



HAL
open science

Stochastic Tangential Pareto Dynamics Provably Samples the Whole Pareto Set

Zachary Jones, Pietro Marco Congedo, Olivier Le Maitre

► **To cite this version:**

Zachary Jones, Pietro Marco Congedo, Olivier Le Maitre. Stochastic Tangential Pareto Dynamics Provably Samples the Whole Pareto Set. 2025. hal-04954830

HAL Id: hal-04954830

<https://inria.hal.science/hal-04954830v1>

Preprint submitted on 18 Feb 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

1 **STOCHASTIC TANGENTIAL PARETO DYNAMICS PROVABLY**
2 **SAMPLES THE WHOLE PARETO SET ***

3 ZACHARY JONES, PIETRO MARCO CONGEDO, OLIVIER LE MAÎTRE

4 **Abstract.** The framework of stochastic multi-objective programming allows for the inclusion of
5 uncertainties in multi-objective optimization problems at the cost of transforming the set of objectives
6 into a set of expectations of random quantities. The stochastic multigradient descent algorithm
7 (SMGDA) gives a solution to these types of problems using only noisy gradient information. However,
8 a bias in the algorithm causes it to converge to only a subset of the whole Pareto front, limiting its
9 use. We analyze the source of this bias and prove the convergence of SMGDA to a stationary point
10 in the nonconvex L-lipschitz smooth case. First, based on this analysis, we propose to reduce the
11 bias of the stochastic multi-gradient calculation using an exponential smoothing technique. We then
12 propose a novel approach to exploring the whole Pareto set by combining the debiased stochastic
13 multigradient with an additive non-vanishing noise that guides the dynamics of the iterates tangential
14 to the Pareto set. We finish by proving that our algorithm, Stochastic Tangential Pareto Dynamics
15 (STPD), generates samples concentrated on the whole Pareto set.

16 **Key words.** Stochastic Programming, Multi-Objective Optimization, Stochastic Multi-Objective
17 Optimization, Gradient, Sampling, Pareto Front, Pareto Set

18 **MSC codes.** 68Q25, 68R10, 68U05

19 **1. Introduction.** Finding all possible solutions to an optimization problem with
20 competing objectives and uncertainty requires specialized methods. In general there is
21 no unique solution to an optimization problem with competing objectives but rather
22 a collection of possible tradeoffs. A Pareto optimal point is then a design point
23 which cannot be improved along any objective without degrading the performance of
24 others. *the Pareto set* is the set of all Pareto optimal points. Introducing stochasticity
25 into the objective functions complicates the process of finding the Pareto set. The
26 objective functions then become random variables making both optimization and the
27 determination of Pareto optimality into probabilistic tasks.

28 To handle the challenges of stochastic multi-objective optimization we work within

*Submitted to the editors DATE.

Funding: This work was funded by Horizon NEXTAIR.

29 the framework of stochastic programming. We formulate our stochastic multi-objective
30 optimization problem as the minimization of a set of expected values of random quan-
31 tities of interest. This changes the stochastic problem into a deterministic one. Unfor-
32 tunately, we rarely have access to analytic expressions for the expected values of the
33 quantities of interest. Furthermore, estimating them can prove to be computationally
34 costly. Therefore, we turn to stochastic approximation.

35 The stochastic multi-gradient descent algorithm (SMGDA) is the extension of
36 stochastic approximation methods to the multi-objective context. It has attracted
37 attention in engineering, operations research, and machine learning fields [5, 11, 9, 4].
38 At each iteration SMGDA uses stochastic gradient information from each objective to
39 compute a direction of descent, the stochastic multigradient. This approach has two
40 main benefits. SMGDA iterations can be performed with as few as a single gradient
41 sample from each objective, making estimation unnecessary. It also has the advantage
42 of provably converging to the Pareto set [9]. However, the stochastic multigradient
43 has a known bias [5, 11]. Because of this bias, previous analysis of the convergence
44 of SMGDA beyond the strongly convex setting have depended on strong assumptions
45 or even concluded that the algorithm does not converge [11, 5, 4].

46 Our first contribution is to analyze the source and effect of this bias. We prove
47 that, despite bias in the stochastic multigradient, SMGDA converges to a *subset*
48 of the Pareto set in the convex L -Lipschitz smooth setting. We additionally show
49 convergence of SMGDA in the strongly convex and nonconvex L -Lipschitz smooth
50 settings.

51 Our second contribution is a sampling approach to generate candidate Pareto
52 optimal points, stochastic tangential Pareto dynamics (STPD). It has two main in-
53 gredients, a debiased stochastic multigradient step and an additive noise term. To
54 debias the stochastic multigradient, based on our analysis of SMGDA, we propose
55 a method using exponential smoothing. We then combine this debiased stochastic
56 multigradient with a nonvanishing noise term which adds noise perpendicular to the
57 expected stochastic multigradient. This noise term guides the sampling procedure
58 tangentially along the Pareto set. We go on to prove that STPD generates samples

59 which lie arbitrarily close to the whole Pareto set in finite time.

60 Our companion paper then tackles the complementary problem of probabilistic
61 inference of the Pareto set using a noisy fixed-size sample set [CITATION NEEDED].

62 We begin the manuscript by recalling some basic concepts in both deterministic
63 and stochastic multi-objective optimization in Section 2. We then introduce the scalar
64 potential for multi-objective optimization in Section 3.1, and use it to prove the
65 convergence of MGDA in Section 3.2. We use the foundation built in Section 3 to
66 make new, more general, proof of convergence of SMGDA in Section 4.1. We also see
67 that SMGDA converges only to a subset of the Pareto front. We examine the bias
68 of SMGD in Section 4.2. We introduce our sampling approach, Stochastic Tangential
69 Pareto Dynamics, in Section 5. We continue to show that it samples the whole
70 of the Pareto front under mild assumptions and offer discussion about the potential
71 applications, pitfalls, improvements, and variations to our proposed approach. Finally,
72 we give our conclusions and future outlook in Section 6

73 **2. Background on Gradient Based (Stochastic) Multi-Objective Op-**
74 **timization.** We start by introducing stochastic programming before reviewing the
75 multigradient descent algorithm and Pareto concepts.

76 Let us look at the stochastic multi-objective programming problem. Given a
77 probability space $(\Theta, \mathcal{F}, \mu)$ and a set of k quantities of interest collected as a vector
78 valued function $F := [f_1(\mathbf{x}, \mathbf{W}(\theta)), \dots, f_k(\mathbf{x}, \mathbf{W}(\theta))] : \mathcal{X} \subseteq \mathbb{R}^d \times \mathbf{W}(\Theta) \mapsto \mathcal{Y} \subseteq \mathbb{R}^k$, our
79 goal is to find *all* \mathbf{x}^* such that

$$80 \quad (2.1) \quad \mathbf{x}^* \in \underset{\mathbf{x} \in \mathcal{X}}{\operatorname{argmin}} \{ \mathbb{E}[f_1(\mathbf{x}, \mathbf{W}(\theta))], \dots, \mathbb{E}[f_k(\mathbf{x}, \mathbf{W}(\theta))] \}.$$

81 This formulation is sufficiently general to incorporate mean/variance minimization,
82 conditional value at risk, and quantile functions as objectives. Notably, the problem
83 is deterministic. If we have analytical expressions for $G(\mathbf{x}) = [g_1(\mathbf{x}), \dots, g_k(\mathbf{x})]$, $\mathcal{X} \mapsto$
84 $\mathcal{Y} \subseteq \mathbb{R}^k$ with $g_i(\mathbf{x}) = \mathbb{E}[f_i(\mathbf{x}, \mathbf{W}(\theta))]$, we can treat this as a deterministic problem

85 and solve it as

$$86 \quad (2.2) \quad \mathbf{x}^* \in \underset{\mathbf{x} \in \mathcal{X}}{\operatorname{argmin}} \{g_1(\mathbf{x}), \dots, g_k(\mathbf{x})\}.$$

87 Clearly we will rarely have a unique solution to problem 2.1. Instead, given two
 88 solutions, $\mathbf{y}', \mathbf{y} \in \mathcal{Y}$, \mathbf{y}' *dominates* \mathbf{y} ($\mathbf{y}' \prec \mathbf{y}$) if it is lower componentwise in all
 89 dimensions, e.g. $\forall j \in \{1, \dots, k\} \mathbf{y}'_j < \mathbf{y}_j$. If no \mathbf{y}' exists which dominates a point \mathbf{y} , we
 90 say it is *undominated*. The Pareto front is defined as the set of undominated points.

91 DEFINITION 2.1 (Pareto Front).

$$92 \quad (2.3) \quad P(\mathcal{Y}) := \{\mathbf{y} \in \mathcal{Y} \mid \nexists \mathbf{y}' \in \mathcal{Y} \setminus \mathbf{y} \text{ s.t. } \mathbf{y}'_j \leq \mathbf{y}_j \forall j \in \{1, \dots, k\}\}.$$

93 This set gives a mathematical description of Pareto optimality. If $\mathbf{y} \in P$ it means
 94 that, in at least one pair of criteria, there is no global performance improvement. We
 95 can define the set of solutions in design space, the Pareto set, as the pre-image of the
 96 Pareto front of G .

97 DEFINITION 2.2 (Pareto Set).

$$98 \quad (2.4) \quad P^{-1}(\mathcal{X}) := (P \circ G)(\mathcal{X}) = \{\mathbf{x} \in \mathcal{X} \mid G(\mathbf{x}) \in P(\mathcal{Y})\}.$$

99 Points \mathbf{x}^* in 2.2 can be found using the multigradient descent algorithm.

100 It is an extension of gradient descent to the multi-objective context using the
 101 multigradient.

102 At each iteration we calculate

$$103 \quad (2.5) \quad \nabla_x^C \{G\}(\mathbf{x}_t) = \operatorname{argmax}_{\mathbf{d}} \max_i \langle \nabla g_i(\mathbf{x}_t), \mathbf{d} \rangle + \frac{1}{2} \|\mathbf{d}\|_2^2.$$

104 Under sufficient conditions on G , the sequence of points $\mathbf{x}_{t+1} = \mathbf{x}_t - \varepsilon \nabla_x^C \{G\}(\mathbf{x}_t)$
 105 will then converge to a Pareto stationary point

106 DEFINITION 2.3 (Pareto Stationary Point). \mathbf{x}^* such that

107 (2.6)
$$\nabla_x^c \{G\}(\mathbf{x}^*) = \mathbf{0}$$

108 In the case where all objective functions are convex this forms a sufficient condition
109 for Pareto optimality [6].

110 In later sections, we will use Pareto stationarity as a tool to search for Pareto
111 optimal points.

112 In practice, solving Eq. 2.5 in its primal form is cumbersome, and we will use its
113 dual formulation, which reduces to a quadratic subproblem in the number of objec-
114 tives,

115 (2.7)
$$\nabla_x^c \{G\}(\mathbf{x}) = \sum_i \alpha_i^* \nabla g_i(\mathbf{x})$$

116
$$\text{s.t. } \alpha^* = \operatorname{argmin}_{\alpha \in \Delta^{k-1}} \left\| \sum_i \alpha_i \nabla g_i(\mathbf{x}) \right\|_2^2,$$

117

118 where Δ^{k-1} is the $k - 1$ dimensional simplex. This subproblem has a closed form
119 solution in the bi-objective case and can be efficiently solved for more objectives
120 using the Franke-Wolfe algorithm [3]. The multigradient, $\nabla_x^c \{G\}(\mathbf{x})$, can be concisely
121 understood as the minimum-norm convex combination of gradients, as seen in Fig. 1.

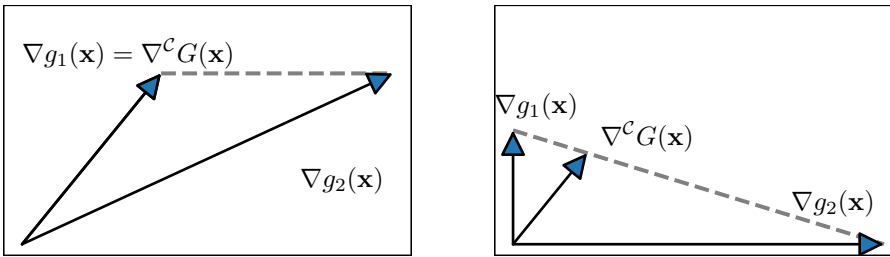


Fig. 1: Examples of Calculation of Direction of Descent $\nabla_x^c \{G\}(\mathbf{x})$ for two objectives.

