

FROM CONFLICTS TO CONVERGENCE: A ZERO-ORDER METHOD FOR MULTI-OBJECTIVE LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Multi-objective learning (MOL) is a popular paradigm for learning problems under multiple criteria, where various dynamic weighting algorithms (e.g., MGDA and MODO) have been formulated to find an updated direction for avoiding conflicts among objectives. Recently, increasing endeavors have struggled to tackle the black-box MOL when the gradient information of objectives is unavailable or difficult to attain. Albeit the impressive success of zeroth-order method for single-objective black-box learning, the corresponding MOL algorithm and theoretical understanding are largely absent. Unlike single-objective problems, the errors of MOL introduced by zeroth-order gradients can simultaneously affect both the gradient estimation and the gradient coefficients λ , leading to further error amplification. To address this issue, we propose a Stochastic Zeroth-order Multiple Objective Descent algorithm (SZMOD), which leverages function evaluations to approximate gradients and develops a new decomposition strategy to handle the complicated black-box multi-objective optimization. Theoretically, we provide convergence and generalization guarantees for SZMOD in both general non-convex and strongly convex settings. Our results demonstrate that the proposed SZMOD enjoys a promising generalization bound of $\mathcal{O}(n^{-\frac{1}{2}})$, which is comparable to the existing results of first-order methods requiring additional gradient information. Experimental results validate our theoretical analysis.

1 INTRODUCTION

Multi-objective learning (MOL) aims to learn a single model that can optimize multiple potentially conflicting objectives simultaneously. An unconstrained multi-objective optimization problem can be defined as

$$\min_{x \in \mathbb{R}^d} F_S(x) := [f_{S,1}(x), \dots, f_{S,M}(x)], \quad (1)$$

where $S = \{z_i\}_{i=1}^n$ is the training dataset, $f_{S,m}(x)$ is the m -th empirical objective for $m \in [M] =: \{1, 2, \dots, M\}$. Usually, we can set $f_{S,m}(x) = \sum_{i=1}^n f_{z_i,m}(x)$ as the empirical risk on the entire training dataset S , where $f_{z,m} : \mathbb{R}^d \mapsto \mathbb{R}$ measures the performance of a model $x \in \mathbb{R}^d$ on a datum z for the m -th objective.

Multi-objective learning has gained increasing attention, due to the complex decision-making processes involved in many challenging tasks, e.g., managing traffic systems (Felten et al., 2024), electricity grids (Lu et al., 2022), and taxation policy design (Zheng et al., 2022). These burgeoning fields in practice, which require trading off multiple conflict objectives, underscore the significance of research in MOL. Specifically, balancing bias and variance (Neal et al., 2018), or accuracy and calibration (Guo et al., 2017), are well-known common objectives in machine learning that need to be optimized. To tackle these problems, this paper pays particular attention to multi-objective gradient methods that aim to find a common descent direction for all objectives. Désidéri (2012) initially introduced the concept of a Pareto stationary and the multi-gradient descent (MGDA) algorithm. Since then, stochastic variants such as MOCO (Fernando et al., 2023) and MODO (Chen et al., 2024) have been proposed. Those first-order multi-objective algorithms have great performed in the white-box problem.

However, when we consider the black box problem, where obtaining explicit gradients is either unattainable or too expensive, these algorithms are no longer applicable. For instance, in the field

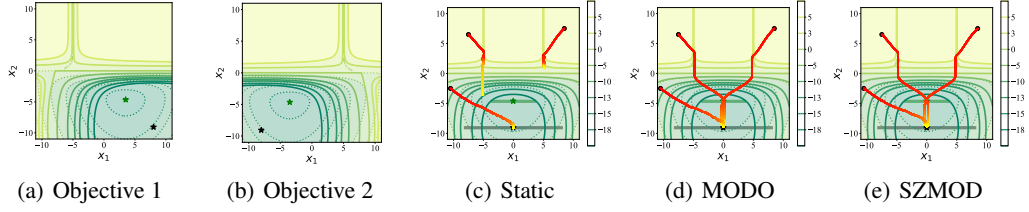


Figure 1: An example from (Liu et al., 2021) involves two objectives in Figure 1(a) and 1(b) to demonstrate the conflict between objectives. Figures 1(c)-1(e) show the optimization trajectories, where the black dots indicate the initialization points of the trajectories, with the colors transitioning from red (start) to yellow (end). The background solid/dotted contours represent the landscape of the average empirical and population objectives, respectively. The gray/green bars mark the empirical/population Pareto fronts, while the black \star green \star marks the solution to the average objectives.

of multiple-objective reinforcement learning (Hu et al., 2023; Felten et al., 2024; Terry et al., 2021; Gupta et al., 2017), agents often can only learn strategies through interaction and external reward signals, without access to the internal state or dynamics of the environment. Similarly, in most attack scenarios (Akhtar & Mian, 2018; Liu et al., 2022; Papernot et al., 2017; 2016), the attacker’s knowledge of the classifier is very limited, which causes the attacker only to execute a black-box attack. Liang et al. (2022) state that the black-box attacks can manipulate model outputs by adjusting the trade-offs between true and false positives without direct access to the model’s internals. Williams & Li (2023) consider a novel multi-objective sparse attack that can simultaneously reduce the number and the individual size of modified pixels during the attack process.

Most of the black-box MOL scenarios discussed above are traditionally optimized using the hypervolume indicator (Felten et al., 2024) as the standard performance metric and are typically solved using methods such as evolutionary algorithms (Zhou et al., 2024; Mathai et al., 2020; Liu et al., 2024). Unfortunately, these methods impose strict constraints on problem dimensionality. In contrast, zeroth-order (ZO) optimization algorithms demonstrate greater versatility in handling higher-dimensional problems and can achieve excellent performance, often comparable to or even surpassing that of white-box models where gradients are explicitly available. (Sun et al., 2022; Papernot et al., 2017). Unfortunately, there has been no endeavor to apply the zeroth-order optimization to multi-objective optimization.

To fill this gap, we present the Stochastic Zeroth-order Multiple Objective Descent algorithm (SZMOD), which integrates coordinate-based zeroth-order gradient estimations and employs a consistent directional selection strategy during the λ iteration process. Specifically, by using the same direction for gradient approximation throughout the iterations, SZMOD ensures that the update direction of the dynamic weigh λ_t is updated in alignment with the chosen direction, thereby maintaining stability and reducing variance in the optimization process. Combining coordinate zeroth-order techniques and unified directional updates enhances the algorithm’s ability to effectively address black-box multi-objective learning problems.

- **Gradient Direction Conflict:** In first-order multi-objective optimization algorithms, the gradients of multiple objective functions are computed to determine a suitable direction for optimization. However, in zeroth-order multi-objective problems, we rely on zeroth-order gradient estimates, where the direction estimation depends entirely on a random vector u (determined by the zeroth-order estimation process). This dependence makes it challenging to identify an appropriate CA direction (the proper direction to update λ , will be defined in section 2.4), complicating the optimization process.
- **Excessive Error Risk:** Zeroth-order gradient estimation inherently introduces errors, which also propagate into the iterative updates of λ . These compounded errors affect the term of the CA direction, increasing the risk of divergence during the iteration of x . Therefore, it is crucial to control these errors effectively to ensure convergence and maintain the stability of the optimization process.

2 PRELIMINARIES

In this section, we first introduce MOL’s problem formulation, the analysis target, and the metric to measure its optimization, generalization, and CA direction.

2.1 NOTATION

Denote the vector-valued objective function on datum z as $F_z(x) = [f_{z,1}(x), \dots, f_{z,M}(x)]$. The training and testing performance of x can then be measured by the empirical objective $F_S(x)$ and the population objective $F(x)$ which are, respectively, defined as $F_S(x) := \frac{1}{n} \sum_{i=1}^n F_{z_i}(x)$ and $F(x) := \mathbb{E}_{z \sim \mathcal{D}} [F_z(x)]$. Their corresponding gradients are denoted as $\nabla F_S(x)$ and $\nabla F(x) \in \mathbb{R}^{d \times M}$.

2.2 METHOD OF MOL

Analogous to the stationary solution and optimal solution in single-objective learning, we define the Pareto stationary point and Pareto optimal solution for MOL problem $\min_{x \in \mathbb{R}^d} F(x)$ as follows.

Definition 1 (Pareto stationary and Pareto optimal). *If there exists a convex combination of the gradient vectors that equals to zero, i.e., there exists $\lambda \in \Delta^M$ such that $\nabla F(x)\lambda = 0$, then $x \in \mathbb{R}^d$ is Pareto stationary. If there is no $x \in \mathbb{R}^d$ and $x \neq x^*$ such that, for all $m \in [M]$, $f_m(x) \leq f_m(x^*)$, with $f_{m'}(x) < f_{m'}(x^*)$ for at least one $m' \in [M]$, then x^* is Pareto optimal. If there is no $x \in \mathbb{R}^d$ such that for all $m \in [M]$, $f_m(x) < f_m(x^*)$, then x^* is weakly Pareto optimal.*

By definition, at a Pareto stationary solution, there is no common descent direction for all objectives. A necessary and sufficient condition for x being Pareto stationary for smooth objectives is that $\min_{\lambda \in \Delta^M} \|\nabla F(x)\lambda\| = 0$. Therefore, $\min_{\lambda \in \Delta^M} \|\nabla F(x)\lambda\|$ can be used as a measure of Pareto stationarity (PS). We will refer to the aforementioned quantity as the PS population risk henceforth and its empirical version as PS empirical risk or PS optimization error. We next introduce the target of our analysis based on the above definitions.

2.3 ZERO-ORDER GRADIENT ESTIMATION

Coordinate-wise Gradient Estimation. When only function evaluations are available, here, we employ the deterministic coordinate-wise direction to derive the decent direction. Specifically, for the smoothing constant v and vector u_i (u_i represents the unit vector where the i -th element is 1 and the remaining elements are 0), the directional derivative of $f_{z,m}$ in the direction u for the smooth function $f_i, i \in [n]$, can be estimated as:

$$\hat{\nabla} f_{z,m}(x, u, v) = \sum_{j=1}^d \frac{f_{z,m}(x + vu_j) - f_{z,m}(x)}{v} u_j. \quad (2)$$

as the approximation of the full directional gradient. Since the smoothing constant v is fixed, for simplicity, we leave out v in these gradient estimations and set

$$\hat{\nabla} f_{z,m}(x, u) := \hat{\nabla} f_{z,m}(x, u, v). \quad (3)$$

Denote the vector-valued objective function on datum z as $F_z(x) = [f_{z,1}(x), \dots, f_{z,M}(x)]$. The training and testing performance of x can then be measured by the empirical objective $F_S(x)$ and the population objective $F(x)$ which are, respectively, defined as $F_S(x) := \frac{1}{n} \sum_{i=1}^n F_{z_i}(x)$ and $F(x) := \mathbb{E}_{z \sim \mathcal{D}} [F_z(x)]$. Their corresponding estimate gradients are denoted as $\hat{\nabla} F_S(x)$ and $\hat{\nabla} F(x) \in \mathbb{R}^{d \times M}$. Thus the zeroth-order estimate for all objectives on datum z should be written as $\hat{\nabla} F_z(x) = [\hat{\nabla} f_{z,1}(x), \dots, \hat{\nabla} f_{z,M}(x)]$.

2.4 PROBLEM SETUP

Proposition 1 ((Tanabe et al., 2019) Lemma 2.2). *If $f_m(x)$ are convex or strongly convex for all $m \in [M]$, and $x \in \mathbb{R}^d$ is a Pareto stationary point of $F(x)$, then x is weakly Pareto optimal or Pareto optimal.*

Next, we proceed to decompose the PS population risk.

Error Decomposition. Given a model x , the PS population risk can be decomposed into

$$\underbrace{\min_{\lambda \in \Delta^M} \|\nabla F(x)\lambda\|}_{\text{PS population risk } R_{\text{pop}}(x)} = \underbrace{\min_{\lambda \in \Delta^M} \|\nabla F(x)\lambda\| - \min_{\lambda \in \Delta^M} \|\nabla F_S(x)\lambda\|}_{\text{PS generalization error } R_{\text{gen}}(x)} + \underbrace{\min_{\lambda \in \Delta^M} \|\nabla F_S(x)\lambda\|}_{\text{PS optimization error } R_{\text{opt}}(x)}, \quad (4)$$

where the optimization error quantifies the training performance, i.e., how well does model x perform on the training data; and the generalization error (gap) quantifies the difference between the testing performance on new data sampled from \mathcal{D} and the training performance, i.e., how well the model x performs on unseen testing data compared to the training data.

The zeroth-order optimization is a gradient-based black-box optimization that utilizes the difference information of function values to approximate the true gradient. Furthermore, this method does not alter the optimization objective, only the optimization process differs from the first-order one. As for MOL black-box problems, the optimization objective of the SZMOD remains $\min_{\lambda \in \Delta^M} \|\nabla F(x)\lambda\| = 0$.

Let $A : \mathcal{Z}^n \mapsto \mathbb{R}^d$ denote a randomized MOL algorithm. Given training data S , we are interested in the expected performance of the output model $x = A(S)$, which is measured by $\mathbb{E}_{A,S} [R_{\text{pop}}(A(S))]$. From equation 4 and linearity of expectation, it holds that

$$\mathbb{E}_{A,S} [R_{\text{pop}}(A(S))] = \mathbb{E}_{A,S} [R_{\text{gen}}(A(S))] + \mathbb{E}_{A,S} [R_{\text{opt}}(A(S))]. \quad (5)$$

Distance to CA direction. Consider an update direction $d = -\nabla F_S(x)\lambda$, where λ is the dynamic weights from a simplex $\lambda \in \Delta^M := \{\lambda \in \mathbb{R}^M \mid \mathbf{1}^\top \lambda = 1, \lambda \geq 0\}$. To obtain such a steepest CA direction in unconstrained learning that maximizes the minimum descent of all objectives, we can solve the following problem (Fliege et al., 2019)

$$\text{CA direction } d(x) = \arg \min_{d \in \mathbb{R}^d} \max_{m \in [M]} \left\{ \langle \nabla f_{S,m}(x), d \rangle + \frac{1}{2} \|d\|^2 \right\} \quad (6)$$

$$\stackrel{\text{equivalent to}}{\iff} d(x) = -\nabla F_S(x)\lambda^*(x) \text{ s.t. } \lambda^*(x) \in \arg \min_{\lambda \in \Delta^M} \|\nabla F_S(x)\lambda\|^2. \quad (7)$$

Defining $d_\lambda(x) = -\nabla F_S(x)\lambda$ given $x \in \mathbb{R}^d$ and $\lambda \in \Delta^M$, we measure the distance to $d(x)$ via (Fernando et al., 2023)

$$\text{CA direction error } \mathcal{E}_{\text{ca}}(x, \lambda) := \|d_\lambda(x) - d(x)\|^2. \quad (8)$$

With the above definitions of measures that quantify the performance of algorithms in different aspects, we then introduce a stochastic gradient algorithm for MOL that is analyzed in this work.

