision making—methods and applications: a state-of-the-art survey*, volume 164. Springer Science & Business Media, 2012.

Michael Jordan, Guy Kornowski, Tianyi Lin, Ohad Shamir, and Manolis Zampetakis. Deterministic nonsmooth nonconvex optimization. In *The Thirty Sixth Annual Conference on Learning Theory*, pages 4570–4597. PMLR, 2023.

Mohammad Mahdi Kamani, Rana Forsati, James Z Wang, and Mehrdad Mahdavi. Pareto efficient fairness in supervised learning: From extraction to tracing. *arXiv preprint arXiv:2104.01634*, 2021.

Ken Kobayashi, Naoki Hamada, Akiyoshi Sannai, Akinori Tanaka, Kenichi Bannai, and Masashi Sugiyama. Bézier simplex fitting: Describing Pareto fronts of simplicial problems with small samples in multi-objective optimization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 2304–2313, 2019.

Guy Kornowski and Ohad Shamir. Oracle complexity in nonsmooth nonconvex optimization. *Advances in Neural Information Processing Systems*, 34:324–334, 2021.

Aditya Kulkarni, Maximilian Kohns, Michael Bortz, Karl-Heinz Küfer, and Hans Hasse. Regularities of pareto sets in low-dimensional practical multi-criteria optimisation problems: analysis, explanation, and exploitation. *Optimization and Engineering*, pages 1–22, 2022.

William G La Cava. Optimizing fairness tradeoffs in machine learning with multiobjective meta-models. *arXiv preprint arXiv:2304.12190*, 2023.

Kenneth Lange, David R Hunter, and Ilsoon Yang. Optimization transfer using surrogate objective functions. *Journal of computational and graphical statistics*, 9(1):1–20, 2000.

Jiajin Li, Anthony Man-Cho So, and Wing-Kin Ma. Understanding notions of stationarity in nonsmooth optimization: A guided tour of various constructions of subdifferential for nonsmooth functions. *IEEE Signal Processing Magazine*, 37(5):18–31, 2020.

Xi Lin, Hui-Ling Zhen, Zhenhua Li, Qing-Fu Zhang, and Sam Kwong. Pareto multi-task learning. In *Advances in Neural Information Processing Systems*, volume 32, 2019.

Xingchao Liu, Xin Tong, and Qiang Liu. Profiling Pareto front with multi-objective Stein variational gradient descent. In *Advances in Neural Information Processing Systems*, volume 34, pages 14721–14733, 2021.

Zhengliang Liu and Matthias Ehrgott. Primal and dual algorithms for optimization over the efficient set. *Optimization*, 67(10):1661–1686, 2018.

Debabrata Mahapatra and Vaibhav Rajan. Multi-task learning with user preferences: Gradient descent with controlled ascent in Pareto optimization. In *International Conference on Machine Learning*, pages 6597–6607. PMLR, 2020.

Natalia Martinez, Martin Bertran, and Guillermo Sapiro. Minimax Pareto fairness: A multi objective perspective. In *International Conference on Machine Learning*, pages 6755–6764. PMLR, 2020.

Naoki Marumo, Takayuki Okuno, and Akiko Takeda. Majorization-minimization-based Levenberg–Marquardt method for constrained nonlinear least squares. *Computational Optimization and Applications*, pages 1–42, 2023.

I Maruşciac. On Fritz John type optimality criterion in multi-objective optimization. *Mathematica-Revue d'analyse numérique et de théorie de l'approximation. L'analyse numérique et la théorie de l'approximation*, pages 109–114, 1982.

Kaisa Miettinen. *Nonlinear multiobjective optimization*, volume 12. Springer Science & Business Media, 1999.

Yu Nesterov. Gradient methods for minimizing composite functions. *Mathematical programming*, 140(1):125–161, 2013.

Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2003.

Johan Philip. Algorithms for the vector maximization problem. *Mathematical programming*, 2: 207–229, 1972.

Matteo Pirotta and Marcello Restelli. Inverse reinforcement learning through policy gradient minimization. In *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.

Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., USA, 1st edition, 1994.

Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. *Advances in neural information processing systems*, 31, 2018.

Hila Sheftel, Oren Shoval, Avi Mayo, and Uri Alon. The geometry of the Pareto front in biological phenotype space. *Ecology and evolution*, 3(6):1471–1483, 2013.

Steve Smale. Global analysis and economics I: Pareto optimum and a generalization of Morse theory. In *Dynamical systems*, pages 531–544. Elsevier, 1973.

Michael Spivak. *Calculus on manifolds: a modern approach to classical theorems of advanced calculus*. CRC press, 2018.

Ralph E Steuer. The Tchebycheff procedure of interactive multiple objective programming. In *Multiple criteria decision making and risk analysis using microcomputers*, pages 235–249. Springer, 1989.

Pham Dinh Tao et al. Numerical solution for optimization over the efficient set by dc optimization algorithms. *Operations Research Letters*, 19(3):117–128, 1996.

PT Thach, H Konno, and D Yokota. Dual approach to minimization on the set of Pareto-optimal solutions. *Journal of optimization theory and applications*, 88:689–707, 1996.

Yoshitsugu Yamamoto. Optimization over the efficient set: overview. *Journal of Global Optimization*, 22:285–317, 2002.

Mao Ye and Qiang Liu. Pareto navigation gradient descent: a first-order algorithm for optimization in Pareto set. In *Uncertainty in Artificial Intelligence*, pages 2246–2255. PMLR, 2022.

Jingzhao Zhang, Hongzhou Lin, Stefanie Jegelka, Suvrit Sra, and Ali Jadbabaie. Complexity of finding stationary points of nonconvex nonsmooth functions. In *International Conference on Machine Learning*, pages 11173–11182. PMLR, 2020.