122

123 **3. Solving Deterministic Multi-Objective Problems with MGDA.** In

124 this section we reformulate and solve the deterministic multi-objective problem. To

125 reformulate the multi-objective problem we make an equivalent problem with only one
 126 objective, a scalar potential. This provides a new way prove convergence of MGDA
 127 and to view multi-objective optimization. In Section 4 we will extend these techniques
 128 to the stochastic case and in Section 5 we will exploit this framework to prove that
 129 our proposed approach samples points along the whole Pareto set.

130 **3.1. Potential function for Multi-Objective Optimization.** To reformu-
 131 late the multi-objective problem we make an equivalent problem with only one objec-
 132 tive, a scalar potential. Previous attempts have used the coefficients, α^* , defined in
 133 Eq. 2.7 to define a pseudo-objective $\Psi(\mathbf{x}) = \sum_i \alpha_i^* g_i(\mathbf{x})$ [5]. However, the parameters
 134 α^* are not constant with respect to \mathbf{x} and so the gradient, $\nabla_x^{\mathcal{C}}\{G\}(\mathbf{x})$, is *not* $\nabla\Psi(\mathbf{x})$.
 135 Consequently $\Psi(\mathbf{x})$ is minimized only at the minima of each individual objective and
 136 not over the whole front.

137 We propose a potential function that has minima over the whole of the Pareto
 138 set. We define the potential, $\Phi_G(\mathbf{x}) : \mathbb{R}^d \mapsto \mathbb{R}$, as follows:

139 DEFINITION 3.1 (Equivalent Scalar Potential Function).

$$140 \quad (3.1) \quad \Phi_G(\mathbf{x}) = \int_0^1 \langle \nabla_x^{\mathcal{C}}\{G\}(s\mathbf{x}), \mathbf{x} \rangle ds = \int_0^1 \sum_i \alpha_i^*(s\mathbf{x}) \langle \nabla g_i(s\mathbf{x}), \mathbf{x} \rangle ds.$$

141 The potential $\Phi_G(\mathbf{x})$ is \mathcal{C}^1 smooth, locally Lipschitz continuous, and it can be seen
 142 that $\nabla\Phi(\mathbf{x}) = \nabla_x^{\mathcal{C}}\{G\}(\mathbf{x})$ by the gradient theorem. Clearly this function has minima
 143 at Pareto stationary points.

144 It can be readily seen in Fig. 2d that the minima of Φ_G lies on the Pareto set.
 145 In the case where there are nonconvex objectives, the set of Pareto stationary points
 146 can form a saddle point in $\Phi_G(\mathbf{x})$. The global minima remains on the Pareto set. We
 147 remark here that we have used the curve from 0 to \mathbf{x} parameterized by t to perform
 148 the line integral that defines the potential function. This is done for convenience only.
 149 Any continuous curve in \mathcal{X} with endpoint \mathbf{x} can be used so long as the initial point
 150 is consistent and \mathcal{X} is pathwise connected.

151 Using 3.1, we can make claims about the set of objective functions using a single
 152 scalar value.

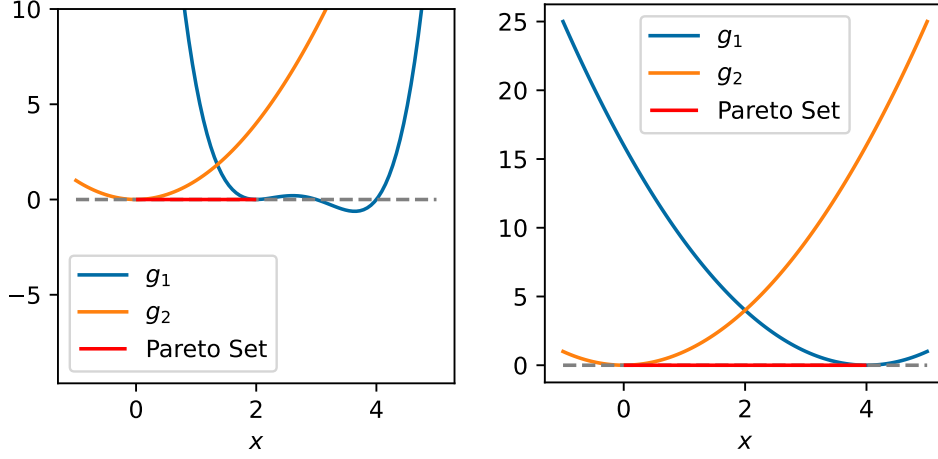
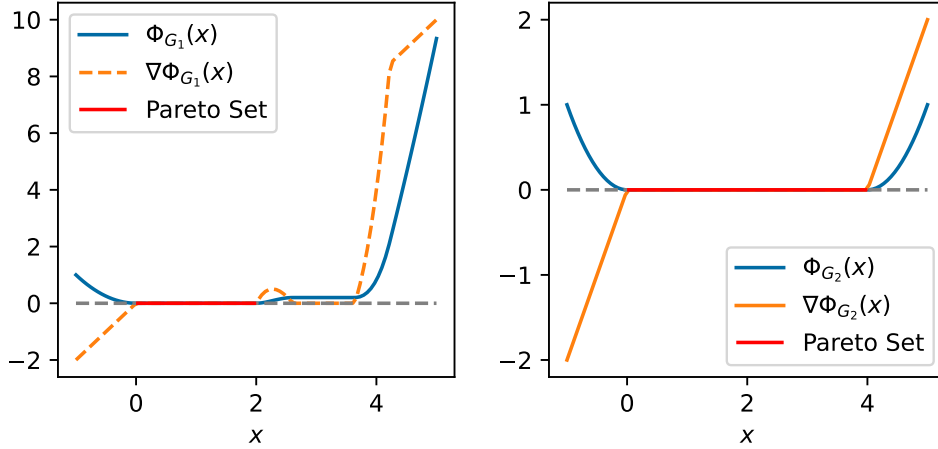
(a) $g_1(\mathbf{x}) = (x-2)^2(x-3)(x-4)$, $g_2(\mathbf{x}) = x^2$.(b) $g_1(\mathbf{x}) = (x-4)^2$, $g_2(\mathbf{x}) = x^2$.(c) $g_1(\mathbf{x}) = (x-2)^2(x-3)(x-4)$, $g_2(\mathbf{x}) = x^2$.(d) $g_1(\mathbf{x}) = (x-4)^2$, $g_2(\mathbf{x}) = x^2$.

Fig. 2: Top row: Competing objective functions in the cases with local minima (left) and bi-convex objectives(right). Bottom row: Minima of the equivalent scalar potential function lie on the Pareto front, with Pareto stationary nonoptimal points revealing themselves as saddle points (left).

153 **3.2. Convergence of MGDA.** We use the scalar potential to make a new proof
 154 of the convergence of MGDA. In the next section, we will use a similar strategy to
 155 show the conditions required for SMGDA to converge.

156 The minima of $\Phi(\mathbf{x})$ are the set of Pareto stationary points, and performing
 157 gradient descent on $\Phi(\mathbf{x})$ yields Pareto stationary points.

158 THEOREM 3.2. Let $\Phi_G(\mathbf{x})$ be a potential function defined as in 3.1. With $\Phi_G :$
 159 $\mathcal{X} \mapsto \mathcal{Y}$, \mathcal{X} compact. Then the sequence

$$160 \quad (3.2) \quad \mathbf{x}_{t+1} = \mathbf{x}_t - \varepsilon \nabla \Phi_G(\mathbf{x}).$$

161 converges to a Pareto stationary point.

162 *Proof.* From the assumption that Φ_G is L-Lipschitz smooth

$$163 \quad (3.3) \quad \Phi_G(\mathbf{x}_{t+1}) \leq \Phi_G(\mathbf{x}_t) + \langle \nabla \Phi_G(\mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{x}_t \rangle + \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2$$

$$164 \quad (3.4) \quad \leq \Phi_G(\mathbf{x}_t) - \varepsilon \left(1 - \frac{L\varepsilon}{2}\right) \|\nabla_x^c \{G\}(\mathbf{x}_t)\|^2.$$

166 Picking $\varepsilon \leq \frac{2}{L}$, summing, using the fact that $\nabla \Phi(\mathbf{x}) = \nabla_x^c \{G\}(\mathbf{x})$, and dividing by T
 167 gives us the result

$$168 \quad (3.5) \quad \frac{1}{T} \sum_{t=1}^T \|\nabla_x^c \{G\}(\mathbf{x}_t)\|^2 \leq \frac{\Phi_G(\mathbf{x}_0) - \Phi_G(\mathbf{x}_{T+1})}{T\varepsilon}. \quad \square$$

169 **4. Solving Stochastic Multi-Objective Problems with SMGDA.** We now
 170 turn our attention to the stochastic formulation of the multigradient descent algo-
 171 rithm, SMGDA. It has already been proven that SMGDA converges to the Pareto
 172 front almost surely under strong assumptions. Previous approaches either require the
 173 individual objective functions are almost surely convex with at least one being almost
 174 surely strongly convex or they assume that the bias in the stochastic multigradient
 175 vanishes in the long time limit.[9, 4, 5]. These cases are not realistic. It has also
 176 been shown that there exist cases where there is a nonvanishing bias in the stochastic
 177 multigradient along the whole Pareto set [11]. The authors went on to conclude that
 178 SMGDA does not converge in these cases. However, this is not necessarily true.

179 Here, we extend the tools we have developed in the previous section to the sto-
 180 chastic case. This allows us to give proof of convergence of SMGDA in the nonconvex
 181 and biased case. We then go on to show that, when Pareto stationary points for which
 182 the stochastic gradient estimator is unbiased exist, that SMGDA converges to one of

183 those points.

184 In the stochastic programming setting the set of objective functions $\{g_i\}_{i=1,\dots,k}$
 185 now have the functional form of an expectation $g_i(\mathbf{x}) = \mathbb{E}[f_i(\mathbf{x}, \mathbf{W}(\theta))]$. Unfortu-
 186 nately, we do not often have analytic expressions for our objectives and must work
 187 with realizations of the quantities of interest, $\{f_i(\mathbf{x}, \mathbf{W}(\theta))\}_{1,\dots,k}$ and their gradients.
 188 It is possible to estimate the value of our objectives and plug them in to a determinis-
 189 tic optimization algorithm. However, this strategy introduces residual error and can
 190 require costly resampling at each iteration. Ignoring the random nature of the prob-
 191 lem by naïvely substituting samples of the quantities of interest, $f_i(\mathbf{x}_t, \mathbf{W}(\theta)_t)$, for
 192 their true value in a deterministic multi-objective method could lead to convergence
 193 to a suboptimal point.

194 Instead, we estimate the gradient of each objective, $\nabla g_i(\mathbf{x})$, using a single real-
 195 ization of the gradient of the corresponding quantity of interest $\nabla f_i(\mathbf{x}, \mathbf{W}(\theta))$. These
 196 estimates are then plugged into Eq. 2.7 to calculate the stochastic multigradient
 197 $\nabla_x^C\{F\}(\mathbf{x}, \mathbf{W}(\theta))$. It is assumed that $\nabla f_i(\mathbf{x}, \mathbf{W}(\theta))$ is an unbiased estimator of $\nabla g(\mathbf{x})$
 198 with finite second moment, e.g. $\mathbb{E}[\nabla f(\mathbf{x}, \mathbf{W}(\theta))] = \nabla g(\mathbf{x})$ and $\mathbb{V}(\nabla f(\mathbf{x}, \mathbf{W}(\theta))) < \infty$
 199 with $\mathbb{V}(\cdot)$ as the variance functional. The heart of the stochastic multigradient descent
 200 algorithm is then the sequence of points generated by the relation

$$201 \quad (4.1) \quad \mathbf{x}_{t+1} = \mathbf{x}_t - \varepsilon_t \nabla_x^C\{F\}(\mathbf{x}, W).$$

202 We require that ε_t meet the requirement that

$$203 \quad (4.2) \quad \frac{\sum_t^T \varepsilon_t^2}{\sum_t^T \varepsilon_t} \xrightarrow{T \rightarrow \infty} 0.$$

204 This is a standard requirement in stochastic optimization [1, 9].

205 It is straightforward then to extend the potential function from Section 3 to the
 206 stochastic multi-objective case.