3 A STOCHASTIC ALGORITHM FOR BLACK-BOX MOL

In this section, we first introduce our main algorithm, Stochastic Zeroth-order Multiple Objective Descent (SZMOD).

At each iteration t , α_t, γ_t are step sizes, and $\Pi_{\Delta^M}(\cdot)$ denotes Euclidean projection to the simplex Δ^M . Denoting $z_{t,s}$ as an independent sample from S with $s \in [3]$, and $\hat{\nabla} F_{z_{t,s}}(x_t)$ as the gradient estimate of $\nabla F_{z_{t,s}}(x_t)$.

Remark 1. In the iteration process of λ_t , gradient direction conflicts prevent us from achieving convergence. To ensure the algorithm converges, SZMOD requires that $\hat{\nabla} f_{z,1}(x)$ and $\hat{\nabla} f_{z,2}(x)$ use the same stochastic direction. By this method, we have

$$\mathbb{E}_{z_{t,1}, z_{t,2}} \left[\hat{\nabla} F_{z_{t,1}}(x_t)^\top \hat{\nabla} F_{z_{t,2}}(x_t) \lambda_t \right] = \nabla F_S(x_t)^\top \nabla F_S(x_t) \lambda_t + \mathcal{O}(v),$$

which means that we can stabilize the updates and control the error through v .

Algorithm 1 Stochastic Zeroth-order Multiple Objective Descent (SZMOD)

Input: Training data S , initial model x_0 , weighting co-efficient λ_0 , and their learning rates $\{\alpha_t\}_{t=0}^T, \{\gamma_t\}_{t=0}^T$.

Output: x_T

```

1: for  $t = 0, \dots, T - 1$  do
2:   for  $m = 1, \dots, M$  do
3:     Compute zeroth-order gradients  $\hat{\nabla} f_{m,z_t,s}(x_t)$  using same  $u, s \in [2]$ 
4:     Compute zeroth-order gradients  $\hat{\nabla} f_{m,z_t,3}(x_t)$  with coordinate
5:   end for
6:   Compute dynamic weight  $\lambda_{t+1}$  following
7:   Compute  $\lambda_{t+1} = \Pi_{\Delta^M} \left( \lambda_t - \gamma_t \hat{\nabla} F_{z_t,1}(x_t)^\top \hat{\nabla} F_{z_t,2}(x_t) \lambda_t \right)$ 
8:   Compute  $x_{t+1} = x_t - \alpha_t \hat{\nabla} F_{z_t,3}(x_t) \lambda_{t+1}$ 
9: end for

```

In the iteration process of x_t , the zeroth-order method will also lead to excessive error risk, which is caused by the error of λ_{t+1} and $\hat{\nabla} F_{z,3}$. The error of λ_{t+1} can be control by remark 1. Here, we choose to use the coordinate zeroth-order estimate to minimize the error of $\hat{\nabla} F_{z,3}$.

4 OPTIMIZATION OF SZMOD

In this section, we bound the multi-objective PS optimization error $\min_{\lambda \in \Delta^M} \|\nabla F_S(x) \lambda\|$ (Fernando et al., 2023; Fliege et al., 2019; Désidéri, 2012). As discussed in Section 2.2, this measure being zero implies the model x achieves a Pareto stationarity for the empirical problem.

Below, we list the standard assumptions used to derive the optimization error, which has been widely used for theoretical analysis for (Chen et al., 2024; Lei, 2023; Fliege et al., 2019).

Assumption 1 (Lipschitz continuity of $F_z(x)$). *For all $m \in [M]$, $f_{z,m}(x)$ are ℓ_f -Lipschitz continuous for all z . Then $F_z(x)$ are ℓ_F -Lipschitz continuous in Frobenius norm for all z with $\ell_F = \sqrt{M} \ell_f$.*

Assumption 2 (Lipschitz continuity of $\nabla F_z(x)$). *For all $m \in [M]$, $\nabla f_{z,m}(x)$ is $\ell_{f,1}$ -Lipschitz continuous for all z . And $\nabla F_z(x)$ is $\ell_{F,1}$ -Lipschitz continuous in Frobenius norm for all z .*

Assumption 3. *For all $m \in [M]$, $z \in \mathcal{Z}$, $f_{z,m}(x)$ is μ -strongly convex w.r.t. x with $\mu > 0$.*

Note that in the strongly convex case, the gradient norm $\|\nabla F_z(x)\|_F$ can be unbounded in \mathbb{R}^d . Therefore, one cannot assume Lipschitz continuity of $f_{z,m}(x)$ w.r.t. $x \in \mathbb{R}^d$. We address this challenge by showing that $\{x_t\}$ generated by the SZMOD algorithm is bounded as stated in Lemma 1. Notably, combined with Assumption 1, we can derive that the gradient norm $\|\nabla F_z(x_t)\|_F$ is also bounded.

Lemma 1 (Boundedness of x_t for strongly convex and smooth objectives). *Suppose Assumptions 2, 3 hold. For $\{x_t\}, t \in [T]$ generated by SZMOD algorithm or other dynamic weighting algorithm with weight $\lambda \in \Delta^M$, step size $\alpha_t = \alpha$, and $0 \leq \alpha \leq \ell_{f,1}^{-1}$, there exists a finite positive constant c_x such that $\|x_t\| \leq c_x$. And there exists finite positive constants $\ell_f, \ell_F = \sqrt{M} \ell_f$, such that for all $\lambda \in \Delta^M$, we have $\|\nabla F(x_t) \lambda\| \leq \ell_f, \|\nabla F(x_t)\|_F \leq \ell_F$.*

4.1 DISTANCE TO CA DIRECTION

Theorem 1 (Distance to CA direction). *Suppose either: 1) Assumptions 1, 3 hold; or 2) Assumptions 1, 2 hold, with ℓ_f and ℓ_F defined in Lemma 1. Consider $\{x_t\}, \{\lambda_t\}$ generated by the SZMOD algorithm. For all $\lambda \in \Delta^M$, it holds that:*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left[\|d_{\lambda_t}(x_t) - d(x_t)\|^2 \right] \leq \frac{4}{\gamma T} + 6 \sqrt{M \ell_{f,1} \ell_f^2 \frac{\alpha}{\gamma}} + \gamma M \ell_f^4 + e \quad (9)$$

Here $e = \frac{\ell_{f,1}^2 v^2 d}{4} \mathbb{E}_A \|\lambda_t - \lambda\|_1 + \frac{\ell_{f,1} v}{2} \mathbb{E}_A (\|\lambda_t - \lambda\|_1 \|\nabla F_S \lambda\|_1 + d \|\nabla F_S(\lambda_t - \lambda)\|_1)$ caused by zeroth-order error. We should mention that e can be seen as $\mathcal{O}(v)$. Analyzing convergence to

the CA direction using the measure introduced in Section 2.4. By, e.g., choosing $\alpha = \Theta\left(T^{-\frac{3}{4}}\right)$, $\gamma = \Theta\left(T^{-\frac{1}{4}}\right)$ and $v = \gamma/10$, the RHS of equation 9 converges in a rate of $\mathcal{O}\left(T^{-\frac{1}{4}}\right)$.

4.2 PS OPTIMIZATION ERROR

Theorem 2. (PS optimization error of SZMOD). Suppose either 1) Assumptions 1, 3 hold or 2) Assumptions 1, 2 hold, with ℓ_f defined in Lemma 1. Define c_F such that $\mathbb{E}_A[F_S(x_0)\lambda_0] - \min_{x \in \mathbb{R}^d} \mathbb{E}_A[F_S(x)\lambda_0] \leq c_F$. Considering $\{x_t\}$ generated by SZMOD (Algorithm 1), with $\alpha_t = \alpha \leq 1/(2\ell_{f,1})$, $\gamma_t = \gamma$, then under either condition 1) or 2), it holds that

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left[\min_{\lambda \in \Delta^M} \|\nabla F_S(x_t) \lambda\| \right] \leq \sqrt{\frac{c_F}{\alpha T}} + \sqrt{\frac{3}{2} \gamma M \ell_f^4} + \sqrt{\frac{1}{2} \alpha \ell_{f,1} \ell_{f,d}^2} + e. \quad (10)$$

The choice of step sizes $\alpha = \Theta(T^{-\frac{3}{4}})$, $\gamma = \Theta(T^{-\frac{1}{4}})$, and smoothing constant $v = \gamma/10$ to ensure convergence to CA direction is suboptimal for the convergence to Pareto stationarity. Then the RHS of equation 10 converges in a rate of $\mathcal{O}\left(T^{-\frac{1}{8}}\right)$.

5 GENERALIZATION OF SZMOD

In the following, we provide uniform stability for the black-box MOL algorithm, whose expected PS generalization error can be further bounded under several convexity scenarios.

Proposition 2 ((Chen et al., 2024), Proposition 2). With $\|\cdot\|_F$ denoting the Frobenious norm, $R_{\text{gen}}(A(S))$ in (2.2) can be bounded by

$$\mathbb{E}_{A,S} [R_{\text{gen}}(A(S))] \leq \mathbb{E}_{A,S} [\|\nabla F(A(S)) - \nabla F_S(A(S))\|_F]. \quad (11)$$

With Proposition 2, we introduce the concept of MOL uniform stability tailored for MOL problems. Then, we analyze their bounds in the general nonconvex and strongly convex cases, respectively.

Definition 2 (MOL uniform stability). A randomized algorithm $A : \mathcal{Z}^n \mapsto \mathbb{R}^d$, is MOL-uniformly stable with ϵ_F iff for all neighboring datasets S, S' that differ in at most one sample, we have

$$\sup_z \mathbb{E}_A \left[\|\nabla F_z(A(S)) - \nabla F_z(A(S'))\|_F^2 \right] \leq \epsilon_F^2.$$

Next, we show the relation between the upper bound of PS generalization error in 4 and MOL uniform stability in Proposition 3.

Proposition 3 ((Chen et al., 2024), proposition 3). Assume for any z , the function $F_z(x)$ is differentiable. If a randomized algorithm $A : \mathcal{Z}^n \mapsto \mathbb{R}^d$ is MOL-uniformly stable with ϵ_F , then

$$\mathbb{E}_{A,S} [\|\nabla F(A(S)) - \nabla F_S(A(S))\|_F] \leq 4\epsilon_F + \sqrt{n^{-1} \mathbb{E}_S [\mathbb{V}_{z \sim \mathcal{D}} (\nabla F_z(A(S)))]}. \quad (12)$$

where $\mathbb{V}_{z \sim \mathcal{D}} (\nabla F_z(A(S))) = \mathbb{E}_{z \sim \mathcal{D}} [\|\nabla F_z(A(S)) - \mathbb{E}_{z \sim \mathcal{D}} [\nabla F_z(A(S))]\|_F^2]$ is the variacne.

Proposition 3 establishes a connection between the upper bound of the PS generalization error and the MOL uniform stability.

Theorem 3 (PS generalization error of SZMOD in nonconvex case). If $\sup_z \mathbb{E}_A [\|\nabla F_z(A(S))\|_F^2] \leq G^2$ for any S , then the MOL uniform stability, i.e., ϵ_F^2 in Definition 2 is bounded by $\epsilon_F^2 \leq 4G^2T/n$. And the PS generalization error $\mathbb{E}_{A,S} [R_{\text{gen}}(A(S))] = \mathcal{O}\left(T^{\frac{1}{2}} n^{-\frac{1}{2}}\right)$.

Remark 2. The proof process of non-convex generalization does not involve parameter updates. Therefore, zeroth-order gradient approximation does not affect the generalization results. At this point, the generalization results of the first-order and zeroth-order methods are naturally the same.

With Lemma 1 and Lemma 6, the stability bound and PS generalization is provided below.

Theorem 4 (PS generalization error of in strongly convex case). *Suppose Assumptions 2 and 3 hold. Let A be the SZMOD algorithm (Algorithm 1). For the MOL uniform stability ϵ_F of algorithm A in Definition 2, if the step sizes satisfy $0 < \alpha_t \leq \alpha \leq 1/(2\ell_{f,1})$, $0 < \gamma_t \leq \gamma \leq \min \left\{ \frac{\mu^2}{484\ell_{f,d}^2\ell_{g,1}}, \frac{1}{8(3\ell_{f,d}^2+2\ell_{g,1})} \right\} / T$, and smooth constant $v \leq \min \left\{ \frac{1}{nd}, \frac{1}{nd(2\ell_{g,1}+\ell_{g,1}^2)} \right\}$ then it holds that*

$$\epsilon_F^2 \leq \frac{48}{\mu n} \ell_{f,d}^2 \ell_{F,1}^2 \left(\alpha + \frac{12 + 4M\ell_{f,d}^2}{\mu n} + \frac{10M\ell_f^4 \gamma}{\mu} \right) + \frac{4}{\mu n} \ell_{F,1}^2 \left(\frac{10\alpha M\ell_{f,d}^2 \gamma + \mu \gamma}{\mu \alpha} + \alpha \ell_{f,1} + \frac{2\alpha \ell_{f,1}^2}{n} \right). \quad (13)$$

and $\mathbb{E}_{A,S} [R_{\text{gen}}(A(S))] = \mathcal{O}(n^{-\frac{1}{2}})$.