207 DEFINITION 4.1 (Stochastic Equivalent Scalar Potential Function).

208 (4.3)
$$\Phi_F(\mathbf{x}) = \int_0^1 \langle \mathbb{E}[\nabla_x^{\mathcal{C}}\{F\}(s\mathbf{x})], \mathbf{x} \rangle ds.$$

209 The gradient of this potential function can be approximated using individual realiza-
 210 tions of $\mathbf{W}(\theta)$ as $\nabla_x^{\mathcal{C}}\{F\}(\mathbf{x}, \mathbf{W}(\theta))$.

These quantities are summarized in Tab. 1.

	Objectives	Gradients	Potential	Direction of Ascent	step size
MGDA	$g(\mathbf{x})$	$\nabla g(\mathbf{x})$	Φ_G	$\nabla_x^{\mathcal{C}}\{G\}(\mathbf{x})$	ε
SMGDA	$\mathbb{E}[f(\mathbf{x}, \mathbf{W}(\theta))]$	$\nabla f(\mathbf{x}, \mathbf{W}(\theta)_i)$	Φ_F	$\nabla_x^{\mathcal{C}}\{F\}(\mathbf{x}, \mathbf{W}(\theta))$	ε_t

Table 1: Quantities used in MGDA and their stochastic programming equivalents in SMGDA.

211

212 **4.1. Convergence of SMGDA.** We can now show new proof of the conver-
 213 gence of the SMGD algorithm to a Pareto stationary point in the nonconvex case
 214 with biased gradients. Proofs of convergence for convex and strongly convex ob-
 215 jective functions in the biased setting with bounded gradients can be found in the
 216 appendix.

217 We first make an assumption on the bias of the gradient estimator,

218 ASSUMPTION 4.2. $\exists a > 0, b \geq 0$ such that

219 (4.4)
$$\langle \mathbb{E}[\nabla_x^{\mathcal{C}}\{F\}(\mathbf{x}, \mathbf{W}(\theta))], \nabla_x^{\mathcal{C}}\{G\}(\mathbf{x}) \rangle \geq a \|\nabla_x^{\mathcal{C}}\{G\}(\mathbf{x})\|_2^2 - b.$$

220 *almost everywhere.*

221 This assumption intuitively states that the stochastic gradient is "in-line" with the
 222 true gradient up to a constant bound, b . Many general proofs of the convergence of
 223 SGD in the single objective case make a similar assumption, setting $b = 0$ directly [1].
 224 This assumption is more general. We can always find $b \geq 0$ such that this relationship
 225 holds with arbitrarily high probability. As we will see later, if we can find a pair a, b
 226 such that $b = 0$ and Assumption 4.2 holds almost surely over all $\mathbf{x} \in \mathcal{X}$, then SMGD
 227 converges to a Pareto stationary point. The interested reader can pursue more details

228 in the appendix.

229 We also assume that

230 ASSUMPTION 4.3. $\exists a_G > 2a - 1$ and $b_G \geq 0$ such that

$$231 \quad (4.5) \quad \mathbb{E}[\|\nabla_x^C\{F\}(\mathbf{x}, \mathbf{W}(\theta))\|_2^2] \leq a_G \|\nabla_x^C\{G\}(\mathbf{x})\|_2^2 + b_G$$

232 almost everywhere.

233 Assumptions 4.2 and 4.3 are not strong and imply two things. The first, that the

234 effect of the bias is lower bounded:

$$235 \quad (4.6) \quad \mathbb{E}[\langle \nabla_x^C\{F\}(\mathbf{x}, \mathbf{W}(\theta)), \nabla_x^C\{G\}(\mathbf{x}) \rangle - \|\nabla_x^C\{G\}(\mathbf{x})\|_2^2] \geq (a - 1) \|\nabla_x^C\{G\}(\mathbf{x})\|_2^2 - b.$$

236 Second, setting $M_V \geq (a_g + 1 - 2a)$ and $M_0 \geq 2b + b_G$, we see that assumptions 4.2

237 and 4.3 imply 4.7.

$$238 \quad (4.7) \quad \mathbb{E}[\|\nabla_x^C\{F\}(\mathbf{x}, \mathbf{W}(\theta)) - \nabla_x^C\{G\}(\mathbf{x})\|_2^2] \leq M_V \|\nabla_x^C\{G\}(\mathbf{x})\|_2^2 + M_0$$

239 Readers may recognize Eq. 4.7 as a common assumption in stochastic approximation

240 [1]. We now prove the convergence of SMGDA to a Pareto stationary point in the

241 nonconvex setting.

242 THEOREM 4.4. Let $F = \{f_i(\mathbf{x}, \mathbf{W}(\theta))\}_{i=1, \dots, k}$ be a vector valued function $F : \mathcal{X} \mapsto \mathcal{Y} \subseteq \mathbb{R}^k$ with each $g_i(\mathbf{x}) = \mathbb{E}[f_i(\mathbf{x}, \mathbf{W}(\theta))]$ L -Lipschitz continuous. If there is

243 pair of constants $a > 0$, $b = 0$ such that $\langle \nabla \Phi_F, \nabla \Phi_G \rangle \geq a \|\nabla \Phi_G\|_2^2 - b$ and ε_t meeting

244 requirement 4.2, then the sequence of iterates $\{\mathbf{x}_i\}_{i \in \mathbb{N}}$ generated through the relation

245

$$246 \quad (4.8) \quad \mathbf{x}_{t+1} = \mathbf{x}_t - \varepsilon_t \nabla_x^C\{F\}(\mathbf{x}_t, \mathbf{W}(\theta)_t)$$

247 converges almost surely to a Pareto stationary point.

248 *Proof.* Let $B_t = \mathbb{E}[\nabla\Phi_F(\mathbf{x}_t) - \nabla\Phi_G(\mathbf{x}_t)]$. We start from L-Lipschitz smoothness.

$$249 \quad \Phi_G(\mathbf{x}_{t+1}) \leq \Phi_G(\mathbf{x}_t) + \langle \nabla\Phi_G(\mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{x}_t \rangle + \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|_2^2.$$

250 Taking the conditional expectation, and using the definition of SMGD iterations, we
251 have the expression:

$$252 \quad \Phi_G(\mathbf{x}_{t+1}) \leq \Phi_G(\mathbf{x}_t) - \mathbb{E}[\varepsilon_t \langle \nabla\Phi_G(\mathbf{x}_t), \nabla_x^c \{F\}(\mathbf{x}_t, \mathbf{W}(\theta)_t) \rangle] + \mathbb{E}\left[\frac{L\varepsilon_t^2}{2} \|\nabla_x^c \{F\}(\mathbf{x}_t, \mathbf{W}(\theta)_t)\|_2^2\right].$$

253 Adding and subtracting $\varepsilon_t \|\nabla\Phi_G(\mathbf{x}_t)\|_2^2$, substituting our expression for the bias, and
254 using assumptions 4.2 and 4.3, we see that

$$255 \quad \Phi_G(\mathbf{x}_{t+1}) \leq \Phi_G(\mathbf{x}_t) - \varepsilon_t \|\nabla\Phi_G(\mathbf{x}_t)\|_2^2 - \varepsilon_t \langle B_t, \nabla\Phi_G(\mathbf{x}_t) \rangle + \frac{L\varepsilon_t^2}{2} (b + (a_g + 1 - 2a) \|\nabla\Phi_G(\mathbf{x}_t)\|_2^2)$$

$$256 \quad \leq \Phi_G(\mathbf{x}_t) - \varepsilon_t (a - \varepsilon_t \frac{L}{2} M_V) \|\nabla\Phi_G(\mathbf{x}_t)\|_2^2 + \varepsilon_t b + \frac{L\varepsilon_t^2}{2} M_0.$$

258 Requiring that $\varepsilon_t < \frac{2a}{LM_V} \forall t$, taking the full expectation, summing, and re-arranging,
259 we get

$$260 \quad \sum_t^T \varepsilon_t \|\nabla\Phi_G(\mathbf{x}_t)\|_2^2 \leq \Phi_G(\mathbf{x}_0) - \mathbb{E}[\Phi_G(\mathbf{x}_T)] + b \sum_t \varepsilon_t + \frac{LM_0}{2} \sum_t \varepsilon_t^2.$$

261 Finally, defining $E_T = \sum_t \varepsilon_t$, we can treat ε_t/E_T as a probability measure. Choosing
262 index t in $1, \dots, T$ and applying Markov's inequality yields

$$263 \quad \mathbb{P}(\|\nabla\Phi_G(\mathbf{x})\|_2^2 \geq \epsilon) \leq \frac{\mathbb{E}_{E_t}[\|\nabla\Phi_G(\mathbf{x})\|_2^2]}{\epsilon} \leq \frac{\Phi_G(\mathbf{x}_0) - \mathbb{E}[\Phi_G(\mathbf{x}_T)]}{\epsilon E_T} + \frac{b}{\epsilon} + \frac{LM_0}{2\epsilon E_T} \sum_t \varepsilon_t^2.$$

264 Taking the limit as $T \rightarrow \infty$ gives the result

$$265 \quad \mathbb{P}(\|\nabla\Phi_G(\mathbf{x})\|_2^2 \geq \epsilon) \rightarrow \frac{b}{\epsilon}.$$

266 Therefore, if we can find a pair a, b such that $a > 0$ and $b = 0$ in assumption 4.2 then
267 SMGD converges almost surely to a Pareto stationary point. \square

268 Then, even in the case where the bias does not disappear, the SMGD algorithm can be
 269 shown to converge to a point on the Pareto set. However, the bias still has an effect.
 270 It is straightforward to show that if there is a point, \mathbf{x} where $\mathbb{E}[\nabla_x^C\{F\}(\mathbf{x}, \mathbf{W}(\theta))] =$
 271 $\Phi_G(\mathbf{x}) = \mathbf{0}$ and $b = 0$ then SMGD converges to that point.

272 **THEOREM 4.5.** *Let F , G , ε_t , a , and b be defined as in Thm. 4.4. If there is an*
 273 *$a > 0$ and $b = 0$ meeting the requirements of assumption 4.2 and at least one point,*
 274 *$\mathbf{x}_* \in \mathcal{X}^* := \{\mathbf{x} | \nabla\Phi_F(\mathbf{x}) = \nabla\Phi_G(\mathbf{x}) = \mathbf{0}\}$, then SMGD converges to a point in \mathcal{X}^* .*

275 Let $d_t = \mathbf{x}_t - \mathbf{x}^*$ for some $\mathbf{x}^* \in \mathcal{X}^*$ and $B_t = \nabla\Phi_F(\mathbf{x}_t) - \nabla\Phi_G(\mathbf{x}_t)$. To prove Thm.
 276 4.5 it is enough to show that $\langle B_t, d_t \rangle$ goes to zero. Since $B_t \perp d_t$ on the Pareto set
 277 implies that $\mathcal{X}^* = \emptyset$ there are two remaining cases. Either $B_t = 0$, indicating that
 278 $\nabla\Phi_F(\mathbf{x}_t) = \nabla\Phi_G(\mathbf{x}_t) \rightarrow 0$ and that $\mathbf{x}_t \in \mathcal{X}^*$, or $d_t = 0$ in which case $\mathbf{x}_t = \mathbf{x}_*$.

279 *Proof.* We start from the inequality

$$280 \quad d_{t+1}^2 \leq d_t^2 - 2\varepsilon_t \langle \nabla_x^C\{F\}(\mathbf{x}_t, \mathbf{W}(\theta)_t), d_t \rangle + \varepsilon_t^2 \|\nabla_x^C\{F\}(\mathbf{x}_t, \mathbf{W}(\theta)_t)\|_2^2.$$

281 Then, using our definition of B_t , we have the expression

$$282 \quad d_{t+1}^2 \leq d_t^2 - 2\varepsilon_t \langle B_t, d_t \rangle + \varepsilon_t^2 (M_v + 1) \|\nabla\Phi_G(\mathbf{x}_t)\|_2^2 + \varepsilon_t^2 M_0.$$

283 Rearranging and summing we see that

$$284 \quad (4.9) \quad \sum_t^T 2\varepsilon_t \langle B_t, d_t \rangle \leq d_0^2 - d_T^2 + (M_v + 1) \sum_t^T \varepsilon_t^2 \|\nabla\Phi_G(\mathbf{x}_t)\|_2^2 + M_0 \sum_t^T \varepsilon_t^2.$$

285 As before, defining $E_T = \sum_t^T \varepsilon_t$, we can use Markov's inequality to show

$$286 \quad (4.10) \quad \mathbb{P}_T(\langle B_t, d_t \rangle \geq \epsilon) \leq \frac{d_0^2 - d_T^2}{2\epsilon E_T} + (M_v + 1) \frac{\sum_t^T \varepsilon_t^2 \|\nabla\Phi_G(\mathbf{x}_t)\|_2^2}{2\epsilon E_t} + M_0 \frac{\sum_t^T \varepsilon_t^2}{2\epsilon E_t}.$$

287 Since $\|\nabla\Phi[G](\mathbf{x})\|_2^2 \rightarrow 0$ almost surely, we can conclude that the left hand side of Eq.

288 4.10 goes to zero. □

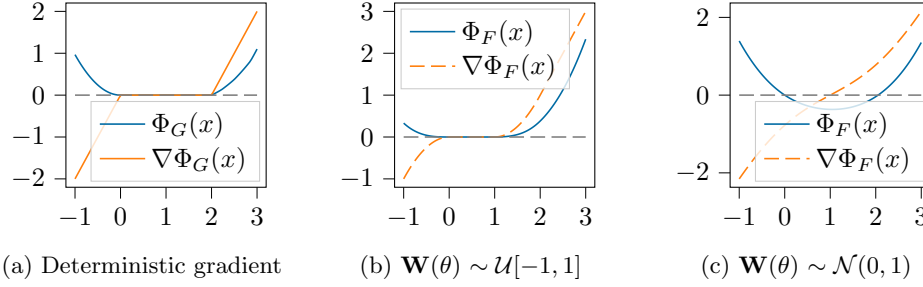


Fig. 3: Comparison of Potential Functions with $f_1(\mathbf{x}, \mathbf{W}(\theta)) = (x - 2 + \mathbf{W}(\theta))^2$ and $f_2(\mathbf{x}, \mathbf{W}(\theta)) = (x + \mathbf{W}(\theta))^2$. SMGDA converges to a point where $\mathbb{E}[\nabla\Phi_F] - \nabla\Phi_G = 0$. Deterministic gradient calculated as $\nabla g_i(\mathbf{x}) = \mathbb{E}[f_i(\mathbf{x}, \mathbf{W}(\theta))]$.

289 **4.2. The Bias of the Stochastic Multigradient.** Why is the stochastic multi-
 290 gradient biased when the stochastic gradients used to compute it are not? By assump-
 291 tion we have access to samples of the jacobian of F , $J_F(\mathbf{x}, \mathbf{W}(\theta))$ which is an unbiased
 292 estimator of the Jacobian of G , $J_G(\mathbf{x})$, e.g. $\mathbb{E}[J_F(\mathbf{x}, \mathbf{W}(\theta))] = J_G(\mathbf{x})$.

293 We look to the minimizer $\alpha^*(\mathbf{x}, \mathbf{W}(\theta)) = \operatorname{argmin}_{\alpha \in \Delta^{k-1}} \alpha^\top J_F(\mathbf{x}, \mathbf{W}(\theta)) J_F(\mathbf{x}, \mathbf{W}(\theta))^\top \alpha$
 294 and notice that it is a quadratic function of the Jacobian. By Jensen's inequality
 295 $\mathbb{E}[J_F(\mathbf{x}, \mathbf{W}(\theta)) J_F(\mathbf{x}, \mathbf{W}(\theta))^\top] \succeq \mathbb{E}[J_F(\mathbf{x}, \mathbf{W}(\theta))] \mathbb{E}[J_F(\mathbf{x}, \mathbf{W}(\theta))]^\top$, and so $\mathbb{E}[\alpha^*(\mathbf{x}, \mathbf{W}(\theta))] \neq$
 296 $\alpha^*(\mathbf{x})$.

297 α^* could be debiased using Monte-Carlo estimation, resampling and averaging
 298 the Jacobian and using $\widehat{J}_G(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N J_F(\mathbf{x}, \mathbf{W}(\theta)_i)$ in place of $J_G(\mathbf{x})$ in Eq. 2.7.
 299 However, not only would it be computationally taxing to compute N gradients at each
 300 iteration, but there would still be residual variance in our estimation rendering our
 301 computational efforts somewhat pyrrhic. Other approaches attempt to smooth the
 302 minimizers α^* using exponential smoothing [11]. This approach effectively reduces
 303 the variance of α^* and enforces that the estimator α^* be smooth across iterations.
 304 In contrast, our approach to the problem is to mitigate the bias at its source, in the
 305 calculation of α^* itself.

306 **5. Proposed Approach: Stochastic Tangential Pareto Dynamics.** The
 307 SMGD algorithm, even if run from distinct starting points, converges to a subset of the
 308 Pareto set. Then, even if completely debiased, in order to estimate the whole Pareto

309 set one has to restart the algorithm from several points in design space. Not only is
 310 this approach not guaranteed to uncover the whole Pareto set but it is computationally
 311 taxing. To ameliorate this, we propose to modify the stochastic multigradient descent
 312 iterations in two ways. First, we introduce a novel debiasing strategy for the stochastic
 313 multigradient. This allows the algorithm to converge to the whole of the Pareto front.
 314 Then, we add a noise term which guides the dynamics generated by the algorithm
 315 tangentially along the Pareto set, promoting exploration as well as overcoming any
 316 residual bias. Because of the added noise term, this approach efficiently explores
 317 the area around the Pareto set and is less likely to stay stuck in a saddle point.
 318 Furthermore, iterations can be performed with single samples of the quantities of
 319 interest and their gradients at each round.

320 Dubbed stochastic tangential Pareto dynamics (STPD), the points generated by
 321 this approach create a dense sample set concentrated on the Pareto set.

322 **5.1. Debiasing the Stochastic Multigradient.** We would like to use the
 323 deterministic Jacobian matrix, J_G , to calculate the multigradient. This would be
 324 equivalent to having direct access to $\mathbb{E}[J_F(\mathbf{x}, \mathbf{W}(\theta))] \mathbb{E}[J_F(\mathbf{x}, \mathbf{W}(\theta))]^\top$. However, we
 325 only have access to samples of the matrix $J_F(\mathbf{x}, \mathbf{W}(\theta))J_F(\mathbf{x}, \mathbf{W}(\theta))^\top$. We notice that
 326 the two quantities are related,

(5.1)

$$327 \quad \mathbb{E}[J_F(\mathbf{x}, \mathbf{W}(\theta))J_F(\mathbf{x}, \mathbf{W}(\theta))^\top] = \Sigma_{J_F(\mathbf{x}, \mathbf{W}(\theta))} + \mathbb{E}[J_F(\mathbf{x}, \mathbf{W}(\theta))] \mathbb{E}[J_F(\mathbf{x}, \mathbf{W}(\theta))]^\top,$$

328 where $\Sigma_{J_F(\mathbf{x}, \mathbf{W}(\theta))}$ denotes the covariance matrix of the gradients of F .

329 We propose to debias the stochastic multigradient using exponential smoothing.
 330 We first compute an *online* estimate of $\Sigma_{J_F(\mathbf{x}, \mathbf{W}(\theta))}$ using the recursive estimate

$$331 \quad (5.2) \quad \hat{\mu}_{t+1} = \gamma_t \hat{\mu}_t + (1 - \gamma_t) J_f(\mathbf{x}_t, \mathbf{W}(\theta)_t)$$

$$332 \quad (5.3) \quad \hat{\Sigma}_{J_F, t+1} = \gamma_t \hat{\Sigma}_{J_F, t} + (1 - \gamma_t) (J_F(\mathbf{x}_t, \mathbf{W}(\theta)_t) J_F(\mathbf{x}_t, \mathbf{W}(\theta)_t)^\top - \hat{\mu}_t \hat{\mu}_t^\top)$$

334 using individual samples of $J_F(\mathbf{x}, \mathbf{W}(\theta))$ and $\gamma_t \in (0, 1)$. We then calculate α^* using

335 the following modified formulation.

$$336 \quad (5.4) \quad \hat{\alpha}^*(\mathbf{x}) = \underset{\alpha \in \Delta^{k-1}}{\operatorname{argmin}} \alpha^\top (J_F(\mathbf{x}, \mathbf{W}(\theta)) J_F(\mathbf{x}, \mathbf{W}(\theta))^\top - \widehat{\Sigma}_{J_F(\mathbf{x}, \mathbf{W}(\theta))}) \alpha.$$

337 We can then calculate a debiased form of the stochastic gradient

$$338 \quad (5.5) \quad \widehat{\nabla_x^C \{F\}}(\mathbf{x}, \mathbf{W}(\theta)) := \sum_i \hat{\alpha}_i^* \nabla f_i(\mathbf{x}, \mathbf{W}(\theta))$$

339 This approach has the advantage of not requiring excess computation in individual
 340 rounds and places less assumptions on the required behavior of α^* while also effectively
 341 reducing the bias in the calculation of $\nabla_x^C \{F\}(\mathbf{x}, \mathbf{W}(\theta))$.

342 **5.2. Stochastic Tangential Pareto Dynamics.** To find and explore the Pareto ■

343 set we generate a sequence of points using the recurrence relation,

$$344 \quad (5.6) \quad \mathbf{x}_{t+1} = \mathbf{x}_t - \varepsilon_t \widehat{\nabla_x^C \{F\}}(\mathbf{x}_t) + \sqrt{2\varepsilon_t \beta_t} \Pi_t Z,$$

345 with $\widehat{\nabla_x^C \{F\}}(\mathbf{x}_t)$ a direction of descent, $Z \sim \mathcal{N}(0, \mathbb{1}_{d \times d})$, Π_t a projection operator that
 346 projects the noise perpendicular to the direction of descent, ε_t a stepsize parameter,
 347 and β_t the scale of the noise. The subscript t in this case should be understood as
 348 being adapted to the filtration generated by \mathbf{x}_t . This algorithm can be intuitively
 349 understood as gradient descent with an additive exploratory noise term. To keep
 350 the samples from converging to a limiting distribution centered on a single point we
 351 set the stepsize parameter, ε_t , and the moving average parameter, γ_t as a pair of
 352 nonincreasing but *strictly positive* numbers:

$$353 \quad (5.7) \quad \begin{aligned} \{\varepsilon_t\}_{t=0, \dots, \infty} &:= \varepsilon_0 \geq \varepsilon_1 \geq \dots \geq \varepsilon_\infty > 0, \\ \{\gamma_t\}_{t=0, \dots, \infty} &:= \gamma_0 \geq \gamma_1 \geq \dots \geq \gamma_\infty > 0. \end{aligned}$$

354 The moving average parameter, as we see below, is used in the estimation of the
 355 debiased stochastic multigradient and its expected value.

356 Ideally, we would have access to true values of the gradients for each objective.

357 It would then be possible to use the deterministic multigradient, $\nabla_x^C\{G\}(\mathbf{x})$, as an
 358 unbiased direction of descent. However, having only samples of the gradients of each
 359 objective, we approximate the direction of descent. Here, the expected value of the
 360 debaised stochastic multigradient as given in Eq. 5.5 defines the gradient of the
 361 potential, $\nabla\Phi_F(\mathbf{x})$. Since we do not have access to its expected value, we approximate
 362 it using single samples.