Remark 3. *Theorem 3, 4 implies setting proper step sizes for different convexity helps to improve the generalization. Under strong convexity conditions, the proof process involving parameter updates will inevitably introduce the cumulative error brought by zeroth-order estimation. We must constrain the smoothness parameter v to achieve the same generalization convergence rate as the first-order method.*

6 CONNECTION BETWEEN OPTIMIZATION, CONFLICT AVOIDANCE AND GENERALIZATION

In this section, we combine the proof process and theoretical results on optimization error, generalization bounds, and the distance to the CA direction to discuss the impact of introducing zeroth-order gradient approximations on multi-objective algorithms. Summarizing the findings from Sections 4 and 5, we derive the PS population risk. With $A_t(S) = x_t$ denoting the output of algorithm A at the t -th iteration, we can decompose the PS population risk $R_{\text{pop}}(A_t(S))$ as (cf. equation 4, equation 11)

$$\mathbb{E}_{A,S} [R_{\text{pop}}(A_t(S))] \leq \mathbb{E}_{A,S} \left[\min_{\lambda \in \Delta^M} \|\nabla F_S(A_t(S)) \lambda\| \right] + \mathbb{E}_{A,S} [\|\nabla F(A_t(S)) - \nabla F_S(A_t(S))\|_F]$$

Theorem 5 (The general nonconvex case). *Suppose Assumptions 1, 2 hold. By the optimization error in Theorem 2 and the generalization error bound in Theorem 3, the PS population risk of the output of SZMOD can be bounded by*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{A,S} [R_{\text{pop}}(A_t(S))] = \mathcal{O} \left(\alpha^{-\frac{1}{2}} T^{-\frac{1}{2}} + \alpha^{\frac{1}{2}} + \gamma^{\frac{1}{2}} + T^{\frac{1}{2}} n^{-\frac{1}{2}} \right) + \mathcal{O}(v).$$

Remark 4. *By selecting step sizes of $\alpha = \Theta(T^{-\frac{1}{2}})$ and $\gamma = \Theta(T^{-\frac{1}{2}})$, with the number of steps $T = \Theta(n^{\frac{2}{3}})$, we can choose a smoothing parameter of $v = \Theta(n^{-\frac{1}{6}})$, which effectively limits the impact of the zeroth-order approximation on optimization convergence. Under these conditions, the expected PS population risk is $\mathcal{O}(n^{-\frac{1}{6}})$.*

Theorem 6 (The strongly convex case). *Suppose Assumptions 2, 3 hold. By the optimization error and the generalization error given in Theorems 2 and 4, SZMOD's PS population risk can be bounded by*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{A,S} [R_{\text{pop}}(A_t(S))] = \mathcal{O} \left(\alpha^{-\frac{1}{2}} T^{-\frac{1}{2}} + \alpha^{\frac{1}{2}} + \gamma^{\frac{1}{2}} + n^{-\frac{1}{2}} \right) + \mathcal{O}(v).$$

Remark 5. *Choosing step sizes $\alpha = \Theta(T^{-\frac{1}{2}})$, $\gamma = o(T^{-1})$. Under strongly convex and smooth conditions, generalization analysis requires smoothing parameter size of $v = \Theta((nd)^{-1})$. And*

number of steps $T = \Theta(n^2)$. We have the expected PS population risk in gradients is $\mathcal{O}(n^{-\frac{1}{2}})$, aligning with the upper bound for the PS population risk in general nonconvex first-order methods as shown in Chen et al. (2024).

Zeroth-order method demonstrates the connection between optimization, conflict avoidance, and generalization.

The core of the SZMOD algorithm lies in its dynamic weighting mechanism, which uses approximate gradient information to update λ . A high-quality λ is essential for balancing conflicts among multiple objectives. The distance to the CA direction is a critical metric for assessing the quality of these updates and plays a pivotal role in ensuring algorithmic convergence. In SZMOD, the deviation from the CA direction arises from the data and limited iterations and the cumulative error e introduced by the zeroth-order method. This CA direction error transfers the cumulative error e into an optimization error. Theoretical results indicate that in corresponding first-order algorithms, the relationship between CA direction error and optimization error is not as inherently inheritable and may exhibit a degree of antagonism (Chen et al., 2024). Thus, zeroth-order optimization opens a window into understanding the interaction between CA direction and optimization. Due to the propagation of cumulative error, optimization error imposes constraints on the smooth parameter v to ensure convergence. Furthermore, under strongly convex and smooth conditions, achieving generalization depends on controlling the size of v . Therefore, determining the appropriate value of v requires balancing the demands of both generalization and optimization.

7 EMPIRICAL VALIDATION

In this section, we systematically evaluate the performance of our proposed SZMOD algorithm on toy examples and CIFAR-10 datasets. The experiments are designed to mimic a variety of multi-objective landscapes with adjustable complexity levels. We employ synthetic datasets and realistic image data that encapsulate the essential characteristics of multi-objective problems for evaluating the optimization accuracy, generalization capability, conflict avoidance, and convergence performance of our proposal SZMOD algorithm.

7.1 SYNTHETIC EXPERIMENT

In the following content, we explore the subtleties of the SZMOD algorithm’s efficacy across a spectrum of hyperparameters, particularly emphasizing the trade-offs between optimization, generalization capabilities, and the mitigation of conflicting objectives. The synthetic experiments have been meticulously crafted to emulate a multi-objective optimization context, which successfully evaluates the influence exerted by diverse hyperparameters.

Strongly Convex Scenario: Inspired by (Chen et al., 2024), the following formulation is exploited to generate the MOL examples, whose m -th objective function is

$$f_{z,m}(x) = \frac{1}{2}b_{1,m}x^\top Ax - b_{2,m}z^\top x,$$

where $b_{1,m} > 0$ for all $m \in [M]$, and $b_{2,m}$ is another scalar. We set $M = 3$, $b_1 = [b_{1,1}; b_{1,2}; b_{1,3}] = [1; 2; 1]$, and $b_2 = [b_{2,1}; b_{2,2}; b_{2,3}] = [1; 3; 2]$. Each experimental setting has been repeated ten times, where the average results with standard deviation information are recorded in Figure 7.1. The detailed experimental settings for nonconvex cases are left in Appendix A.

The number of iterations, T , plays a pivotal role in the convergence properties of the SZMOD algorithm. As depicted in Figure 2a, we maintain $\alpha = 0.05$ and $\gamma = 0.001$ while varying T . The results indicate that an increase in T brings a decrease in both the optimization error and the distance to the conflict-avoidant (CA) direction, aligning with our theoretical predictions in Theorem 1, 2. This observation underscores the importance of sufficient training duration to achieve optimal solutions in multi-objective landscapes.

The step size for model parameters, α , is another critical hyperparameter that influences the algorithm’s ability to navigate the multi-objective space. In Figure 2b, we fix $T = 500$ and $\gamma = 0.001$ while adjusting α . The findings reveal an initial decrease in the optimization error as α increases, while further enlarging α does not yield significant improvements. This non-linear relationship

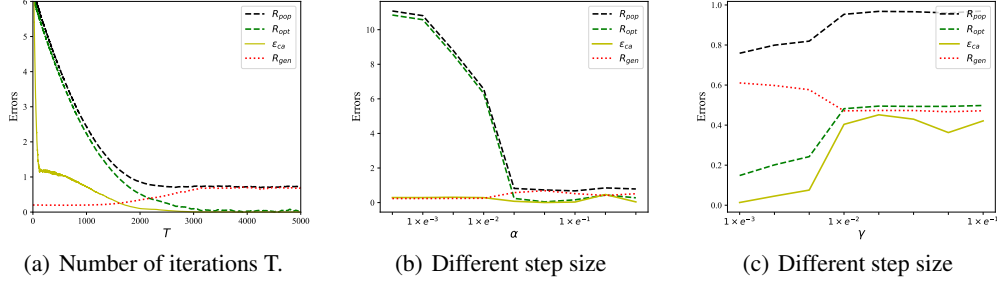


Figure 2: Optimization, generalization, and CA direction errors of SZMOD in the strongly convex case under different T, α, γ . The default parameters are $T = 500, \alpha = 0.05, \gamma = 0.001, v = 0.0001$.

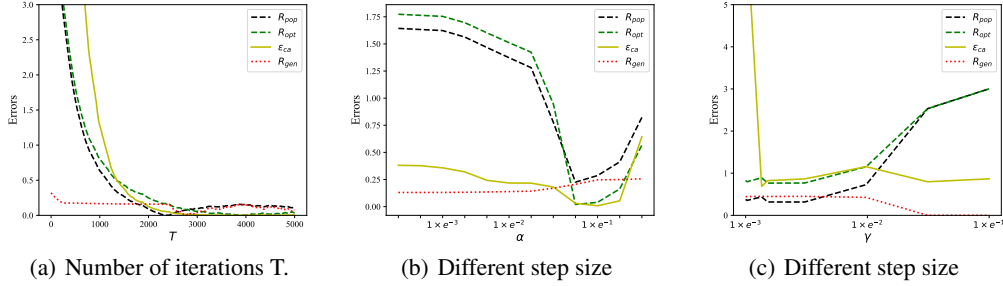


Figure 3: Optimization, generalization, and CA direction errors of SZMOD in the nonconvex case for MNIST image classification under different T, α, γ . The default parameters are $T = 500, \alpha = 0.05, \gamma = 0.001, v = 0.0001$.

between α and the optimization error highlights the need to carefully tune this hyperparameter to balance rapid convergence and potential overshooting of optimal solutions.

The weight step size, γ , is a unique aspect of SZMOD, controlling the update pace of the weighting parameters. In Figure 2c, with $T = 500$ and $\alpha = 0.05$, one can observe that the increasing γ leads to a decrease in the distance to the CA direction, suggesting that a more aggressive update of weights can be beneficial for navigating conflicting objectives. However, too large γ might lead to instability in convergence, indicating a delicate balance is required to harness the full potential of dynamic weighting.

The synthetic experiments provide valuable insights into the role of hyperparameters in shaping the trade-offs between optimization, generalization, and conflict avoidance in multi-objective learning. By systematically varying T, α , and γ , we have demonstrated the nuanced interplay between these parameters and their impact on the algorithm’s performance. These findings serve as a foundation for developing more sophisticated hyperparameter tuning strategies and provide empirical evidence to support theoretical analyses presented in prior sections. It is worth noting that, unlike the first-order MODO algorithm, the trends of $R_{opt}(\gamma)$ and $E_{ca}(\gamma)$ are not always opposite. This is due to the error caused by $\varepsilon_{ca}(\gamma)$, which is related to γ . When the trends are aligned, the graph of $R_{opt}(\gamma)$ always shows similar changes after changes occur in the graph of $\varepsilon_{ca}(\gamma)$. This is precisely due to error propagation, which nicely validates our theory.

7.2 ATTACK EXPERIMENT ON CIFAR-10

Adversarial attacks trick machine learning models by adding carefully designed subtle perturbations to inputs, leading to mispredictions. Black-box adversarial attacks occur when attackers can’t access a model’s internals and must deduce its behavior from inputs and outputs. The Black-box attack method is closer to real-world attack scenarios. Therefore, we consider a multi-objective adversarial attack.

Table 1: Results for multi-objective black-box adverbial attacks

model	Pixel ratio	ASR	L0_avg	L2_avg	AST_avg	SSIM_avg
CNN	2%	0.99	0.019	357.87	13.98	0.9
CNN	5%	0.98	0.049	572.78	8.47	0.78
CNN	10%	0.98	0.097	746.87	7.18	0.65
VGG16	2%	0.99	0.02	25.92	2.46	0.92
VGG16	5%	0.98	0.049	40.23	3.52	0.82
VGG16	10%	1	0.097	477.15	2.3	0.64
Alexnet	2%	0.99	0.019	250.94	7.09	0.85
Alexnet	5%	1	0.049	394.19	7.75	0.71
Alexnet	10%	1	0.097	342.58	4.8	0.62
Densenet	2%	0.91	0.019	22.71	10.7	0.88
Densenet	5%	0.92	0.049	18.26	13.98	0.83
Densenet	10%	0.86	0.097	12.22	13.18	0.87
Res-net18	2%	0.99	0.019	6.81	11.69	0.95
Res-net18	5%	0.98	0.049	3.85	11.04	0.97
Res-net28	10%	0.98	0.097	4.96	18.86	0.95

Define the loss function $\mathcal{L}(\mathbf{x} + \delta)$. We aim to generate a δ that solves the following optimization problem:

$$\min_{\vec{\delta}} F(\mathbf{x} + \vec{\delta}) \quad \text{s.t.} \quad \|\vec{\delta}\|_0 \leq \epsilon, \quad 0 \leq \mathbf{x} + \vec{\delta} \leq 1,$$

where $F(\mathbf{x} + \vec{\delta}) = (\mathcal{L}(\mathbf{x} + \vec{\delta}), \|\vec{\delta}\|_2, \|\vec{\delta}\|_0)^\top$ is the objective vector. $\vec{\delta}$ is the universal perturbation that we seek to optimize. We use the pre-trained model on the CIFAR-10 dataset, we attacked five classifiers: CNN, VGG16, AlexNet, DenseNet, and ResNet. Two types of attacks were implemented: targeted and non-targeted attacks. In the targeted attack, the cross-entropy loss function was used to misclassify the model into a specific target class, while the non-targeted attack employed margin loss to force the model's output to differ from the actual class. Additionally, the algorithm restricted perturbations to the discrete value set $\{-1, 1, 0\}$, which helped reduce the l_2 norm and ensured sparsity, enhancing both the effectiveness and stealth of the attack. **Metrics to evaluate the performance of attack methods include:** Average Attack Success Rate (ASR_avg), which measures the average success rate of misclassification due to adversarial attacks; Attack Success Rate (ASR), indicating the proportion of successful misclassifications; l_0 and l_2 norms, where l_0 counts the modified pixels and l_2 assesses perturbation magnitude; and Structural Similarity Index (SSIM), evaluating the similarity between the adversarial example and the original image, with values closer to 1 indicating less perceptible modifications.

We set $M = 2$, $\alpha = 0.1$, $\gamma = 0.001$, $v = 0.0001$, the maximum number of attack attempts 1000, and maximum modification per pixel 0.5. The corresponding results in Table 1 imply that the higher accuracy of the model could bring better effectiveness of the attack, which aligns with the principles of the zeroth-order multi-objective algorithm (*the more accurate the loss, the more accurate the gradient based on the loss*). Moreover, our attack success rate is generally above 90 percent, further demonstrating the advantages of our algorithm.

8 CONCLUSION

In this paper, we introduce the SZMOD algorithm, designed explicitly for black-box multi-objective learning. Theoretically, we establish the statistical guarantees for optimization error, generalization bound, and distance to conflict avoidance directions comparable to the relevant first-order method. Furthermore, we discover that zeroth-order methods could bridge the above three evaluation criteria of SZMOD. Experimentally, we validate SZMOD's performance in terms of optimization accuracy, generalization capability, and conflict avoidance. Additionally, we demonstrate the effectiveness of our algorithm in practical black-box attack scenarios, as evidenced by high attack success rates and low modification rates.

REFERENCES

- Naveed Akhtar and Ajmal Mian. Threat of adversarial attacks on deep learning in computer vision: A survey. *Ieee Access*, 6:14410–14430, 2018.
- Lisha Chen, Heshan Fernando, Yiming Ying, and Tianyi Chen. Three-way trade-off in multi-objective learning: Optimization, generalization and conflict-avoidance. *Advances in Neural Information Processing Systems*, 36, 2024.
- Jean-Antoine Désidéri. Multiple-gradient descent algorithm (mgda) for multiobjective optimization. *Comptes Rendus Mathématique*, 350(5-6):313–318, 2012.
- Florian Felten, Umut Ucak, Hicham Azmani, Gao Peng, Willem Röpke, Hendrik Baier, Patrick Mannion, Diederik M Roijers, Jordan K Terry, El-Ghazali Talbi, et al. Momaland: A set of benchmarks for multi-objective multi-agent reinforcement learning. *arXiv preprint arXiv:2407.16312*, 2024.
- Heshan Fernando, Han Shen, Miao Liu, Subhajit Chaudhury, Keerthiram Murugesan, and Tianyi Chen. Mitigating gradient bias in multi-objective learning: A provably convergent approach. *International Conference on Learning Representations*, 2023.
- Jörg Fliege, A Ismael F Vaz, and Luís Nunes Vicente. Complexity of gradient descent for multiobjective optimization. *Optimization Methods and Software*, 34(5):949–959, 2019.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pp. 1321–1330. PMLR, 2017.
- Jayesh K Gupta, Maxim Egorov, and Mykel Kochenderfer. Cooperative multi-agent control using deep reinforcement learning. In *Autonomous Agents and Multiagent Systems: AAMAS 2017 Workshops, Best Papers, São Paulo, Brazil, May 8-12, 2017, Revised Selected Papers 16*, pp. 66–83. Springer, 2017.
- Tianmeng Hu, Biao Luo, Chunhua Yang, and Tingwen Huang. Mo-mix: Multi-objective multi-agent cooperative decision-making with deep reinforcement learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(10):12098–12112, 2023.
- Yunwen Lei. Stability and generalization of stochastic optimization with nonconvex and nonsmooth problems. In *The Thirty Sixth Annual Conference on Learning Theory*, pp. 191–227. PMLR, 2023.
- Siyuan Liang, Longkang Li, Yanbo Fan, Xiaojun Jia, Jingzhi Li, Baoyuan Wu, and Xiaochun Cao. A large-scale multiple-objective method for black-box attack against object detection. In *European Conference on Computer Vision*, pp. 619–636. Springer, 2022.
- Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for multi-task learning. *Advances in Neural Information Processing Systems*, 34:18878–18890, 2021.
- Shengcai Liu, Ning Lu, Wenjing Hong, Chao Qian, and Ke Tang. Effective and imperceptible adversarial textual attack via multi-objectivization. *ACM Transactions on Evolutionary Learning and Optimization*, 4(3):1–23, 2024.
- Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and black-box attacks. In *International Conference on Learning Representations*, 2022.
- Junlin Lu, Patrick Mannion, and Karl Mason. A multi-objective multi-agent deep reinforcement learning approach to residential appliance scheduling. *IET Smart Grid*, 5(4):260–280, 2022.
- Alex Mathai, Shreya Khare, Srikanth Tamilselvam, and Senthil Mani. Adversarial black-box attacks on text classifiers using multi-objective genetic optimization guided by deep networks. *arXiv preprint arXiv:2011.03901*, 2020.
- Brady Neal, Sarthak Mittal, Aristide Baratin, Vinayak Tantia, Matthew Scicluna, Simon Lacoste-Julien, and Ioannis Mitliagkas. A modern take on the bias-variance tradeoff in neural networks. *arXiv preprint arXiv:1810.08591*, 2018.

- Nicolas Papernot, Patrick McDaniel, and Ian Goodfellow. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *arXiv preprint arXiv:1605.07277*, 2016.
- Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram Swami. Practical black-box attacks against machine learning. In *Proceedings of the 2017 ACM on Asia conference on computer and communications security*, pp. 506–519, 2017.
- Tianxiang Sun, Yunfan Shao, Hong Qian, Xuanjing Huang, and Xipeng Qiu. Black-box tuning for language-model-as-a-service. In *International Conference on Machine Learning*, pp. 20841–20855. PMLR, 2022.
- Hiroki Tanabe, Ellen H Fukuda, and Nobuo Yamashita. Proximal gradient methods for multiobjective optimization and their applications. *Computational Optimization and Applications*, 72:339–361, 2019.
- Jordan Terry, Benjamin Black, Nathaniel Grammel, Mario Jayakumar, Ananth Hari, Ryan Sullivan, Luis S Santos, Clemens Dieffendahl, Caroline Horsch, Rodrigo Perez-Vicente, et al. Pettingzoo: Gym for multi-agent reinforcement learning. *Advances in Neural Information Processing Systems*, 34:15032–15043, 2021.
- Phoenix Neale Williams and Ke Li. Black-box sparse adversarial attack via multi-objective optimisation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12291–12301, 2023.
- Stephan Zheng, Alexander Trott, Sunil Srinivasa, David C Parkes, and Richard Socher. The ai economist: Taxation policy design via two-level deep multiagent reinforcement learning. *Science advances*, 8(18):eabk2607, 2022.
- Shasha Zhou, Mingyu Huang, Yanan Sun, and Ke Li. Evolutionary multi-objective optimization for contextual adversarial example generation. *Proceedings of the ACM on Software Engineering*, 1 (FSE):2285–2308, 2024.

Appendix

A ADDITIONAL EXPERIMENTS AND IMPLEMENTATION DETAILS

Implementation. Experiments are done on a machine with GPU NVIDIA RTX 4060. We use Python 3.8, CUDA 11.8, Pytorch 1.8.0 for all experiments. Unless otherwise stated, all experiments are repeated with five random seeds. Their average performance and standard deviations are reported throughout the whole manuscript.

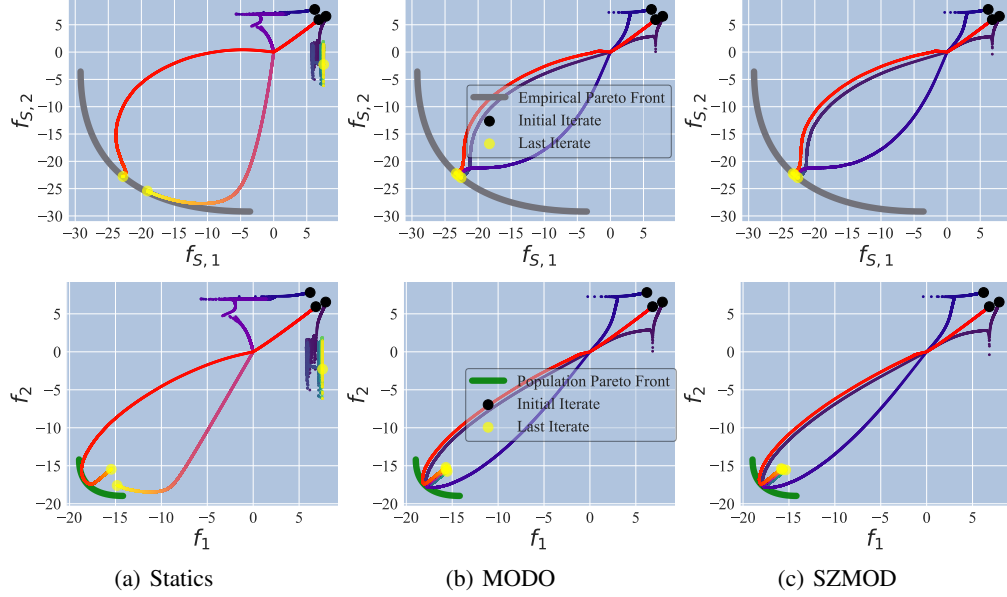


Figure 4: Convergence of static weighting, MoDo and SZMOD to the empirical (gray, upper) and population (green, lower) Pareto fronts. The horizontal and vertical axes in the figures in the first/second row are the values of the two empirical / population objectives. Three colormaps are used for the trajectories from three initializations, respectively, where the same colormaps represent the trajectories of the same initializations, darker colors in one colormap indicate earlier iterations, and lighter colors indicate later iterations.

A.1 EXPERIMENTS ON NONCONVEX OBJECTIVES

Implementation details. The toy example is modified from (Liu et al., 2021) to consider stochastic data. Denote the model parameter as $x = [x_1, x_2]^T \in \mathbb{R}^2$, stochastic data as $z = [z_1, z_2]^T \in \mathbb{R}^2$ sampled from the standard multi-variate Gaussian distribution. The individual empirical objectives are defined as:

$$\begin{aligned}
 f_{z,1}(x) &= c_1(x)h_1(x) + c_2(x)g_{z,1}(x) \quad \text{and} \quad f_{z,2}(x) = c_1(x)h_2(x) + c_2(x)g_{z,2}(x) \quad \text{where} \\
 h_1(x) &= \log(\max(|0.5(-x_1 - 7) - \tanh(-x_2)|, 0.000005)) + 6, \\
 h_2(x) &= \log(\max(|0.5(-x_1 + 3) - \tanh(-x_2) + 2|, 0.000005)) + 6 \\
 g_{z,1}(x) &= ((-x_1 + 3.5)^2 + 0.1 * (-x_2 - 1)^2)/10 - 20 - 2 * z_1x_1 - 5.5 * z_2x_2, \\
 g_{z,2}(x) &= ((-x_1 - 3.5)^2 + 0.1 * (-x_2 - 1)^2)/10 - 20 + 2 * z_1x_1 - 5.5 * z_2x_2, \\
 c_1(x) &= \max(\tanh(0.5 * x_2), 0) \quad \text{and} \quad c_2(x) = \max(\tanh(-0.5 * x_2), 0).
 \end{aligned}$$

Since z is zero-mean, the individual population objectives are correspondingly:

$$f_1(x) = c_1(x)h_1(x) + c_2(x)g_1(x) \text{ and } f_2(x) = c_1(x)h_2(x) + c_2(x)g_2(x), \text{ where}$$

$$g_1(x) = ((-x_1 + 3.5)^2 + 0.1 * (-x_2 - 1)^2)/10 - 20,$$

$$g_2(x) = ((-x_1 - 3.5)^2 + 0.1 * (-x_2 - 1)^2)/10 - 20.$$

The training dataset size is $n = |S| = 20$. For all methods, i.e., static weighting, MoDo, SZMOD, the number of iterations is $T = 10000$. The initialization of λ is $\lambda_0 = [0.5, 0.5]^\top$.

In Figure 4 and Figure 1, the trajectories of various methods from different initializations to the empirical and population Pareto fronts (PF) are shown. In Figure 4a (first row), the static weighting method with uniform weights shows one trajectory successfully converging to the center of the empirical PF. In contrast, the other two trajectories oscillate around suboptimal parameters, forming clusters of scattered points, with one failing to reach the empirical PF altogether. Only one empirically suboptimal solution (shown by the red-to-yellow trajectory) achieves low population risk in the second row. In Figures 4b and 4c, MODO and SZMOD demonstrate identical convergence rates in the first row, with both methods converging to the center of the empirical PF, representing the optimal solution for the uniform average of the two objectives. In the second row, all three solutions for both MODO and SZMOD achieve relatively low population risk, highlighting their strong generalization ability. Comparing Figures 4b and 4c, we observe that MODO and SZMOD exhibit nearly identical convergence trajectories under the same parameter settings and initializations, confirming that SZMOD maintains strong performance even without accurate gradients.

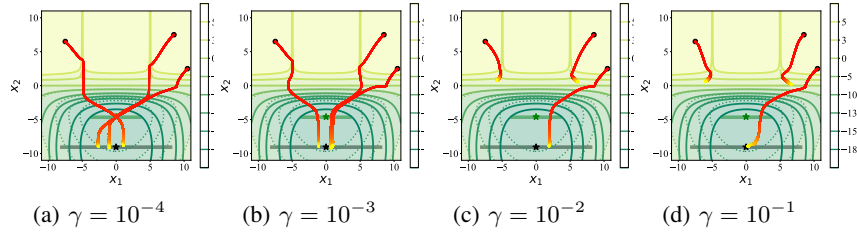


Figure 5: Trajectories of SZMOD under different γ on the contour of the average of objectives. The black \cdot marks initializations of the trajectories, colored from red (start) to yellow (end). The background solid/dotted contours display the landscape of the average empirical/population objectives. The gray/green bar marks empirical/population Pareto front, and the black \star / green \star marks solution to the average objectives.

To demonstrate how the choice of γ impacts the performance of SZMOD, we further conduct experiments with different SZMOD. We should mention that when $\gamma \geq 10^{-2}$, SZMOD did not converge to the Pareto front. This is because the error term of SZMOD is related to the value of γ . If the γ is large enough, it will cause the error term to be too large, resulting in the algorithm not converging.

B PRELIMINARIES FOR PROOF

B.1 ZERO-TH-ORDER GRADIENT ESTIMATION

When only function evaluations are available, we employ the deterministic coordinate-wise direction to derive the decent direction. Specifically, for the smoothing constant v and vector u_i (u_i represents the unit vector where the i -th element is 1 and the remaining elements are 0), the directional derivative of $f_{z,m}$ in the direction u for the smooth function $f_i, i \in [n]$, can be estimated as:

$$\hat{\nabla} f_{z,m}(x, u, v) = \sum_{j=1}^d \frac{f_{z,m}(x + v u_j) - f_{z,m}(x)}{v} u_j.$$

As the approximation of the full directional gradient. Since the smoothing constant v is fixed, for simplicity, we leave out v in these gradient estimations and set

$$\hat{\nabla} f_{z,m}(x, u) := \hat{\nabla} f_{z,m}(x, u, v)$$

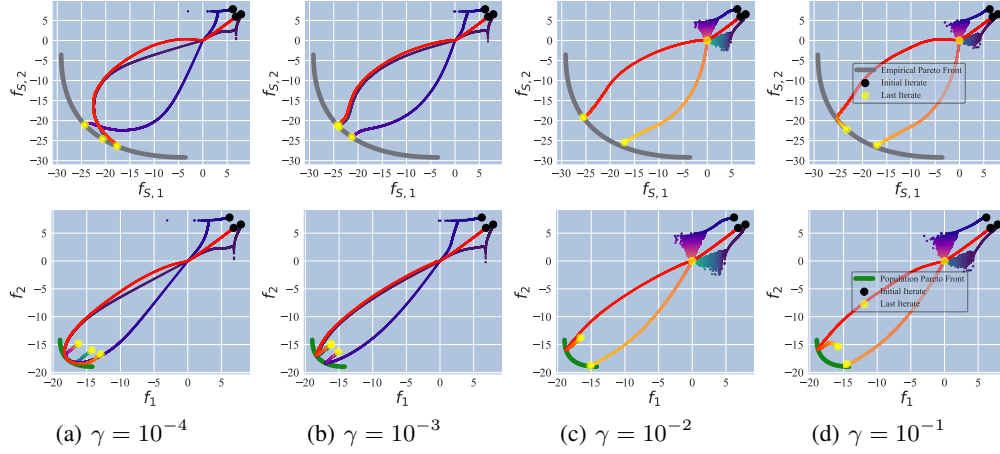


Figure 6: Convergence of SZMOD to the empirical (gray, upper) and population (green, lower) Pareto fronts under different γ . The horizontal and vertical axes in the first/second-row figures are the values of the two empirical / population objectives. Three colormaps are used for the trajectories from three initializations, respectively, where the same colormaps represent the trajectories of the same initializations; darker colors in one colormap indicate earlier iterations, and lighter colors indicate later iterations.