363 The performance of the estimation of the debaised stochastic multigradient de-
 364 pends on the noise in the stochastic gradients and the strength of the additive noise
 365 term. If the noise in the individual stochastic gradients is low then exponential
 366 smoothing with a large value of γ_∞ will yield a good approximation of the gradi-
 367 ent. As either the strength of the additive noise term or the noise in the gradients
 368 increases the autocorrelation between iterates falls and the exponential smoothing
 369 procedure breaks down. This is viewed with more detial in our companion paper.

370 Π_t removes the component of the additive noise that lies parallel to the multigradi-
 371 dent. On the Pareto set its value is unity and the algorithm produces iterates that
 372 move randomly. Near the Pareto set, the direction orthogonal to the multigradient
 373 lies tangential to the Pareto set. The tangential motion of the iterates generated by
 374 STPD can then be understood as a tangential drift along the Pareto set. If Π_t was
 375 set permanently to unity then this approach would degrade into a type of stochastic
 376 Langevin dynamics. This would be unsatisfactory, however, as the variance of samples
 377 orthogonal to the Pareto set would be much higher.

378 To calculate Π_t we need access to the expected value of the direction of descent.
 379 However, we do not in general have access to an analytic expression for the expected
 380 value of the debaised multigradient. Therefore it must be estimated during the course
 381 of optimization. This estimation is also biased due to the nonzero stepsize.

382 At each iteration we update the expected value of the stochastic multigradient

$$383 \quad \overline{\nabla_x^C\{F\}}_{t+1} = \gamma_t \overline{\nabla_x^C\{F\}}_t + (1 - \gamma_t) \widehat{\nabla_x^C\{F\}}(\mathbf{x}_t). \quad 384$$

385 To calculate the projection operator we use the exponentially smoothed debaised

386 stochastic multigradient,

$$387 \quad (5.9) \quad \Pi_{t,i,j} = \delta_{i,j} - \frac{1}{\|\widehat{\nabla_x^C \{F\}}_t\|_2^2} \widehat{\nabla_x^C \{F\}}_{t,i} \widehat{\nabla_x^C \{F\}}_{t,j},$$

388 where $\delta_{i,j}$ denotes the dirac delta function.

389 The value of β_t determines the relative strength of the additive noise compared to
 390 the direction of descent and can be viewed as a type of temperature parameter. We
 391 only require that β_t be a positive upper bounded function adapted to the filtration
 392 generated by \mathbf{x}_t . It must be positive since a value of zero would remove the noise
 393 term, and it must be upper bounded in order to allow for convergence. We discuss
 394 possible function choices in Section 5.4. The effect of the relative strength of the noise
 395 term is studied numerically in our companion work.

This entire process is summarized in algorithm 5.1.

Algorithm 5.1 Summary of Stochastic Tangential Pareto Dynamics

Input $\{\varepsilon_t\}, \{\gamma_t\}, \mathbf{x}_0$.

Initialize $\widehat{\mu}_0 = \mathbf{0}, \widehat{\Sigma}_{J_F,0} = \mathbb{1}_{d \times d}, t = 1, \widehat{\nabla_x^C \{F\}}_0 = 0$.

while Running **do**

Query Stochastic Oracle for $J_F(\mathbf{x}_t, \mathbf{W}(\theta)_t)$ ▷ Referred to as J_{F_t} below.

$\widehat{\mu}_t \leftarrow \gamma_t \widehat{\mu}_{t-1} + (1 - \gamma_t) J_{F,t}$

$\widehat{\Sigma}_{J_F,t} \leftarrow \gamma_t \widehat{\Sigma}_{J_F,t-1} + (1 - \gamma_t) (J_{F,t} - \widehat{\mu}_t)(J_{F,t} - \widehat{\mu}_t)^\top$

$\widehat{\alpha}^* \leftarrow \operatorname{argmin}_{\alpha \in \Delta^{k-1}} \alpha^\top (J_{F_t} J_{F_t}^\top - \widehat{\Sigma}_{J_F,t}) \alpha$ ▷ Eq. ??

$\widehat{\nabla_x^C \{F\}}_t \leftarrow J_{F_t}^\top \widehat{\alpha}^*$ ▷ Eq. ??

$\widehat{\nabla_x^C \{F\}}_t \leftarrow \gamma_t \widehat{\nabla_x^C \{F\}}_{t-1} + (1 - \gamma_t) \widehat{\nabla_x^C \{F\}}_t$ ▷ Eq. 5.8

$\Pi_t \leftarrow \mathbb{1}_{d \times d} - \frac{\widehat{\nabla_x^C \{F\}}_t \widehat{\nabla_x^C \{F\}}_t^\top}{\|\widehat{\nabla_x^C \{F\}}_t\|_2^2} \widehat{\nabla_x^C \{F\}}_t \widehat{\nabla_x^C \{F\}}_t^\top$ ▷ Eq. 5.9

$Z \sim \mathcal{N}(\mathbf{0}, \mathbb{1}_{d \times d})$

$\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \varepsilon_t \widehat{\nabla_x^C \{F\}}_t + \sqrt{2\varepsilon_t \beta_t} \Pi_t Z$ ▷ Eq. 5.6

$t \leftarrow t + 1$

end while

396

5.3. Stochastic Tangential Pareto Dynamics Samples the Whole Pareto

397 **Front.** Here we show that, given enough time, stochastic tangential Pareto dynamics
 398 samples arbitrarily close to all points on the Pareto set.

400 In addition to assumption 4.2 and L-Lipschitz smoothness of Φ_G , we assume

401 coerciveness of Φ_G .

402 ASSUMPTION 5.1. $\exists c_a, c_b \in \mathbb{R}_+, \forall \mathbf{x}$, such that

$$403 \quad (5.10) \quad \|\nabla \Phi_G(\mathbf{x})\|^2 \geq c_a \Phi_G(\mathbf{x}) - c_b, \quad \|\mathbf{x}\|^2 \leq c_a \Phi_G(\mathbf{x}) + c_b \quad .$$

405 This assumption is coercive in the sense that the strength the gradient increases the
 406 further away \mathbf{x} is from the minima, "pushing" iterates back towards the minima as
 407 the value of the scalar potential gets larger. Assumption 5.1 is necessary to show the
 408 ergodicity of the SDE $dX = \nabla F(X) + dW$, and so it is appropriate in this context as
 409 well [10, 8, 2].

410 We can now show the main theorem.

411 THEOREM 5.2. *Given a set of constants $\{\varepsilon_t\}_{t \in \mathbb{N}}$, assuming $\Phi_G : \mathbb{R}^d \mapsto \mathbb{R}$ is L -
 412 Lipschitz smooth and obeys conditions 5.1 and 4.2 almost everywhere, and assuming
 413 lower bound $B := L\beta^{-1}d + b + \frac{\varepsilon_0 LM_0}{2} + d + c_b \leq c_1 \inf \Phi_G(\mathbf{x})$, then for any $\mathbf{y}^* \in P^{-1}$
 414 and $\epsilon, p > 0$*

$$415 \quad (5.11) \quad \mathbb{P}(\|\mathbf{x}_t - \mathbf{y}^*\| \leq \epsilon \text{ for some } t < \infty) \geq 1 - p.$$

416 Our framework for proving Theorem 5.2 for STPD follows closely from [8, 2, 10].
 417 We first show recurrence, that there is a compact sublevel set, $\Phi_G(\mathbf{x}) \leq M$, that is
 418 reached infinitely many times by STPD. We then show reachability, on a compact
 419 sublevel set, there is a nonzero probability to reach any target point $\mathbf{y}^* \in \mathcal{P}$. Once
 420 we have established both recurrence and reachability, it is straightforward to see that
 421 STPD reaches an arbitrary point \mathbf{y}^* in finite time with probability arbitrarily close
 422 to 1. Since \mathbf{y}^* is arbitrary, it follows that it reaches the whole Pareto set.

423 First we show recurrence.

424 LEMMA 5.3. *Recurrent visits to the sublevel set M .*

425 *Assume that Φ_G is L -Lipschitz smooth, obeys assumptions 5.1 and 4.2, and let $\{\varepsilon_t\}_{t \in \mathbb{N}}$
 426 be a nonincreasing sequence such that $\varepsilon_0 \geq \varepsilon_1, \dots, \geq \varepsilon_\infty > 0$. Let $M > 0$ and define the
 427 constant $B := L\beta^{-1}d + b + \frac{\varepsilon_0 LM_0}{2} + c_b \leq c_a \inf \Phi_G(\mathbf{x})$. Given a sequence of stopping*

428 times $\tau_{k+1} = \inf\{t : t > \tau_k, \Phi_G(\mathbf{x}_t) \leq M\}$ then

$$429 \quad (5.12) \quad a) \quad \tau_0 = \frac{\log \left[\left[\frac{\Phi_G(\mathbf{x}_0)}{M - \frac{\varepsilon_0 B}{\varepsilon_\infty a}} \right]^2 \right]}{2c_a \varepsilon_\infty}$$

430 and

$$431 \quad (5.13) \quad b) \quad \mathbb{E}[\tau_j] = \tau_0 + (j+1)M$$

432 For proof see the appendix.

433 We then have reachability.

434 LEMMA 5.4 (Reachability). Given that $\Phi_G(\mathbf{x}_{\tau_k}) \leq M$, let $\mathbb{E}[\nabla\Phi_F(\mathbf{x})] \leq \mathbf{D}$ and
 435 $c_t > 0$. We have, for $\mathbf{y}^* \in \{\mathbf{y} : \Phi_G(\mathbf{y}) \leq M\} \cap \{\mathbf{y} : \mathbf{y} \in P^{-1}\}$

$$436 \quad (5.14) \quad \mathbb{P}(\|\mathbf{x}_{\tau_j+t} - \mathbf{y}^*\|^2 \leq \epsilon) = c_t > 0$$

437 For proof see the appendix.

438 Having shown both recurrence and reachability established in lemmas 5.3 and
 439 5.4, we now prove Theorem 5.2

440 *Proof.* First, define a sequence of stopping times and a $\tau_{j+1} = \inf\{t : t >$
 441 $\tau_j, \Phi_G(\mathbf{x}_{\tau_j+1}) \leq M\}$. Then, also define a stopping time $\tau^* = \inf t : \|\mathbf{x}_t - \mathbf{y}^*\| \leq \epsilon$
 442 for $\epsilon > 0$. We have that

$$443 \quad \mathbb{P}(\tau^* \geq T) = \mathbb{P}(\tau^* \geq T, \tau_j \geq T) + \mathbb{P}(\tau^* \geq T, \tau_j \leq T)$$

$$444 \quad = \mathbb{P}(\tau_j \geq T) + \mathbb{P}(\tau^* \geq T, \|\mathbf{x}_{\tau_k} - \mathbf{y}^*\| \geq \epsilon \forall k = \{1, \dots, j\})$$

$$445 \quad \leq \mathbb{E}[\tau_j] + \prod_{k=1}^j \mathbb{P}(\|\mathbf{x}_{\tau_k} - \mathbf{y}^*\| \geq \epsilon) \quad (\text{Lem. 5.3})$$

$$446 \quad \leq \frac{(j+1)M + \tau_0}{T} + \prod_{k=1}^j (1 - c_j) \quad (\text{Lem. 5.4})$$

$$447 \quad \leq \frac{(j+1)M + \tau_0}{T} + (1 - c^*)^j.$$

448

449 Where we have used the boundedness of c_j in the final line. Setting $\frac{(j+1)M+\tau_0}{T} = \alpha p$
 450 and $(1 - c^*)^j = (1 - \alpha)p$ for $\alpha \in (0, 1)$ and solving for j and T we see that setting

$$451 \quad j \geq \frac{\ln((1 - \alpha)p)}{1 - c^*} \quad (\text{and}) \quad T \geq \frac{\left(\frac{\ln((1 - \alpha)p)}{1 - c^*} + 1\right)M + \tau_0}{\alpha p}$$

453 gives us $\mathbb{P}(\tau^* \geq T) \leq p$. Since this is true for any p ,

$$454 \quad \mathbb{P}(\tau^* \leq T) \geq 1 - p,$$

455 which can be made arbitrarily close to 1. □

456 **5.4. Discussion.** Here we consider implementation choices, the effect of non-
 457 stationarity of the sample distribution, parameter settings, and dealing with saddle
 458 points and nonconvexity.

459 We have presented here a basic version of our approach with the aim of imparting
 460 understanding about the general working of the algorithm. There are other techniques
 461 which may be included in the algorithm. They do not affect our theoretical results
 462 and may improve the sampling properties of the algorithm in practice.

463 We first consider variance reduction. Variance reduction techniques can be used
 464 on both $\nabla f_i(\mathbf{x}, \mathbf{W}(\theta))$ and the debiased stochastic multigradient $\widehat{\nabla_x^c \{F\}}(\mathbf{x})$. Reduc-
 465 ing the variance of the objective gradients would further decrease the bias of the
 466 calculation of $\widehat{\nabla_x^c \{F\}}(\mathbf{x})$. It would also cause the points generated by the algorithm
 467 to be more concentrated on the Pareto set. Common approaches include the pop-
 468 ular SAGA, SVRG, and minibatching [1]. The increased precision comes at a cost,
 469 for a given stepsize lower gradient variance will also decrease the speed at which the
 470 algorithm diffuses along the Pareto set. The use of variance reduction is ultimately
 471 problem dependent and will affect parameter choices for stepsize and additive noise
 472 strength.

473 It is also possible to include preconditioners in our approach, i.e. take steps

$$474 \quad (5.15) \quad \mathbf{x}_{t+1} = \mathbf{x}_t - \varepsilon_t H_t \Phi_F(\mathbf{x}) + \sqrt{2\beta_t \varepsilon_t} H_t \Pi_t Z$$

475 where H_t is a positive definite matrix. There are many options for the preconditioner,
 476 as discussed in the context of single objective stochastic Langevin dynamics in [7].

477 The samples generated by STPD do not reach a stationary distribution. This is
 478 because of the nonvanishing stepsize. As a result the estimates for the jacobian, its
 479 covariance, and the debiased stochastic multigradient do not converge. This appears
 480 in the algorithm as a bias in both the direction of descent and the projection operator.
 481 This bias is ultimately overcome by the variation introduced by the additive noise.

482 The bias in the direction of descent lies tangential to the Pareto set. Clearly, im-
 483 mediately projecting normal to this direction would be counterproductive. Averaging
 484 over the history of the direction of descent allows us to estimate a true tangential di-
 485 rection. This average, however, still contains a trace of the bias, causing the projected
 486 noise to have a component normal to the Pareto set. This manifests as an increased
 487 residual variance of the samples around the Pareto set.

488 Projecting normal to the average history of the direction of descent is not the only
 489 way to calculate the projection operator. It is also possible to average the history of
 490 the directions perpendicular to the direction of descent, e.g. calculate Π_t as:

$$491 \quad \Pi_{t+1} = \gamma_t \Pi_t + \left(\mathbb{1} - \frac{1}{\|\widehat{\nabla_x^c \{F\}}(\mathbf{x}_t)\|_2^2} \widehat{\nabla_x^c \{F\}}(\mathbf{x}_t) \widehat{\nabla_x^c \{F\}}(\mathbf{x}_t)^\top \right)$$

492 We have chosen exponential smoothing to model several quantities. It has the
 493 advantage of requiring little storage and being easy to calculate. However, its ef-
 494 fectiveness depends strongly on the correlation time of the sample gradients. The
 495 correlation time, in turn, is influenced by the noise of the stochastic gradients along
 496 the Pareto set and the relative strength of the additive noise term. These interactions
 497 are not straightforward. More complex models can be used in the place of simple
 498 exponential smoothing. However, they will similarly show a dependence on the corre-
 499 lation time of the sample gradients and must be calibrated. This can be done simply
 500 by minimizing the predictive error of the estimation of a quantity of interest, e.g. the
 501 Jacobian.

502 It has been previously mentioned that the additive noise term β_t need only be a

503 positive upper bounded function adapted to the filtration generated by \mathbf{x}_t . This leaves
 504 flexibility in its choice beyond a constant function. Since the additive noise term may
 505 interfere with initial convergence of the iterates to the Pareto set it might be beneficial
 506 to damp it when iterates are far from the Pareto set. The strength of the additive noise
 507 can be damped far from the Pareto set by selecting $\beta_t = Ce^{-\overline{\nabla_x^c\{F\}}^\top A \overline{\nabla_x^c\{F\}} + D}$ with
 508 $C, D \in \mathbb{R}^{d \times d}$, $A \in \mathbb{R}_+^{d \times d}$. The dynamics generated by STPD are also strongly affected
 509 by the relative strength of the additive noise term. This suggests a scaling to augment
 510 the strength of the additive noise in proportion to the variation of the objective
 511 functions by selecting $\beta_t = \sum_i \text{VAR}(\nabla f_i(\mathbf{x}_t, \mathbf{W}(\theta)_t))$. Vector valued functions can
 512 also be considered.

513 We take a moment to discuss saddle points in the scalar potential. We have
 514 proven that the iterates generated by STPD cover the whole Pareto set under mild
 515 assumptions. This does not preclude cases with non-convex objectives. It has been
 516 shown in the single objective case that adding Gaussian noise to stochastic gradients
 517 can allow gradient descent algorithms to escape shallow second order minima [2]. In
 518 the multi-objective case, local minima in the objective functions correspond to saddle
 519 points in the scalar potential which, if mild, may be similarly escaped. However,
 520 a potential with a split Pareto set (a 'W' shaped potential) may lead to complica-
 521 tions. Either multiple restarts or an alternative approach more suited to nonconvex
 522 optimization may be more effective for such situations.

523 Diffusing across the Pareto set is not the only way to exploit the stochastic multi-
 524 gradient. An alternative approach would be to select an ensemble of random points
 525 in design space and optimize them individually using debiased SMGDA. This ap-
 526 proach may be more effective on non-convex problems and problems with discontin-
 527 uous Pareto sets. However, it has two corresponding drawbacks. One is the loss of
 528 information. Two points from distinct regions in design space may converge to similar
 529 locations on the Pareto set, making one redundant. Another drawback is resolution.
 530 Picking an ensemble of points pre-determines the resolution at which one can deter-
 531 mine the Pareto set. This has a similar flavor to picking the number of iterations of
 532 STPD. The computational cost, however, is not the same. STPD scales linearly with

533 the number of samples taken, whereas adding another point in an ensemble algorithm
534 scales multiplicatively with the number of iterations used.

535 **6. Conclusion.** We have introduced a new approach to solving multi-objective
536 optimization problems using noise added debiased stochastic multigradients, in which
537 the whole of the Pareto front is of interest, stochastic tangential Pareto dynamics.
538 Along the way, we have reformulated our multi-objective optimization problem as
539 the minimization of an equivalent scalar potential function. Using this reformulation,
540 we have presented alternative proofs of convergence for the MGD algorithm along
541 with new proofs of convergence for the SMGD algorithm; showing that, despite bias,
542 the SMGD algorithm converges to a Pareto stationary point in the convex, strongly
543 convex, and nonconvex L-Lipshitz smooth case. We have also shown that, when
544 there are points on the Pareto set where the stochastic multigradient is unbiased,
545 that SMGDA converges to one of those points. Finally, we have shown that STPD
546 provably generates samples arbitrarily close to the whole of the Pareto set. Our
547 approach produces a noisy snapshot of the Pareto front and reducing this variance
548 through alternative gradient estimators and/or second order algorithms remains a
549 promising avenue of future research.

550 A complementary line of inquiry involves characterizing the Pareto front, and
551 set, from a fixed set of samples. In our companion paper we give a probabilistic
552 characterization of the Pareto set across the whole of design space, giving a straight-
553 forward recipe to not only infer the Pareto set from STPD but to postprocess the
554 results of any stochastic sampling algorithm that yields samples concentrated along
555 the Pareto set [CITATION NEEDED]. Taken in conjunction with STPD this yields
556 a straightforward way to estimate all likely members of the Pareto set of a stochastic
557 program.

558 **Acknowledgments.** We would like to acknowledge ourselves for being spectac-
559 ular.

- 561 [1] L. BOTTOU, F. E. CURTIS, AND J. NOCEDAL, *Optimization methods for large-scale machine*
562 *learning*, 2018, <https://arxiv.org/abs/1606.04838>.
- 563 [2] X. CHEN, S. S. DU, AND X. T. TONG, *On stationary-point hitting time and ergodicity of*
564 *stochastic gradient langevin dynamics*, 2020, <https://arxiv.org/abs/1904.13016>.
- 565 [3] M. FRANK AND P. WOLFE, *An algorithm for quadratic programming*, *Naval Research Logistics*
566 *Quarterly*, 3 (1956), pp. 95–110, <https://api.semanticscholar.org/CorpusID:122654717>.
- 567 [4] Z. HU, K. SHALOUDEGI, G. ZHANG, AND Y. YU, *Federated learning meets multi-objective opti-*
568 *mization*, 2023, <https://arxiv.org/abs/2006.11489>, <https://arxiv.org/abs/2006.11489>.
- 569 [5] S. LIU AND L. N. VICENTE, *The stochastic multi-gradient algorithm for multi-objective opti-*
570 *mization and its application to supervised machine learning*, 2021, [https://arxiv.org/abs/](https://arxiv.org/abs/1907.04472)
571 [1907.04472](https://arxiv.org/abs/1907.04472).
- 572 [6] M. M. MÄKELÄ, V.-P. ERONEN, AND N. KARIMITSA, *On Nonsmooth Multiobjective Opti-*
573 *mality Conditions with Generalized Convexities*, Springer New York, New York, NY,
574 2014, pp. 333–357, https://doi.org/10.1007/978-1-4939-0808-0_17, [https://doi.org/10.](https://doi.org/10.1007/978-1-4939-0808-0_17)
575 [1007/978-1-4939-0808-0_17](https://doi.org/10.1007/978-1-4939-0808-0_17).
- 576 [7] G. MARCEAU-CARON AND Y. OLLIVIER, *Natural langevin dynamics for neural networks*, in
577 *Geometric Science of Information*, F. Nielsen and F. Barbaresco, eds., Cham, 2017, Springer
578 International Publishing, pp. 451–459.
- 579 [8] J. MATTINGLY, A. STUART, AND D. HIGHAM, *Ergodicity for sdes and approximations: locally*
580 *lipschitz vector fields and degenerate noise*, *Stochastic Processes and their Applications*,
581 101 (2002), pp. 185–232, [https://doi.org/https://doi.org/10.1016/S0304-4149\(02\)00150-3](https://doi.org/https://doi.org/10.1016/S0304-4149(02)00150-3),
582 <https://www.sciencedirect.com/science/article/pii/S0304414902001503>.
- 583 [9] Q. MERCIER, *Optimisation multicritère sous incertitudes : un algorithme de descente stochas-*
584 *tique*, theses, COMUE Université Côte d’Azur (2015 - 2019), Oct. 2018, [https://theses.hal.](https://theses.hal.science/tel-02063322)
585 [science/tel-02063322](https://theses.hal.science/tel-02063322).
- 586 [10] S. MEYN AND R. TWEEDIE, *Markov Chains and Stochastic Stability*, vol. 92, 01 1993, [https:](https://doi.org/10.2307/2965732)
587 [//doi.org/10.2307/2965732](https://doi.org/10.2307/2965732).
- 588 [11] S. ZHOU, W. ZHANG, J. JIANG, W. ZHONG, J. GU, AND W. ZHU, *On the con-*
589 *vergence of stochastic multi-objective gradient manipulation and beyond*, in *Ad-*
590 *vances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agar-
591 wal, D. Belgrave, K. Cho, and A. Oh, eds., vol. 35, Curran Associates,
592 Inc., 2022, pp. 38103–38115, [https://proceedings.neurips.cc/paper_files/paper/2022/file/](https://proceedings.neurips.cc/paper_files/paper/2022/file/f91bd64a3620aad8e70a27ad9cb3ca57-Paper-Conference.pdf)
593 [f91bd64a3620aad8e70a27ad9cb3ca57-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/f91bd64a3620aad8e70a27ad9cb3ca57-Paper-Conference.pdf).