Lemma 2. For the $l_{f,1}$ -smooth function $f_{z,m}$ and any $x \in \mathbb{R}^d$, $i \in [n]$, define $s_z(x, u_j) \in [0, 1]$ and $s_z = [s_z(x, u_1), s_z(x, u_2), \dots, s_z(x, u_d)]$, $s = [s_1; s_2; \dots; s_m]$ the estimator in satisfies:

$$\hat{\nabla} F(x)\lambda = \nabla F(x)\lambda + \frac{l_{f,1}v}{2}s(x, u)\lambda \leq \nabla F(x)\lambda + \frac{l_{f,1}v}{2}\mathbb{1}_d \quad (14)$$

Proof:

$$\begin{aligned} \hat{\nabla} f_{z,m}(x, u) &= \sum_{j=1}^d \frac{f_{z,m}(x + vu_j) - f_{z,m}(x)}{v} u_j \\ &\stackrel{(a)}{=} \sum_{j=1}^d (u_j u_j^\top \nabla f_{z,m}(x) + \frac{v}{2} u_j \nabla^2 f_{z,m}(x) u_j^\top u_j) \\ &\stackrel{(b)}{=} \nabla f_{z,m}(x) + \sum_{j=1}^d \frac{l_{f,1}v}{2} s_{z,j}(x, u) \|u_j\|^2 u_j \\ &= \nabla f_{z,m}(x) + \frac{l_{f,1}v}{2} s_z(x, u) \end{aligned}$$

Here (a) is based on Taylor expansion, and (b) is based on $0 \leq u_j \nabla^2 f_{z,m}(x) u_j^\top \leq l_{f,1}^2$. Then, by the definition of $\hat{\nabla} F(x)$, we have:

$$\begin{aligned} \hat{\nabla} F(x)\lambda &= [\hat{\nabla} f_{z,1}, \hat{\nabla} f_{z,2}, \dots, \hat{\nabla} f_{z,m}]\lambda \\ &= \nabla F(x)\lambda + \frac{l_{f,1}v}{2}s(x, u)\lambda \\ &\leq \nabla F(x)\lambda + \frac{l_{f,1}v}{2}\mathbb{1}_d \end{aligned}$$

The inequality is based on the definition of λ . Here, we complete the proof.

In Eq.(4), $s_i(x, u)$ measures the curvature scaled by L along the specified direction u at the given point x . By taking the maximum value of $s_i(x, u)$ (i.e., $s_i(x, u) := 1$), we derive the upper bound of the distance between estimator $\hat{\nabla} f_i(x, u)$ and the full gradient $\nabla f(x)$ in Eq.(6). This bound

comprises three components: the norm of the true gradient $\nabla f(x)$; the trivial perturbation; and the last error $\|\nabla f_i(x) - \nabla f(x)\|^2$ induced by the random sampling of i . Similarly, the following corollary establishes the nearly unbiased and bounded variance properties of $\hat{\nabla} f(x, u)$.

B.2 LEMMAS FOR PROOF

B.3 PROOF OF LEMMA 1

Lemma 3. *Suppose Assumptions 2, 3 hold. WLOG, assume $\inf_{x \in \mathbb{R}^d} f_{m,z}(x) < \infty$ for all $m \in [M]$ and $z \in \mathcal{Z}$. For any given $\lambda \in \Delta^M$, and stochastic sample $z \in \mathcal{Z}$, define $x_{\lambda,z}^* = \arg \min_{x \in \mathbb{R}^d} F_z(x)\lambda$, then $\inf_{x \in \mathbb{R}^d} F_z(x)\lambda < \infty$ and $\|x_{\lambda,z}^*\| < \infty$, i.e., there exist finite positive constants c_{F^*} and c_{x^*} such that*

$$\inf_{x \in \mathbb{R}^d} F_z(x)\lambda \leq c_{F^*} \text{ and } \|x_{\lambda,z}^*\| \leq c_{x^*}$$

Proofed by (Chen et al., 2024).

Lemma 4. *Suppose Assumptions 2, 3 hold, and define $\kappa = 3\ell_{f,1}/\mu \geq 3$. For any given $\lambda \in \Delta^M$, and a stochastic sample $z \in \mathcal{Z}$, define $x_{\lambda,z}^* = \arg \min_x F_z(x)\lambda$. Then by Lemma 3, there exists a positive finite constant $c_{x,1} \geq c_{x^*}$ such that $\|x_{\lambda,z}^*\| \leq c_{x^*} \leq c_{x,1}$. Recall the multi-objective gradient update is*

$$G_{\lambda,z}(x) = x - \alpha \hat{\nabla} F_z(x)\lambda$$

with step size $0 \leq \alpha \leq \ell_{f,1}^{-1}$. Defining $v' = dv/c_{x,1}$ $c_{x,2} = (1 + \frac{l_{f,1}v'}{2}\sqrt{(2+v')\kappa})c_{x,1}$, we have that

$$\text{if } \|x\| \leq c_{x,2}, \text{ then } \|G_{\lambda,z}(x)\| \leq c_{x,2}$$

Proof. We divide the proof into two cases:

1) when $\|x\| < c_{x,1}$. That is $\|x\| < c_{x,1} \leq c_{x,2}$, then we have

$$\begin{aligned} \|G_{\lambda,z}(x)\| &\leq \|G_{\lambda,z}(x) - x^*\| + \|x^*\| \\ &\stackrel{(a)}{=} \|x - \nabla F_z(x)\lambda - (x^* - \nabla F_z(x^*)\lambda) + \nabla F_z(x)\lambda - \hat{\nabla} F_z(x)\lambda\| + \|x^*\| \\ &\leq \|x - \nabla F_z(x)\lambda - (x^* - \nabla F_z(x^*)\lambda)\| + \|\nabla F_z(x)\lambda - \hat{\nabla} F_z(x)\lambda\| + \|x^*\| \\ &\stackrel{(b)}{\leq} \|x - x^*\| + \|x^*\| + \frac{\ell_{f,1}vd}{2} \\ &\leq \|x\| + 2\|x^*\| + \frac{\ell_{f,1}vd}{2} \leq (3 + \frac{l_{f,1}v'}{2})c_{x,1} \leq (1 + \sqrt{6})c_{x,1} \leq (1 + \frac{l_{f,1}v'}{2}\sqrt{(2+v')\kappa})c_{x,1} \leq c_{x,2} \end{aligned}$$

where (a) follows from $\nabla F_z(x^*)\lambda = 0$, and (b) follows from the non-expansiveness of the gradient update for strongly convex and smooth function.

2) $c_{x,1} \leq \|x\| \leq c_{x,2}$, we first consider $\alpha = \ell_{f,1}^{-1}$. Let $\mu' = \mu/3$. Note that since $F_z(x)\lambda$ is μ -strongly convex, it is also μ' -strongly convex. By strong convexity and smoothness of $F_z(x)\lambda$, the gradients are co-coercive [36, Theorem 2.1.12], i.e., for any x we have

$$(\nabla F_z(x)\lambda)^\top (x - x^*) \geq \frac{\ell_{f,1}^{-1} \|\nabla F_z(x)\lambda\|^2}{1 + \kappa^{-1}} + \frac{\mu' \|x - x^*\|^2}{1 + \kappa^{-1}}.$$

For the ZO version we have:

$$\begin{aligned} (\hat{\nabla} F_z(x)\lambda)^\top (x - x^*) &= (\nabla F_z(x)\lambda)^\top (x - x^*) + \left(\frac{l_{f,1}v}{2} s(x, v)\lambda \right)^\top (x - x^*) \\ &\geq \frac{\ell_{f,1}^{-1} \|\nabla F_z(x)\lambda\|^2}{1 + \kappa^{-1}} + \frac{\mu' \|x - x^*\|^2}{1 + \kappa^{-1}} + \left(\frac{l_{f,1}v}{2} s(x, v)\lambda \right)^\top (x - x^*) \end{aligned}$$

Rearranging and applying Cauchy-Schwartz inequality, we have

$$\begin{aligned} \left(\hat{\nabla} F_z(x) \lambda \right)^\top x &\geq (\nabla F_z(x) \lambda)^\top x^* + \frac{l_{f,1} v}{2} s(x, v)^\top \lambda x + \frac{\ell_{f,1}^{-1} \|\nabla F_z(x) \lambda\|^2}{1 + \kappa^{-1}} + \frac{\mu' \|x - x^*\|^2}{1 + \kappa^{-1}}. \\ &\geq -c_{x,1} \|\nabla F_z(x) \lambda\| - \frac{c_{x,1} l_{f,1} dv}{2} + \frac{\ell_{f,1}^{-1} \|\nabla F_z(x) \lambda\|^2}{1 + \kappa^{-1}} + \frac{\mu' \|x - x^*\|^2}{1 + \kappa^{-1}}. \end{aligned} \quad (15)$$

By the definition of $G_{\lambda,z}(x)$,

$$\|G_{\lambda,z}(x)\|^2 = \left\| x - \frac{1}{\ell_{f,1}} \hat{\nabla} F_z(x) \lambda \right\|^2 = \|x\|^2 + \frac{1}{\ell_{f,1}^2} \left\| \hat{\nabla} F_z(x) \lambda \right\|^2 - \frac{2}{\ell_{f,1}} \left(\hat{\nabla} F_z(x) \lambda \right)^\top x \quad (16)$$

From (15) and (16), we have:

$$\begin{aligned} \|G_{\lambda,z}(x)\|^2 &\leq \|x\|^2 + \frac{1}{\ell_{f,1}^2} \left\| \hat{\nabla} F_z(x) \lambda \right\|^2 + \frac{2}{\ell_{f,1}} (c_{x,1} \|\nabla F_z(x) \lambda\| \\ &\quad + \frac{c_{x,1} l_{f,1} dv}{2} - \frac{\ell_{f,1}^{-1} \|\nabla F_z(x) \lambda\|^2}{1 + \kappa^{-1}} - \frac{\mu' \|x - x^*\|^2}{1 + \kappa^{-1}}) \\ &\leq \|x\|^2 + \frac{2}{\ell_{f,1}} \left(\frac{c_{x,1} l_{f,1} dv}{2} - \frac{\mu' \|x - x^*\|^2}{1 + \kappa^{-1}} \right) \\ &\quad + \frac{2}{\ell_{f,1}} \sup_{\gamma \in \mathbb{R}} \underbrace{\left(c_{x,1} \cdot \gamma - \frac{1}{2\ell_{f,1}} \left(\frac{1 - \kappa^{-1}}{1 + \kappa^{-1}} \right) \gamma^2 \right)}_{I_1} \end{aligned} \quad (17)$$

Since $\kappa \geq 3$, thus $\frac{1 - \kappa^{-1}}{1 + \kappa^{-1}} > 0$, then I_1 is a quadratic function w.r.t. γ , and is strictly concave, thus can be bounded above by

$$\sup_{\gamma \in \mathbb{R}} c_{x,1} \cdot \gamma - \frac{1}{2\ell_{f,1}} \left(\frac{1 - \kappa^{-1}}{1 + \kappa^{-1}} \right) \gamma^2 \leq \frac{c_{x,1}^2 \ell_{f,1}}{2} \frac{1 + \kappa^{-1}}{1 - \kappa^{-1}}$$

Substituting this back into (17) gives that

$$\begin{aligned} \|G_{\lambda,z}(x)\|^2 &\leq \|x\|^2 + \frac{2}{\ell_{f,1}} \left(\frac{c_{x,1} l_{f,1} dv}{2} - \frac{\mu' \|x - x^*\|^2}{1 + \kappa^{-1}} + \frac{c_{x,1}^2 \ell_{f,1}}{2} \frac{1 + \kappa^{-1}}{1 - \kappa^{-1}} \right) \\ &= \|x\|^2 + c_{x,1}^2 \frac{1 + \kappa^{-1}}{1 - \kappa^{-1}} + c_{x,1} dv + \frac{l_{f,1} v^2}{4} - 2 \frac{\kappa^{-1}}{1 + \kappa^{-1}} \|x - x^*\|^2 \\ &\leq \|x\|^2 + c_{x,1}^2 \frac{1 + \kappa^{-1}}{1 - \kappa^{-1}} - 2 \frac{\kappa^{-1}}{1 + \kappa^{-1}} (\|x\| - \|x^*\|)^2 + c_{x,1} v \\ &\leq \underbrace{\|x\|^2 + 2c_{x,1}^2 - \kappa^{-1} (\|x\| - c_{x,1})^2}_{I_2} + c_{x,1} dv \end{aligned}$$

Here $\frac{2}{\ell_{f,1}} \geq \min\{\frac{l_{f,1} v^2}{2}, (\|x\| - c_{x,1})^2\}$ where the last inequality follows from $\kappa \geq 3$, thus $\frac{1 + \kappa^{-1}}{1 - \kappa^{-1}} \leq 2$, $-2 \frac{\kappa^{-1}}{1 + \kappa^{-1}} \leq -\kappa^{-1}$, and $\|x^*\| \leq c_{x,1} \leq \|x\|$ by assumption. For $c_{x,1} \leq \|x\| \leq c_{x,2}$, I_2 is a strictly convex quadratic function of $\|x\|$, which achieves its maximum at $\|x\| = c_{x,1}$ or $\|x\| = c_{x,2}$. Therefore,

$$\begin{aligned} \|G_{\lambda,z}(x)\|^2 &\leq \max \left\{ 3c_{x,1}^2 + c_{x,1} v, c_{x,2}^2 + 2c_{x,1}^2 - \kappa^{-1} (c_{x,2} - c_{x,1})^2 + c_{x,1} v \right\} \\ &\stackrel{(c)}{=} \max \left\{ 3c_{x,1}^2 + c_{x,1} v, c_{x,2}^2 \right\} \stackrel{(d)}{<} c_{x,2}^2 \end{aligned}$$

where (c) follows from the definition that $c_{x,2} = (1 + \frac{\ell_{f,1}v'}{2}\sqrt{(2+v')\kappa})c_{x,1}$; (d) follows from $\kappa \geq 3$, and thus $(3+v')c_{x,1}^2 < (1 + \frac{\ell_{f,1}v'}{2}\sqrt{(2+v')\kappa})^2 c_{x,1}^2 = c_{x,2}^2$. We have proved the case for $\alpha = \ell_{f,1}^{-1}$. The result for $0 \leq \alpha < \ell_{f,1}^{-1}$ follows by observing that,

$$\begin{aligned}\|G_{\lambda,z}(x)\| &= \|x - \alpha \hat{\nabla} F_z(x)\lambda\| \\ &= \|(1 - \alpha \ell_{f,1})x + \alpha \ell_{f,1} (x - \ell_{f,1}^{-1} \hat{\nabla} F_z(x)\lambda)\| \\ &\leq (1 - \alpha \ell_{f,1})\|x\| + \alpha \ell_{f,1} \|x - \ell_{f,1}^{-1} \hat{\nabla} F_z(x)\lambda\| \leq c_{x,2}\end{aligned}$$

The proof is complete.