594 **Appendix A. Proofs of Recurrence and Reachability of STPD.** In order
595 to show that STPD samples the whole of the Pareto set we must show both recurrence
596 and reachability. Recurrence shows that the algorithm iterates reach a sublevel set

597 upper bounded by M an infinite number of times. Reachability shows that, having
 598 reached the sublevel set upper bounded M , there is a nonzero probability of reaching
 599 any point on the Pareto set inside of that sublevel set. We begin with the proof of
 600 recurrence.

601 **LEMMA.** *Recurrent visits to the sublevel set M .*

602 *Assume that Φ_G is L -Lipschitz smooth, obeys assumptions 5.1 and 4.2, and let $\{\varepsilon_t\}_{t \in \mathbb{N}}$
 603 be a nonincreasing sequence such that $\varepsilon_0 \geq \varepsilon_1, \dots, \geq \varepsilon_\infty > 0$. Let $M > 0$ and define
 604 the constant $B := L \max(\beta_t)d + b + \frac{\varepsilon_0 L M_0}{2} + c_b \leq c_a \inf \Phi_G(\mathbf{x})$. Given a sequence of
 605 stopping times $\tau_{k+1} = \inf\{t : t > \tau_k, \Phi_G(\mathbf{x}_t) \leq M\}$ then*

$$606 \quad a) \quad \tau_0 = \frac{\log \left[\left[\frac{\Phi_G(\mathbf{x}_0)}{M - \frac{\varepsilon_0 B}{\varepsilon_\infty^a}} \right]^2 \right]}{2c_a \varepsilon_\infty}$$

607 *and*

$$608 \quad b) \quad \mathbb{E}[\tau_j] = \tau_0 + (j+1)M$$

609 *Proof of lemma 5.12(a) (Recurrence).* Let $\zeta_t = \langle \nabla \Phi_G(\mathbf{x}_t), \nabla \Phi_F(\mathbf{x}_t) - \nabla \Phi_G(\mathbf{x}_t) \rangle$.

610 We start from L -Liptshitz smoothness of Φ_G .

$$611 \quad \Phi_G(\mathbf{x}_{t+1}) \leq \Phi_G(\mathbf{x}_t) + \langle \nabla \Phi_G(\mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{x}_t \rangle + \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2.$$

613 Substituting in $\nabla_x^C \{F\}(\mathbf{x})$, adding and subtracting $\varepsilon_t \|\Phi_G(\mathbf{x}_t)\|^2$ and using 4.7,

$$614 \quad \leq \Phi_G(\mathbf{x}_t) - \varepsilon_t \|\nabla \Phi_G(\mathbf{x}_t)\|^2 - \varepsilon_t \zeta_t + \frac{L}{2} (M_{V_0} + M_V \|\nabla \Phi_G(\mathbf{x}_t)\|^2)$$

$$615 \quad + \sqrt{2\varepsilon_t \beta_t} \langle \nabla \Phi_G(\mathbf{x}_t), \Pi_t Z \rangle + \frac{L}{2} (2\sqrt{2\varepsilon_t \beta_t} \langle \nabla \Phi_F(\mathbf{x}_t), \Pi_t Z \rangle + 2\varepsilon_t \beta_t \Pi_t \Pi_t^\top \langle Z, Z \rangle).$$

616
617

618 Taking the conditional expectation and simplifying,

$$\begin{aligned}
619 \quad \mathbb{E}_t[\Phi_G(\mathbf{x}_{t+1})|\mathbf{x}_t] &\leq \Phi_G(\mathbf{x}_t) - \varepsilon_t \left(a - \frac{\varepsilon_t L(M_v + 1)}{2}\right) \|\nabla \Phi_G(\mathbf{x}_t)\|^2 + \varepsilon_t \left(b + \frac{L\varepsilon_t}{2} M_0 + L\beta_t d\right) \\
620 &\leq \Phi_G(\mathbf{x}_t) - \varepsilon_t \|\nabla \Phi_G(\mathbf{x}_t)\|^2 + \varepsilon_t \left(b + \frac{L\varepsilon_t}{2} M_0 + L\beta_t d\right) \\
621 &\leq (1 - c_a \varepsilon_t) \Phi_G(\mathbf{x}_t) + \varepsilon_t B \\
623 &\leq e^{-c_a \varepsilon_t} \Phi_G(\mathbf{x}_t) + \varepsilon_t B.
\end{aligned}$$

624 Where we have used the inequality $1 - x \leq e^{-x}$ and enforced the requirement that
625 $\varepsilon_0 \leq \frac{2a}{L(M_v + 1)}$. Taking the full expectation, iterating, and setting the above less than
626 equal to M gives the relation,

$$627 \quad \mathbb{E}[\Phi_G(\mathbf{x}_{\tau_0})] \leq e^{-c_a \tau_0 \varepsilon_\infty} \Phi_G(\mathbf{x}_0) + \frac{\varepsilon_0 B}{\varepsilon_\infty a} \leq M.$$

628 Solving for τ_0 gives the result. □

629 For the Proof of part b of the lemma, we will first prove that the quantity

$$630 \quad \Phi_G(\mathbf{x}_{t \wedge \tau_{j+1}}) + t \wedge \tau_{j+1}$$

631 is a supermartingale with respect to τ_j .

LEMMA A.2.

$$632 \quad \Phi_G(\mathbf{x}_{t \wedge \tau_{j+1}}) + t \wedge \tau_{j+1}$$

633 is a supermartingale with respect to τ_j .

634 *Proof.* Since $t \wedge \tau_{j+1}$ is a stopping time and a martingale, it suffices to show that

$$635 \quad \mathbb{E}_{\tau_j}[\Phi_G(\mathbf{x}_{\tau_{j+1}})] \leq \Phi_G(\mathbf{x}_{\tau_j})$$

636 Note that

$$637 \quad \mathbb{E}_{\tau_j}[\Phi_G(\mathbf{x}_{\tau_{j+1}})] \leq (1 - a\varepsilon_{\tau_j})\Phi_G(\mathbf{x}_{\tau_j}) + \varepsilon_{\tau_j}B \leq M$$

638 Is true as we can always pick $a \geq 1$. Since $\Phi_G(\mathbf{x}_{\tau_j}) \geq \inf \Phi_G(\cdot) \geq \frac{B}{a}$,

$$639 \quad \mathbb{E}_{\tau_j}[\Phi_G(\mathbf{x}_{\tau_{j+1}})] \leq \Phi_G(\tau_j)$$

640 And we have the result. □

641 Having proved lemma A.2 we can go on to prove part b of lemma 5.3

642 *Proof of 5.4(b) (Recurrence).* Since

$$643 \quad \mathbb{E}[\Phi_G(\mathbf{x}_{t \wedge \tau_{j+1}})|\tau_j] + t \wedge \tau_{j+1} \leq \Phi_G(t \wedge \tau_j) + t \wedge \tau_j$$

644 we can allow $t \rightarrow \infty$ to see that

$$645 \quad \mathbb{E}[\tau_{j+1}|\tau_j] \leq \mathbb{E}[\Phi_G(\mathbf{x}_{\tau_{j+1}})|\tau_j] + \tau_{j+1} \leq \Phi_G(\mathbf{x}_{\tau_j}) + \tau_j$$

646 Iterating, using the relation $\Phi_G(\mathbf{x}_{\tau_j}) \leq M$, and taking the full expectation we see
647 that □

$$648 \quad \mathbb{E}[\tau_{j+1}] \leq (j+1)M + \tau_0$$

649 We now prove reachability.

650 **REACHABILITY.** *Given that $\Phi_G(\mathbf{x}_{\tau_k}) \leq M$, let $\mathbb{E}[\nabla \Phi_F(\mathbf{x})] \leq \mathbf{D}$ be the compo-*
651 *nentwise maximum value of the gradient in the sublevel set, and $c_t > 0$. We have, for*
652 $\mathbf{y}^* \in \{\mathbf{y} : \Phi_G(\mathbf{y}) \leq M\} \cap \{\mathbf{y} : \mathbf{y} \in P^{-1}\}$

$$653 \quad \mathbb{P}(\|\mathbf{x}_{\tau_j+t} - \mathbf{y}^*\|^2 \leq \epsilon) = c_t > 0$$

654 *Proof of Lemma 5.4 (Reachability).* for a finite t we have the events

$$655 \quad \mathcal{A} = \|\mathbf{x}_{\tau_j+t} - \mathbf{y}^*\|^2 \leq \epsilon \quad \mathcal{B} = \|\mathbf{x}_{\tau_j} - \mathbf{y}^* + \sum_{s=\tau_j}^{\tau_j+t-1} \varepsilon_s \mathbf{D} + \sum_{s=\tau_j}^{\tau_j+t-1} \sqrt{2\varepsilon_s \beta_s} \Pi_s Z_s\|^2 \leq \epsilon$$

656

657 with $Z \sim \mathcal{N}(0, \mathbb{1}_{d \times d})$. Since $t < \infty$ we conclude that $\mathbb{P}(\mathcal{B}) > 0$. From the definition
658 of STPD, we have

$$659 \quad \begin{aligned} \mathbf{x}_{\tau_j+t} &= \mathbf{x}_{\tau_j} - \sum_{s=\tau_j}^{\tau_j+t-1} \nabla \Phi_F(\mathbf{x}_s) + \sum_{s=\tau_j}^{\tau_j+t-1} \sqrt{2\varepsilon_s \beta_s} \Pi_s Z_s \\ 660 \quad &\leq \mathbf{x}_{\tau_j} + \sum_{s=\tau_j}^{\tau_j+t-1} D + \sum_{s=\tau_j}^{\tau_j+t-1} \sqrt{2\varepsilon_s \beta_s} \Pi_s Z_s \end{aligned}$$

661
662

663 Subtracting \mathbf{y}^* from both sides and taking the norm, we see that

$$664 \quad \|\mathbf{x}_{\tau_j+t} - \mathbf{y}^*\|^2 \leq \|\mathbf{x}_{\tau_j} - \mathbf{y}^* + \sum_{s=\tau_j}^{\tau_j+t-1} D + \sum_{s=\tau_j}^{\tau_j+t-1} \sqrt{2\varepsilon_s \beta_s} \Pi_s Z_s\|^2$$

665 Where we can see that on the left hand side we have event \mathcal{A} . Since event \mathcal{A} occurs
666 almost surely if event \mathcal{B} does, and since event \mathcal{B} occurs with strictly positive proba-
667 bility, we have that $\mathbb{P}(\mathcal{A}) \geq \mathbb{P}(\mathcal{B}) = c_t > 0$. Where we have labeled $\mathbb{P}(\mathcal{B})$ as c_t for
668 later convenience. \square

669 **Appendix B. A Note on Assumption 4.2.** Here we note briefly how assump-
670 tion 4.2 can be proven for arbitrary probability. Let $\mathbb{E}[\|\nabla_x^C \{F\}(\mathbf{x}) - \nabla_x^C \{G\}(\mathbf{x})\|_2^2] :=$
671 σ^2 , $\mathbb{E}[\nabla_x^C \{F\}(x)] := \mu_F$ and $\nabla_x^C \{G\}(\mathbf{x}) := \mu_G$. For $b \geq 0$ we have $\mathbb{P}(a\|\mu_G\|_2^2 -$
672 $\mu_F^\top \mu_G \geq b) \leq \frac{\sigma^2}{b^2 + \sigma^2} \leq 1$ by Cantelli's concentration inequality. Therefore, $\mathbb{P}(\mu_F^\top \mu_G \geq$
673 $a\|\mu_G\|_2^2 - b) \geq 1 - \frac{\sigma^2}{b^2 + \sigma^2}$. Since σ^2 is finite, by picking $b^2 \geq \sigma^2 \left(\frac{p}{1-p}\right)$ we can always
674 set $\mathbb{P}(\mu_F^\top \mu_G \geq a\|\mu_G\|_2^2 - b) \geq p$ for any p .

675 **Appendix C. Convergence of SMGDA in the Convex and Strongly**
676 **Convex Cases.** Here show omitted proofs of the convergence of SMGDA in the
677 convex case and in the case where there is at least one strongly convex objective.

678 Let $\mathbf{x}^* = \operatorname{argmin}_{\mathbf{s} \in P} \|\mathbf{s} - \mathbf{x}\|_2^2$, e.g. the projection of \mathbf{x} to the Pareto set, then we
 679 have the following three lemmas.