Lemma 5. Suppose Assumptions 2, 3 hold. For all $\lambda \in \Delta^M$ and $z \in S$, define $x_{\lambda,z}^* = \arg \min_x F_z(x)\lambda$, then there exist finite positive constants c_{F^*} and c_{x^*} such that $F_z(x_{\lambda,z}^*)\lambda \leq c_{F^*}$ and $\|x_{\lambda,z}^*\| \leq c_{x^*}$. And for $x \in \mathbb{R}^d$ such that $\|x\|$ is bounded, i.e., there exists a finite positive constant c_x such that $\|x\| \leq c_x$, then

$$\|\nabla F_z(x)\lambda\| \leq \ell_{f,1}(c_x + c_{x^*}), \quad \text{and} \quad F_z(x)\lambda \leq \frac{\ell_{f,1}}{2}(c_x + c_{x^*})^2 + c_{F^*}$$

Proof. Under Assumptions 2, 3 by Lemma 3, there exist finite positive constants c_{F^*} and c_{x^*} such that $F_z(x_{\lambda,z}^*)\lambda \leq c_{F^*}$ and $\|x_{\lambda,z}^*\| \leq c_{x^*}$. By Assumption 1, the $\ell_{f,1}$ -Lipschitz continuity of the gradient $\nabla F_z(x)\lambda$, we have

$$\begin{aligned}\|\nabla F_z(x)\lambda\| &= \|\nabla F_z(x)\lambda - \nabla F_z(x_{\lambda,z}^*)\lambda\| \\ &\leq \ell_{f,1} \|x - x_{\lambda,z}^*\| \leq \ell_{f,1} (\|x\| + \|x_{\lambda,z}^*\|) \leq \ell_{f,1} (c_x + c_{x^*})\end{aligned}$$

where the first equality uses the fact that $\nabla F_z(x_{\lambda,z}^*)\lambda = 0$. For the function value, by Assumption 2, the $\ell_{f,1}$ -Lipschitz smoothness of $F_z(x)\lambda$, we have

$$\begin{aligned}F_z(x)\lambda &\leq F_z(x_{\lambda,z}^*)\lambda + \langle \nabla F_z(x_{\lambda,z}^*)\lambda, x - x_{\lambda,z}^* \rangle + \frac{\ell_{f,1}}{2} \|x - x_{\lambda,z}^*\|^2 \\ &\leq F_z(x_{\lambda,z}^*)\lambda + \frac{\ell_{f,1}}{2} \|x - x_{\lambda,z}^*\|^2 \\ &\leq c_{F^*} + \frac{\ell_{f,1}}{2} (c_x + c_{x^*})^2\end{aligned}$$

from which the proof is complete.

Corollary 1. Suppose Assumptions 2, 3 hold and $v < \frac{1}{T}$. Define $\kappa = 3\ell_{f,1}/\mu$ and $x_{\lambda,z}^* = \arg \min_x F_z(x)\lambda$ with $\lambda \in \Delta^M$. Then there exists a finite positive constant c_{x^*} such that $\|x_{\lambda,z}^*\| \leq c_{x^*}$. Choose the initial iterate to be bounded, i.e., there exists a finite positive constant c_{x_0} such that $\|x_0\| \leq c_{x_0}$, then for $\{x_t\}$ generated by SZMOD algorithm with $\alpha_t = \alpha$ and $0 \leq \alpha \leq \ell_{f,1}^{-1}$, we have

$$\|x_t\| \leq c_x, \quad \text{with } c_x = \max \left\{ \left(1 + \frac{\ell_{f,1}v'}{2}\sqrt{(2+v')\kappa}\right)c_{x^*} + \frac{\ell_{f,1}v}{4}, c_{x_0} \right\} \quad (18)$$

Proof. Under Assumptions 2, 3, by Lemma 3, $\|x_{\lambda,z}^*\| < \infty$, i.e., there exists a finite positive constant c_{x^*} such that $\|x_{\lambda,z}^*\| \leq c_{x^*}$. Let $c_{x,1} = \max \left\{ \left(1 + \frac{\ell_{f,1}v'}{2}\sqrt{(2+v')\kappa}\right)^{-1}c_{x_0}, c_{x^*} \right\}$, and

$c_{x,2} = (1 + \frac{l_{f,1}v'}{2}\sqrt{(2+v')\kappa})c_{x,1} = \max \left\{ c_{x_0}, (1 + \frac{l_{f,1}v'}{2}\sqrt{(2+v')\kappa})c_{x^*} \right\}$ in Lemma 4. We then consider the following two cases: 1) If $(1 + \frac{l_{f,1}v'}{2}\sqrt{(2+v')\kappa})c_{x^*} \leq c_{x_0}$, then $\|x_{\lambda,z}^*\| \leq c_{x^*} \leq (1 + \frac{l_{f,1}v'}{2}\sqrt{(2+v')\kappa})^{-1}c_{x_0}$. Then it satisfies the condition in Lemma 4 that $\|x_{\lambda,z}^*\| \leq c_{x,1}$ and $\|x_0\| \leq c_{x,2}$. Applying Lemma 4 yields $\|x_1\| \leq c_{x,2}$. 2) If $(1 + \frac{l_{f,1}v'}{2}\sqrt{(2+v')\kappa})c_{x^*} > c_{x_0}$, then $\|x_0\| \leq c_{x_0} < (1 + \frac{l_{f,1}v'}{2}\sqrt{(2+v')\kappa})c_{x^*}$. Then it satisfies the condition in Lemma 4 that $\|x_{\lambda,z}^*\| \leq c_{x,1}$ and $\|x_0\| \leq c_{x,2}$. Applying Lemma 4 yields $\|x_1\| \leq c_{x,2}$. Therefore, equation 18 holds for $t = 1$. We then prove by induction that equation 18 also holds for $t \in [T]$. Assume equation 18 holds at $1 \leq k \leq T-1$, i.e.,

$$\|x_k\| \leq c_x = c_{x,2}$$

Then by Lemma 4, at $k+1$,

$$\|x_{k+1}\| = \|G_{\lambda_{k+1}, z_{k,3}}(x_k)\| \leq c_{x,2}$$

Since $\|x_1\| \leq c_{x,2}$, for $t = 0, \dots, T-1$, we have

$$\|x_{t+1}\| = \|G_{\lambda_{t+1}, z_{t,3}}(x_t)\| \leq c_{x,2}$$

Therefore, by mathematical induction, $\|x_t\| \leq c_{x,2} = c_x$, for all $t \in [T]$. The proof is complete.

Proof of Lemma 1. By Corollary 1, for $\{x_t\}$ generated by SZMOD algorithm with $\alpha_t = \alpha$ and $0 \leq \alpha \leq \ell_{f,1}^{-1}$, we have

$$\|x_t\| \leq c_x, \quad \text{with} \quad c_x = \max \left\{ (2 + \frac{l_{f,1}v'}{2}\sqrt{(2+v')\kappa})c_{x^*} + \frac{l_{f,1}v}{4}, c_{x_0} \right\}$$

According to Lemma 3, define $\ell_f = \ell_{f,1}(c_x + c_{x^*})$, and $\ell_F = \sqrt{M}\ell_f$, then it holds for all $\lambda \in \Delta^M$

$$\|\nabla F(x_t)\lambda\| \leq \ell_f \quad \text{and} \quad \|\nabla F(x_t)\| \leq \|\nabla F(x_t)\|_F \leq \ell_F$$

Lemma 6. Suppose Assumptions 2, 3 hold. For all $\lambda \in \Delta^M$ and $z \in S$, according to Lemma 1, we have $\|\nabla F(x_t)\lambda\| \leq \ell_f$, and $\|\nabla F(x_t)\| \leq \ell_F$, then

$$\|\hat{\nabla} F_z(x)\lambda\| \leq \ell_{f,d}, \quad \text{and} \quad \|\hat{\nabla} F_z(x)\| \leq \ell_{F,d}$$

Where $\ell_{f,d} = \ell_f + \frac{l_{f,1}vd}{2}$, and $\ell_{F,d} = \ell_F + \frac{l_{f,1}v\sqrt{Md}}{2}$.

Proof of Lemma 6 By Lemma 2, for the $l_{f,1}$ -smooth function $f_{z,m}$ and any $x \in \mathbb{R}^d, i \in [n]$, define $s_z(x, u_j) \in [0, 1]$ and $s_z = [s_z(x, u_1), s_z(x, u_2), \dots, s_z(x, u_d)]$, $s = [s_1; s_2; \dots; s_m]$. By Lemma 1, $\|\nabla F(x_t)\lambda\| \leq \ell_f$ and $\|\nabla F(x_t)\| \leq \ell_F$. Then the estimator satisfies:

$$\begin{aligned}
 \|\hat{\nabla} F(x)\lambda\| &= \left\| \nabla F(x)\lambda + \frac{l_{f,1}v}{2}s(x, u)\lambda \right\| \\
 &\leq \|\nabla F(x_t)\lambda\| + \left\| \frac{l_{f,1}v}{2}s(x, u)\lambda \right\| \\
 &\leq \ell_f + \frac{l_{f,1}vd}{2}
 \end{aligned}$$

and

$$\begin{aligned}\|\hat{\nabla} F(x)\| &= \left\| \nabla F(x) + \frac{l_{f,1}v}{2} s(x, u) \right\| \\ &\leq \|\nabla F(x_t)\| + \left\| \frac{l_{f,1}v}{2} s(x, u) \right\| \\ &\leq \sqrt{M} \ell_f + \frac{l_{f,1}v\sqrt{Md}}{2} = \ell_F + \frac{l_{f,1}v\sqrt{Md}}{2}\end{aligned}$$

Defining $\ell_{f,d} = \ell_f + \frac{l_{f,1}vd}{2}$, and $\ell_{F,d} = \ell_F + \frac{l_{f,1}v\sqrt{Md}}{2}$, then it holds for all $\lambda \in \Delta^M$

$$\|\hat{\nabla} F(x_t) \lambda\| \leq \ell_{f,d} \quad \text{and} \quad \|\hat{\nabla} F(x_t)\| \leq \|\nabla F(x_t)\|_F \leq \ell_{F,d}$$

C BOUNDING THE OPTIMIZATION ERROR

Lemma 7. Suppose Assumption 1 holds. Consider the sequence $\{x_t\}, \{\lambda_1\}$ generated by SZMOD in unbounded domain for x . Define

$$\begin{aligned}S_{1,T} &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left\| \nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t) \lambda_t \right\|^2 \\ S_{3,T} &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left\| \nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t) \lambda_t \right\| \left\| \nabla F_S(x_t)^\top \nabla F_S(x_t) \lambda_1 \right\| \\ S_{4,T} &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left\| \hat{\nabla} F_{z_{t,3}}(x_t) \lambda_{t+1} \right\|^2\end{aligned}$$

Then it holds that

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left\| \nabla F_S(x_t) \lambda_t^*(x_t) \right\|^2 \leq \frac{1}{2\alpha T} \mathbb{E}_A [F_S(x_1) - F_S(x_{T+1})] \lambda_1 + \frac{1}{2} \gamma S_{1,T} + \gamma S_{3,T} + \frac{1}{2} \alpha \ell_{f,1} S_{4,T} + e.$$

Proof. By the $\ell_{f,1}$ -Lipschitz smoothness of $F_S(x)\lambda$ for all $\lambda \in \Delta^M$, we have

$$\begin{aligned}F_S(x_{t+1}) \lambda - F_S(x_t) \lambda &\leq (\nabla F_S(x_t) \lambda, x_{t+1} - x_t) + \frac{\ell_{f,1}}{2} \|x_{t+1} - x_t\|^2 \\ &= -\alpha_t \langle \nabla F_S(x_t) \lambda, \nabla F_{z_{t,3}}(x_t) \lambda_{t+1} \rangle + \frac{\ell_{f,1}}{2} \alpha_t^2 \left\| \hat{\nabla} F_{z_{t,3}}(x_t) \lambda_{t+1} \right\|^2.\end{aligned}$$

Taking expectation over $z_{\ell,3}$ on both sides of the above inequality gives

$$\mathbb{E}_{x_{z_{t,3}}, \mathbf{u}} [F_S(x_{t+1})] \lambda - F_S(x_t) \lambda \leq -\alpha_t \left\langle \nabla F_S(x_t) \lambda, \left(\nabla F_S(x_t) + \frac{Lv}{2} s(x, u) \right) \lambda_{t+1} \right\rangle + \frac{\ell_{f,1}}{2} \alpha_t^2 \mathbb{E}_{x_{z_{t,3}}, \mathbf{u}} \left\| \hat{\nabla} F_{z_{t,3}}(x_t) \lambda_{t+1} \right\|^2$$

By Lemma 8, we have

$$\begin{aligned}&2\gamma_t \mathbb{E}_A \left(\lambda_t - \lambda, \left(\nabla F_S(x_t)^\top \nabla F_S(x_t) \right) \lambda_t \right) \\ &\leq \mathbb{E}_A \left\| \lambda_t - \lambda \right\|^2 - \mathbb{E}_A \left\| \lambda_{t+1} - \lambda \right\|^2 + \gamma_t^2 \mathbb{E}_A \left\| \left(\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t) \right) \lambda_t \right\|^2 + e.\end{aligned}$$

Rearranging the above inequality and letting $\gamma_t = \gamma > 0$ gives

$$\begin{aligned} -\mathbb{E}_A \left\langle \lambda, \nabla F_S(x_t)^\top \nabla F_S(x_t) \lambda_t \right\rangle &\leq -\mathbb{E}_A \lambda_t, \left(\nabla F_S(x_t)^\top \nabla F_S(x_t) \right) \lambda_t \Big\rangle + \frac{1}{2\gamma} \mathbb{E}_A \left(\|\lambda_t - \lambda\|^2 - \|\lambda_{t+1} - \lambda\|^2 \right) \\ &\quad + \frac{1}{2} \gamma \mathbb{E}_A \left\| \left(\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t) \right) \lambda_t \right\|^2 + e \\ &\leq -\mathbb{E}_A \|\nabla F_S(x_t) \lambda_t\|^2 + \frac{1}{2\gamma} \mathbb{E}_A \left(\|\lambda_t - \lambda\|^2 - \|\lambda_{t+1} - \lambda\|^2 \right) \\ &\quad + \frac{1}{2} \gamma \mathbb{E}_A \left\| \left(\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t) \right) \lambda_t \right\|^2 + e \end{aligned}$$

Plugging the above inequality into, and setting $\alpha_t = \alpha > 0$, we have

$$\begin{aligned} \mathbb{E}_A [F_S(x_{t+1}) \lambda - F_S(x_t) \lambda] &\leq -\alpha \mathbb{E}_A \left\langle \nabla F_S(x_t) \lambda, \nabla F_S(x_t) \lambda_{t+1} + \frac{Lv}{2} (x, u) \lambda_{t+1} \right\rangle \\ &\leq -\alpha \mathbb{E}_A \|\nabla F_S(x_t) \lambda_t\|^2 + \frac{\alpha}{2\gamma} \mathbb{E}_A \left[\|\lambda_t - \lambda\|^2 - \|\lambda_{t+1} - \lambda\|^2 \right] \\ &\quad + \frac{L^2 v_i^2(x, u)}{8} \|\nabla_t + 1 - \nabla\| + \alpha \mathbb{E}_A \langle \nabla F_S(x_t) \lambda, \nabla F_S(x_t) (\lambda_t - \lambda_{t+1}) \rangle \\ &\quad + \frac{1}{2} \alpha^2 \ell_{f,1} \mathbb{E}_A \left\| \hat{\nabla} F_{z_{t,3}}(x_t) \lambda_{t+1} \right\|^2 + \frac{\ell_{f,1}}{2} \alpha^2 \mathbb{E}_A \left\| \hat{\nabla} F_{z_{t,3}}(x_t) \lambda_{t+1} \right\|^2 \\ &\quad + \frac{1}{2} \alpha \mathbb{E}_A \left\| \left(\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t) \right) \lambda_t \right\|^2 + e \end{aligned}$$

Taking telescope sum and rearranging yields, for all $\lambda \in \Delta^M$,

$$\begin{aligned} &\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \|\nabla F_S(x_t) \lambda_t\|^2 \\ &\leq \frac{1}{2\gamma T} \sum_{t=0}^{T-1} \mathbb{E}_A \left[\|\lambda_t - \lambda\|^2 - \|\lambda_{t+1} - \lambda\|^2 \right] + \frac{1}{\alpha T} \sum_{t=0}^{T-1} \mathbb{E}_A [F_S(x_t) - F_S(x_{t+1})] \lambda \\ &\quad + \frac{1}{2T} \sum_{t=0}^{T-1} \left(\gamma \mathbb{E}_A \left\| \nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t) \lambda_t \right\|^2 + \alpha \ell_{f,1} \mathbb{E}_A \left\| \nabla F_{z_{t,3}}(x_t) \lambda_{t+1} \right\|^2 + e \right. \\ &\quad \left. + 2\mathbb{E}_A \langle \nabla F_S(x_t) \lambda, \nabla F_S(x_t) (\lambda_t - \lambda_{t+1}) \rangle \right) \\ &\leq \frac{1}{2\gamma T} \mathbb{E}_A \left[\|\lambda_1 - \lambda\|^2 - \|\lambda_{T+1} - \lambda\|^2 \right] + \frac{1}{\alpha T} \mathbb{E}_A [F_S(x_1) - F_S(x_{T+1})] \lambda + \frac{1}{2} \gamma S_{1,T} + \gamma S_{3,T} + \frac{1}{2} \alpha \ell_{f,1} S_{4,T} + e. \end{aligned}$$

Setting $\lambda = \lambda_1$ in the above inequality yields

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \|\nabla F_S(x_t) \lambda_t\|^2 \leq \frac{1}{\alpha T} \mathbb{E}_A [F_S(x_1) - F_S(x_{T+1})] \lambda_1 + \frac{1}{2} \gamma S_{1,T} + \gamma S_{3,T} + \frac{1}{2} \alpha \ell_{f,1} S_{4,T} + e$$

Finally, the results follow from the definition of $\lambda_t^*(x_t)$.

Proof of Theorem 1: Then we proceed to bound $S_{1,T}, S_{3,T}, S_{4,T}$. Under either Assumptions 1, 2, or Assumptions 2, 3 with ℓ_f, ℓ_F defined in Lemma 1, we have that for all $z \in S$ and $\lambda \in \Delta^M$, $\|\nabla F_z(x_t) \lambda\| \leq \ell_f$, and $\|\nabla F_z(x_t)\| \leq \ell_F$. Then $S_{1,T}, S_{3,T}, S_{4,T}$ can be bounded below

$$\begin{aligned} S_{1,T} &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left\| \left(\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t) \right) \lambda_t \right\|^2 \leq M \ell_f^4 \\ S_{3,T} &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left\| \nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t) \lambda_t \right\| \left\| \nabla F_S(x_t)^\top \nabla F_S(x_t) \lambda_1 \right\| \leq \ell_F^2 \ell_f^2 = M \ell_f^4 \\ S_{4,T} &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left\| \hat{\nabla} F_{z_{t,3}}(x_t) \lambda_{t+1} \right\|^2 \leq \ell_{f,d}^2 \end{aligned}$$

which proves that

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \|\nabla F_S(x_t) \lambda_t^*(x_t)\|^2 \leq \frac{1}{\alpha T} c_F + \frac{3}{2} \gamma M \ell_f^4 + \frac{1}{2} \alpha \ell_{f,1} \ell_{f,d}^2 + e$$

We arrive at the results by $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \|\nabla F_S(x_t) \lambda_t^*(x_t)\| \leq \left(\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \|\nabla F_S(x_t) \lambda_t^*(x_t)\|^2 \right)^{\frac{1}{2}}$ from the Jensen's inequality and the convexity of the square function, as well as the subadditivity of square root function.

C.1 CA DIRECTION

C.2 LEMMAS

Lemma 8. Consider $\{x_t\}, \{\lambda_t\}$ generated by the SZMOD algorithm. For all $\lambda \in \Delta^M$, it holds that

$$\begin{aligned} & 2\gamma_t \mathbb{E}_A \left\langle \lambda_t - \lambda, \left(\nabla F_S(x_t)^\top \nabla F_S(x_t) \right) \lambda_t \right\rangle \\ & \leq \mathbb{E}_A \|\lambda_t - \lambda\|^2 - \mathbb{E}_A \|\lambda_{t+1} - \lambda\|^2 + \gamma_t^2 \mathbb{E}_A \left\| \left(\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t) \right) \lambda_t \right\|^2, \quad (19) \\ & + \frac{l_{f,1}^2 v^2 d}{4} \mathbb{E}_A \|\lambda_t - \lambda\|_1 + \frac{l_{f,1} v}{2} \mathbb{E}_A (\|\lambda_t - \lambda\|_1 \|\nabla F_S \lambda\|_1 + d \|\nabla F_S(\lambda_t - \lambda)\|_1) \end{aligned}$$

$$\begin{aligned} & \text{and } \gamma_t \mathbb{E}_A \left(\|\nabla F_S(x_t) \lambda_t\|^2 - \|\nabla F_S(x_t) \lambda\|^2 \right) \\ & \leq \mathbb{E}_A \|\lambda_t - \lambda\|^2 - \mathbb{E}_A \|\lambda_{t+1} - \lambda\|^2 + \gamma_t^2 \mathbb{E}_A \left\| \left(\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t) \right) \lambda_t \right\|^2 \quad (20) \\ & + \frac{l_{f,1}^2 v^2 d}{4} \mathbb{E}_A \|\lambda_t - \lambda\|_1 + \frac{l_{f,1} v}{2} \mathbb{E}_A (\|\lambda_t - \lambda\|_1 \|\nabla F_S \lambda\|_1 + d \|\nabla F_S(\lambda_t - \lambda)\|_1) \end{aligned}$$

Proof. By the update of λ , for all $\lambda \in \Delta^M$, we have

$$\begin{aligned} & \|\lambda_{t+1} - \lambda\|^2 \\ & = \left\| \Pi_{\Delta^M} \left(\lambda_t - \gamma_t \left(\hat{\nabla} F_{z_{t,1}}(x_t)^\top \hat{\nabla} F_{z_{t,2}}(x_t) \right) \lambda_t \right) - \lambda \right\|^2 \\ & \leq \left\| \lambda_t - \gamma_t \left(\hat{\nabla} F_{z_{t,1}}(x_t)^\top \hat{\nabla} F_{z_{t,2}}(x_t) \right) \lambda_t - \lambda \right\|^2 \\ & = \|\lambda_t - \lambda\|^2 - 2\gamma_t \left\langle \lambda_t - \lambda, \left(\hat{\nabla} F_{z_{t,1}}(x_t)^\top \hat{\nabla} F_{z_{t,2}}(x_t) \right) \lambda_t \right\rangle + \gamma_t^2 \left\| \left(\hat{\nabla} F_{z_{t,1}}(x_t)^\top \hat{\nabla} F_{z_{t,2}}(x_t) \right) \lambda_t \right\|^2 \quad (21) \end{aligned}$$

Now we focus on $\langle \lambda_t - \lambda, (\hat{\nabla} F_{z_{t,1}}(x_t)^\top \hat{\nabla} F_{z_{t,2}}(x_t)) \lambda_t \rangle$, we have:

$$\begin{aligned}
& \langle \lambda_t - \lambda, (\hat{\nabla} F_{z_{t,1}}(x_t)^\top \hat{\nabla} F_{z_{t,2}}(x_t)) \lambda_t \rangle \\
&= \langle \lambda_t - \lambda, (\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t)) \lambda_t \rangle + \frac{l_{f,1}^2 v^2}{4} (\lambda_t - \lambda)^\top s_{z_{t,1}}(x, u)^\top s_{z_{t,2}}(x, u) \lambda \\
&+ \frac{l_{f,1} v}{2} ((\lambda_t - \lambda)^\top s_{z_{t,1}}(x, u)^\top \nabla f_{z,2} \lambda + (\lambda_t - \lambda)^\top \nabla f_{z_{t,1}}^\top s_{z_{t,2}}(x, u) \lambda) \\
&\geq \langle \lambda_t - \lambda, (\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t)) \lambda_t \rangle - \frac{l_{f,1}^2 v^2 d}{4} \|\lambda_t - \lambda\|_1 \\
&- \frac{l_{f,1} v}{2} (\|\lambda_t - \lambda\|_1 \mathbb{1}_d^\top \nabla f_{z,2} \lambda + (\lambda_t - \lambda)^\top \nabla f_{z_{t,1}}^\top \mathbb{1}_d) \\
&\geq \langle \lambda_t - \lambda, (\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t)) \lambda_t \rangle - \frac{l_{f,1}^2 v^2 d}{4} \|\lambda_t - \lambda\|_1 \\
&- \frac{l_{f,1} v}{2} (\|\lambda_t - \lambda\|_1 \|\nabla f_{z,2} \lambda\|_1 + d \|(\lambda_t - \lambda) \nabla f_{z_{t,1}}\|_1)
\end{aligned} \tag{22}$$

Taking expectation over $z_{t,1}, z_{t,2}$ on both sides and rearranging proves equation 19. By the convexity of the problem, $\min_{\lambda \in \Delta^M} \frac{1}{2} \|\nabla F_S(x_t) \lambda\|^2$, we have

$$\begin{aligned}
& \gamma_t \mathbb{E}_A \left(\|\nabla F_S(x_t) \lambda_t\|^2 - \|\nabla F_S(x_t) \lambda\|^2 \right) \\
&\leq 2\gamma_t \mathbb{E}_A \left\langle \lambda_t - \lambda, (\nabla F_S(x_t)^\top \nabla F_S(x_t)) \lambda_t \right\rangle \\
&\stackrel{(22)}{\leq} 2\gamma_t \mathbb{E}_A \left\langle \lambda_t - \lambda, (\hat{\nabla} F_S(x_t)^\top \hat{\nabla} F_S(x_t)) \lambda_t \right\rangle + \frac{l_{f,1}^2 v^2 d}{4} \|\lambda_t - \lambda\|_1 \\
&+ \frac{l_{f,1} v}{2} (\|\lambda_t - \lambda\|_1 \|\nabla f_{z,2} \lambda\|_1 + d \|(\lambda_t - \lambda) \nabla f_{z_{t,1}}\|_1) \\
&\stackrel{(21)}{\leq} \mathbb{E}_A \|\lambda_t - \lambda\|^2 - \mathbb{E}_A \|\lambda_{t+1} - \lambda\|^2 + \gamma_t^2 \mathbb{E}_A \left\| (\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t)) \lambda_t \right\|^2 \\
&+ \frac{l_{f,1}^2 v^2 d}{4} \mathbb{E}_A \|\lambda_t - \lambda\|_1 + \frac{l_{f,1} v}{2} \mathbb{E}_A (\|\lambda_t - \lambda\|_1 \|\nabla F_S \lambda\|_1 + d \|\nabla F_S(\lambda_t - \lambda)\|_1)
\end{aligned}$$

Rearranging the above inequality proves equation 20

Lemma 9. Given any $\rho > 0$ and $x \in \mathbb{R}^d$, define $\lambda_\rho^*(x) = \arg \min_{\lambda \in \Delta^M} \frac{1}{2} \|\nabla F_S(x) \lambda\|^2 + \frac{1}{2} \rho \|\lambda\|^2$, then the following inequality holds

$$\|\lambda_\rho^*(x) - \lambda_\rho^*(x')\| \leq \rho^{-1} \left\| \nabla F(x)^\top \nabla F(x) - \nabla F(x')^\top \nabla F(x') \right\|$$

Suppose either 1) Assumptions 1, 3 hold, or 2) Assumptions 1, 2 hold, with ℓ_F defined in Lemma 1. Then for $x \in \{x_t\}_{t=1}^T, x' \in \{x'_t\}_{t=1}^T$ generated by MoDo algorithm on training dataset S and S' , respectively, it implies that

$$\|\lambda_\rho^*(x) - \lambda_\rho^*(x')\| \leq 2\rho^{-1} \ell_{F,1} \ell_F \|x - x'\|$$

Proofed in (Chen et al., 2024).