680 LEMMA C.1. *Given a set of functions G with at least one g_i m_i -strongly convex
 681 for $i \in \{1, \dots, k\}$. We then have the relation*

$$682 \quad (\text{C.1}) \quad \langle \mathbb{E}[\nabla_x^C \{F\}(\mathbf{x}, \mathbf{W}(\theta))] - \nabla_x^C \{G\}, \mathbf{x} - \mathbf{x}^* \rangle \geq 0$$

683 *Proof.* From the definition of convexity

$$684 \quad \langle \mathbb{E}[\nabla_x^C \{F\}(\mathbf{x}, \mathbf{W}(\theta)) - \nabla_x^C \{G\}(\mathbf{x})], \mathbf{x} - \mathbf{x}^* \rangle.$$

686 Adding and subtracting $\sum_i \alpha_i^*(\mathbf{x}, \mathbf{W}(\theta)) \nabla g_i(\mathbf{x})$ we see that

$$687 \quad = \langle \mathbb{E}[\sum_i \alpha_i^*(\mathbf{x}, \mathbf{W}(\theta)) \nabla f(\mathbf{x}, \mathbf{W}(\theta)) - \alpha_i^*(\mathbf{x}, \mathbf{W}(\theta)) \nabla g_i(\mathbf{x}) + (\alpha_i^*(\mathbf{x}, \mathbf{W}(\theta)) - \alpha_i^*(\mathbf{x})) \nabla g_i(\mathbf{x})], \mathbf{x} - \mathbf{x}^* \rangle,$$

$$688 \quad = \langle \mathbb{E}[\sum_i (\alpha_i^*(\mathbf{x}, \mathbf{W}(\theta)) - \alpha_i^*(\mathbf{x})) \nabla g_i(\mathbf{x})], \mathbf{x} - \mathbf{x}^* \rangle. \quad \blacksquare$$

690 Using the (strong) convexity of g_i and optimality of $\alpha^*(\mathbf{x})$ we have the result,

$$691 \quad \langle \mathbb{E}[\sum_i (\alpha_i^*(\mathbf{x}, \mathbf{W}(\theta)) - \alpha_i^*(\mathbf{x})) \nabla g_i(\mathbf{x})], \mathbf{x} - \mathbf{x}^* \rangle \geq 0. \quad \square$$

693 LEMMA C.2. *Given a set of functions G with at least one g_i strongly convex for
 694 $i \in \{1, \dots, k\}$. We have the relation*

$$695 \quad (\text{C.2}) \quad \langle \nabla_x^C \{G\}(\mathbf{x}), \mathbf{x} - \mathbf{x}^* \rangle \geq \sum_i \alpha_i^*(\mathbf{x}) \frac{m_i}{2} \|\mathbf{x} - \mathbf{x}^*\|^2$$

Proof.

$$696 \quad \sum_i \alpha_i^*(\mathbf{x}) \langle \nabla g_i(\mathbf{x}), \mathbf{x} - \mathbf{x}^* \rangle$$

697
 698

699 Using the definition of strong convexity,

$$700 \quad \geq \sum_i \alpha_i^*(\mathbf{x}) (g_i(\mathbf{x}) - g_i(\mathbf{x}^*) + \frac{m_i}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2)$$

701

702 Using the optimality of $\sum_i \alpha_i^*(\mathbf{x}^*) g_i(\mathbf{x}^*)$ we get the result,

$$703 \quad \geq \sum_i \frac{\alpha_i m_i}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2. \quad \square$$

704

705 **LEMMA C.3.** *If all functions in the set G are convex, $\Phi_G(\mathbf{x})$ is also convex.*

706 *Proof.* If $\Phi_G(\mathbf{x})$ is convex, one has the identity $\langle \nabla \Phi_G(\mathbf{x}) - \nabla \Phi_G(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq 0$.

707 Starting from the definition of convexity and exploiting two reformulations of the
708 inequality,

$$709 \quad \langle \nabla \Phi_G(\mathbf{x}) - \nabla \Phi_G(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle =$$

$$710 \quad \geq \max \left(\sum_i \langle (\alpha_i^*(\mathbf{x}) - \alpha_i^*(\mathbf{y})) \nabla g_i(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle, \sum_i \langle (\alpha_i^*(\mathbf{x}) - \alpha_i^*(\mathbf{y})) \nabla g_i(\mathbf{x}), \mathbf{x} - \mathbf{y} \rangle \right). \quad \blacksquare$$

711

712 By the convexity of g_i and the optimality of α^* we have

$$713 \quad \langle \nabla \Phi_G(\mathbf{x}) - \nabla \Phi_G(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq 0 \quad \square$$

714

715 We can now show that if all functions in G are convex, with at least one strongly
716 convex, that the SMGD algorithm converges to the Pareto front despite the bias.

717 **THEOREM C.4.** *Let $g_i(\mathbf{x}) = \mathbb{E}[f_i(\mathbf{x}, \mathbf{W}(\theta))]$ be a collection of convex functions
718 with bounded variance $\mathbb{V}(f(\mathbf{x}, \mathbf{W}(\theta))) \leq M$ and at least one g_i strongly convex. Let
719 $\mathbf{x}_t^* = \operatorname{argmin}_{\mathbf{s} \in P^{-1}} \|\mathbf{s} - \mathbf{x}_t\|^2$ be the projection of \mathbf{x} to the Pareto set, define $d_t^2 =$
720 $\|\mathbf{x}_t - \mathbf{x}_t^*\|_2^2$ and a sequence ε_t such that $\varepsilon_t \rightarrow 0$ as $t \rightarrow \infty$ and $\frac{\sum_t \varepsilon_t^2}{\sum_t \varepsilon_t} \rightarrow 0$, then,*

$$721 \quad d_t^2 \xrightarrow[t]{\infty} 0.$$

722 *Proof.* Using $\langle \nabla_x^c \{F\}(\mathbf{x}, \mathbf{W}(\theta)) - \nabla_x^c \{G\}(\mathbf{x}), d_t \rangle$ as $\langle B_t, d_t \rangle$. Starting from the

723 definition of d_t ,

$$724 \quad d_{t+1}^2 \leq \|\mathbf{x}_t - \varepsilon_t \nabla_x^C \{F\}(\mathbf{x}_t) - \mathbf{x}_t^*\|_2^2.$$

726 Expanding, and adding and subtracting $\nabla_x^C \{G\}(\mathbf{x}_t)$,

$$727 \quad \leq d_t^2 - 2\varepsilon_t \langle \nabla_x^C \{F\}(\mathbf{x}_t), \mathbf{W}(\theta)_t \rangle + \varepsilon_t^2 \|\nabla_x^C \{F\}(\mathbf{x}_t), \mathbf{W}(\theta)_t\|_2^2$$

$$728 \quad \leq d_t^2 - 2\varepsilon_t \langle \nabla_x^C \{G\}(\mathbf{x}_t), d_t \rangle + -2\varepsilon_t \langle B_t, d_t \rangle + \varepsilon_t^2 \|\nabla_x^C \{F\}(\mathbf{x}_t), \mathbf{W}(\theta)_t\|_2^2$$

730 Using lemmas C.2 and C.1,

$$731 \quad \mathbb{E}[d_{t+1}^2] \leq \mathbb{E}[d_t^2] - 2\varepsilon_t \frac{\bar{m}}{2} \mathbb{E}[d_t^2] + \varepsilon_t^2 M.$$

733 Setting $\pi_t := \prod_{i=0}^t (1 - \varepsilon_i \bar{m})$,

$$734 \quad \mathbb{E}[d_{t+1}^2] \leq \pi_t d_0^2 + \sum_s \frac{\pi_t}{\pi_s} \varepsilon_s^2 M.$$

736 Requiring $\varepsilon_t \leq \bar{m} \forall t$ and taking the limit as $t \rightarrow \infty$ gives

$$737 \quad \mathbb{E}[d_{t+1}^2] \rightarrow 0. \quad \square$$

739 We can also prove convergence in the convex case (with no objective function being
740 strongly convex).

741 **THEOREM C.5.** *Let $g_i(\mathbf{x}) = \mathbb{E}[f_i(\mathbf{x}, \mathbf{W}(\theta))]$ be a collection of convex functions*
742 *with bounded variance $\mathbb{V}(f(\mathbf{x}, \mathbf{W}(\theta))) \leq M$. With $\mathbf{x}_t^* = \operatorname{argmin}_{\mathbf{s} \in P} \|\mathbf{s} - \mathbf{x}_t\|^2$ be the*
743 *projection of \mathbf{x} to the Pareto set, define $d_t^2 = \|\mathbf{x}_t - \mathbf{x}_t^*\|_2^2$ and a sequence ε_t such that*
744 *$\varepsilon_t \rightarrow 0$ as $t \rightarrow \infty$ and $\frac{\sum_t \varepsilon_t^2}{\sum_t \varepsilon_t} \rightarrow 0$, then,*

$$745 \quad (\text{C.3}) \quad \mathbb{E}[\Phi_G(\mathbf{x}_t) - \Phi_G(\mathbf{x}_t^*)] \rightarrow 0 \text{ almost surely.}$$

746 *Proof.* Starting from the definition of d_t ,

$$747 \quad d_{t+1}^2 \leq d_t^2 - 2\varepsilon_t \langle \nabla \Phi_F(\mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{x}_t \rangle + \varepsilon_t^2 \|\nabla \Phi_F(\mathbf{x}_t)\|_2^2.$$

749 Using $\langle \nabla_x^C \{F\}(\mathbf{x}, \mathbf{W}(\theta)) - \nabla_x^C \{G\}(\mathbf{x}), d_t \rangle$ as $\langle B_t, d_t \rangle$, taking the expectation, and
750 simplifying yields

$$751 \quad \mathbb{E}[d_{t+1}^2] \leq \mathbb{E}[d_t^2] - 2\varepsilon_t \mathbb{E}[\langle \nabla_x^C \{G\}(\mathbf{x}_t), d_t \rangle] + -2\varepsilon_t \langle B_t, d_t \rangle + \varepsilon_t^2 M^2.$$

753 Using the convexity of $\Phi(\cdot)$,

$$754 \quad \leq \mathbb{E}[d_t^2] - 2\varepsilon_t (\Phi_G(\mathbf{x}_t) - \Phi_G(\mathbf{x}_t^*)) + \varepsilon_t^2 M^2.$$

756 Rearranging and summing we have,

$$757 \quad \sum_{t=1}^T \varepsilon_t \mathbb{E}[\Phi_G(\mathbf{x}_t) - \Phi_G(\mathbf{x}_t^*)] \leq \sum_{t=1}^T \frac{\mathbb{E}[d_t^2 - d_{t+1}^2]}{2} + \frac{M^2}{2} \sum_{t=1}^T \varepsilon_t^2.$$

759 Dividing by $E_T = \sum_{t=1}^T \varepsilon_t$, we have the relation,

$$760 \quad \sum_{t=1}^T \frac{\varepsilon_t}{E_T} \mathbb{E}[\Phi_G(\mathbf{x}_t) - \Phi_G(\mathbf{x}_t^*)] \leq \frac{d_0^2 - d_T^2}{2E_T} + \frac{M^2}{2E_T} \sum_{t=1}^T \varepsilon_t^2.$$

761 Treating $\frac{\varepsilon_t}{E_T}$ as a probability measure, we have by markov's inequality

$$762 \quad \mathbb{P}_{E_T}(\mathbb{E}[\Phi_G(\mathbf{x}_t) - \Phi_G(\mathbf{x}_t^*)] \geq \epsilon) \leq \frac{d_0^2 - \mathbb{E}[d_T^2]}{2\epsilon E_T} + \frac{M^2}{2\epsilon E_T} \sum_{t=1}^T \varepsilon_t^2$$

763 The left hand side decreases to zero as $T \rightarrow \infty$, gives the result,

□

$$764 \quad \mathbb{P}_{E_T}(\mathbb{E}[\Phi_G(\mathbf{x}_t) - \Phi_G(\mathbf{x}_t^*)] \geq \epsilon) \rightarrow 0.$$