Lemma 10. Suppose Assumption 2 holds. Let $\{x_t\}, \{\lambda_t\}$ be the sequences produced by the SZMOD algorithm. With a positive constant $\bar{\rho} > 0$, define

$$\begin{aligned}
S_{1,T} &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left\| (\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t)) \lambda_t \right\|^2 \\
S_{2,T} &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \|\nabla F_S(x_{t+1}) + \nabla F_S(x_t)\| \|\nabla F_{z_{t,3}} \lambda_{t+1}\|
\end{aligned}$$

Then it holds that

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left[\|\nabla F_S(x_t) \lambda_t\|^2 - \|\nabla F_S(x_t) \lambda^*(x_t)\|^2 \right] \leq \bar{\rho} + \frac{4}{\gamma T} (1 + \bar{\rho}^{-1} \alpha \ell_{F,1} T S_{2,T}) + \gamma S_{1,T}$$

Proof: Define $\lambda_{\bar{\rho}}^*(x_t) = \arg \min_{\lambda \in \Delta^M} \frac{1}{2} \|\nabla F_S(x_t) \lambda\|^2 + \frac{\bar{\rho}}{2} \|\lambda\|^2$ with $\bar{\rho} > 0$. Note that $\bar{\rho}$ is strictly positive and is used only for analysis but not for algorithm update. Substituting $\lambda = \lambda_{\bar{\rho}}^*(x_t)$ in Lemma 8, we have

$$\begin{aligned} & \gamma_t \mathbb{E}_A \left(\|\nabla F_S(x_t) \lambda_t\|^2 - \|\nabla F_S(x_t) \lambda_{\bar{\rho}}^*(x_t)\|^2 \right) \\ & \leq \mathbb{E}_A \|\lambda_t - \lambda_{\bar{\rho}}^*(x_t)\|^2 - \mathbb{E}_A \|\lambda_{t+1} - \lambda_{\bar{\rho}}^*(x_t)\|^2 + \gamma_t^2 \mathbb{E}_A \left\| \left(\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t) \right) \lambda_t \right\|^2 \\ & \quad + \underbrace{\frac{l_{f,1}^2 v^2 d}{4} \mathbb{E}_A \|\lambda_t - \lambda\|_1 + \frac{l_{f,1} v}{2} \mathbb{E}_A (\|\lambda_t - \lambda\|_1 \|\nabla F_S \lambda\|_1 + d \|\nabla F_S(\lambda_t - \lambda)\|_1)}_e \end{aligned}$$

Setting $\gamma_t = \gamma > 0$, taking expectation and telescoping the above inequality gives

$$\begin{aligned} & \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left[\|\nabla F_S(x_t) \lambda_t\|^2 - \|\nabla F_S(x_t) \lambda_{\bar{\rho}}^*(x_t)\|^2 \right] \\ & \leq \frac{1}{T} \sum_{t=0}^{T-1} \frac{1}{\gamma} \mathbb{E}_A \left[\|\lambda_t - \lambda_{\bar{\rho}}^*(x_t)\|^2 - \|\lambda_{t+1} - \lambda_{\bar{\rho}}^*(x_t)\|^2 \right] + \frac{1}{T} \sum_{t=0}^{T-1} \gamma \mathbb{E}_A \left\| \left(\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t) \right) \lambda_t \right\|^2 + e \\ & = \frac{1}{\gamma T} \underbrace{\left(\sum_{t=0}^{T-1} \mathbb{E}_A \left[\|\lambda_t - \lambda_{\bar{\rho}}^*(x_t)\|^2 - \|\lambda_{t+1} - \lambda_{\bar{\rho}}^*(x_t)\|^2 \right] \right)}_{I_1} + \frac{1}{T} \sum_{t=0}^{T-1} \gamma \mathbb{E}_A \left\| \left(\nabla F_{z_{t,1}}(x_t)^\top \nabla F_{z_{t,2}}(x_t) \right) \lambda_t \right\|^2 + e \end{aligned} \tag{23}$$

where I_1 can be further derived as

$$\begin{aligned} I_1 &= \sum_{t=0}^{T-1} \mathbb{E}_A \|\lambda_t - \lambda_{\bar{\rho}}^*(x_t)\|^2 - \mathbb{E}_A \|\lambda_{t+1} - \lambda_{\bar{\rho}}^*(x_t)\|^2 \\ &= \mathbb{E}_A \|\lambda_0 - \lambda_{\bar{\rho}}^*(x_0)\|^2 - \mathbb{E}_A \|\lambda_T - \lambda_{\bar{\rho}}^*(x_T)\|^2 + \sum_{t=0}^{T-2} \mathbb{E}_A \left[\|\lambda_{t+1} - \lambda_{\bar{\rho}}^*(x_{t+1})\|^2 - \|\lambda_{t+1} - \lambda_{\bar{\rho}}^*(x_t)\|^2 \right] \\ &\leq \mathbb{E}_A \|\lambda_0 - \lambda_{\bar{\rho}}^*(x_0)\|^2 - \mathbb{E}_A \|\lambda_T - \lambda_{\bar{\rho}}^*(x_T)\|^2 \\ &\quad + \sum_{t=0}^{T-2} \mathbb{E}_A \left[\|2\lambda_{t+1} - \lambda_{\bar{\rho}}^*(x_{t+1}) - \lambda_{\bar{\rho}}^*(x_t)\| \|\lambda_{\bar{\rho}}^*(x_{t+1}) - \lambda_{\bar{\rho}}^*(x_t)\| \right] \\ &\leq 4 + 4 \sum_{t=0}^{T-2} \mathbb{E}_A \|\lambda_{\bar{\rho}}^*(x_{t+1}) - \lambda_{\bar{\rho}}^*(x_t)\| \end{aligned}$$

where $\|\lambda_{\bar{\rho}}^*(x_{t+1}) - \lambda_{\bar{\rho}}^*(x_t)\|$, by Lemma 9, can be bounded by

$$\begin{aligned} \|\lambda_{\bar{\rho}}^*(x_{t+1}) - \lambda_{\bar{\rho}}^*(x_t)\| &\leq \bar{\rho}^{-1} \|\nabla F_S(x_{t+1}) + \nabla F_S(x_t)\| \|\nabla F_S(x_{t+1}) - \nabla F_S(x_t)\| \\ &\leq \bar{\rho}^{-1} \ell_{F,1} \|\nabla F_S(x_{t+1}) + \nabla F_S(x_t)\| \|x_{t+1} - x_t\| \\ &\leq \bar{\rho}^{-1} \alpha \ell_{F,1} \|\nabla F_S(x_{t+1}) + \nabla F_S(x_t)\| \|\nabla F_{z_{t,3}} \lambda_{t+1}\| \end{aligned}$$

Hence, it follows that.

$$\begin{aligned} I_1 &\leq 4 + 4\bar{\rho}^{-1}\alpha\ell_{F,1} \sum_{t=0}^{T-1} \mathbb{E}_A \|\nabla F_S(x_{t+1}) + \nabla F_S(x_t)\| \|\nabla F_{z_t,3}\lambda_{t+1}\| \\ &= 4 + 4\bar{\rho}^{-1}\alpha\ell_{F,1}TS_{2,T} \end{aligned}$$

plugging which into (23) gives

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left[\|\nabla F_S(x_t)\lambda_t\|^2 - \|\nabla F_S(x_t)\lambda_{\bar{\rho}}^*(x_t)\|^2 \right] \leq \frac{4}{\gamma T} (1 + \bar{\rho}^{-1}\alpha\ell_{F,1}TS_{2,T}) + \gamma S_{1,T} + e \quad (24)$$

Define $\lambda^*(x_t) \in \arg \min_{\lambda \in \Delta^M} \|\nabla F_S(x_t)\lambda\|^2$. Then

$$\begin{aligned} &\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left[\|\nabla F_S(x_t)\lambda_t\|^2 - \|\nabla F_S(x_t)\lambda^*(x_t)\|^2 \right] \\ &= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left[\|\nabla F_S(x_t)\lambda_t\|^2 - \|\nabla F_S(x_t)\lambda_{\bar{\rho}}^*(x_t)\|^2 + \|\nabla F_S(x_t)\lambda_{\bar{\rho}}^*(x_t)\|^2 - \|\nabla F_S(x_t)\lambda^*(x_t)\|^2 \right] + e \\ &\stackrel{(24)}{\leq} \frac{4}{\gamma T} (1 + \bar{\rho}^{-1}\alpha\ell_{F,1}TS_{2,T}) + \gamma S_{1,T} + \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_A \left[\|\nabla F_S(x_t)\lambda_{\bar{\rho}}^*(x_t)\|^2 - \|\nabla F_S(x_t)\lambda^*(x_t)\|^2 \right] + e \\ &\leq \frac{4}{\gamma T} (1 + \bar{\rho}^{-1}\alpha\ell_{F,1}TS_{2,T}) + \gamma S_{1,T} + \bar{\rho} + e \end{aligned}$$

The proof is complete.

D BOUNDING THE GENERALIZATION ERROR

D.1 PROOF OF THEOREM 3-PS GENERALIZATION ERROR IN NONCONVEX CASE

In this subsection, we prove Theorem 3, which establishes the PS generalization error of SZMOD in the nonconvex case.

Organization of proof. To prove the PS generalization error of SZMOD, we first define the concept of Sampling-determined algorithms in Definition 3. This concept has been described in [22] for the analysis of single-objective learning. Then, we show that SZMOD is sampling-determined in Proposition 4. Finally, combining Propositions 2-4, we can prove Theorem 1, the MOL uniform stability and PS generalization error of SZMOD.

Definition 3 (Sampling-determined algorithm). *Let A be a randomized algorithm that randomly chooses an index sequence $I(A) = \{i_{t,s}\}$ to compute stochastic gradients. We say a symmetric algorithm A is sampling-determined if the output model is fully determined by $\{z_i : i \in I(A)\}$.*

Proposition 4 (SZMOD is sampling determined (Lei, 2023)). *SZMOD (Algorithm 3) is sampling determined. In other words, Let $I(A) = \{i_t\}$ be the sequence of index chosen by these algorithms from training set $S = \{z_1, \dots, z_n\}$, and $z_i \stackrel{i.i.d.}{\sim} \mathcal{P}$ for all $i \in [n]$ to build stochastic gradients, the output $A(S)$ is determined by $\{z_j : j \in I(A)\}$. To be precise, $A(S)$ is independent of z_j if $j \notin I(A)$.*

Proof of Proposition 4. Let $I(A) = \{I_1, \dots, I_T\}$, $I_t = \{i_{t,s}\}_{s=1}^3$ and $i_{t,s} \in [n]$ for all $1 \leq t \leq T$. And $S_{I(A)} = \{z_{i_{t,s}}\}$. By the description in Algorithm 1, $A(S) = G_{z_{I_T}} \circ \dots \circ G_{z_{I_1}}(x_0)$, where $G_z(\cdot)$ is the stochastic update function of the model parameter given random sample z . Therefore, for all possible sample realization z , we have

$$\begin{aligned}
\mathbb{P}(A(S) = x \mid z_j = z, j \notin I(A)) &= \mathbb{P}\left(G_{z_{I_T}} \circ \cdots \circ G_{z_{I_1}}(x_0) = x \mid z_j = z, j \notin I(A)\right) \\
&= \mathbb{P}\left(G_{z_{I_T}} \circ \cdots \circ G_{z_{I_1}}(x_0) = x \mid j \notin I(A)\right) \\
&= \mathbb{P}(A(S) = x \mid j \notin I(A))
\end{aligned}$$

where the last equality holds because $z_j \notin S_{I(A)}$, and z_j is independent of all elements in $S_{I(A)}$ by i.i.d. sampling. Therefore, $A(S)$ is independent of z_j if $j \notin I(A)$. The proof is complete.

Note that, besides SZMOD, other popular stochastic randomized MTL algorithms such as MODO (Chen et al., 2024) and MOCO (Fernando et al., 2023) are also sampling-determined. Therefore, the result is also applicable to these algorithms.

Lemma 11 ((Lei, 2023), Theorem 5 (b)). . *Let A be a sampling-determined random algorithm (Definition 3) and S, S' be neighboring datasets with n data points that differ only in the i -th data point. If $\sup_z \mathbb{E}_A \left[\|\nabla F_z(A(S))\|_F^2 \mid i \in I(A) \right] \leq G^2$ for any S , then*

$$\sup_z \mathbb{E}_A \left[\|\nabla F_z(A(S)) - \nabla F_z(A(S'))\|_F^2 \right] \leq 4G^2 \cdot \mathbb{P}\{i \in I(A)\}$$

Proof of Theorem 3. From Proposition 4, algorithm A , SZMOD is sampling-determined. Then based on Lemma , its MOL uniform stability in Definition 2 can be bounded by

$$\epsilon_F^2 \leq 4G^2 \cdot \mathbb{P}\{i \in I(A)\} \quad (25)$$

Let i_t be the index of the sample selected by A at the t -th step, and i^* be the index of the data point that is different in S and S' . Then

$$\mathbb{P}\{i^* \in I(A)\} \leq \sum_{t=0}^{T-1} \mathbb{P}\{i_t = i^*\} \leq \frac{T}{n} \quad (26)$$

Combining equation 25 and equation 26 gives

$$\epsilon_F^2 \leq \frac{4G^2 T}{n}$$

Then based on Propositions 2-3, we have

$$\begin{aligned}
\mathbb{E}_{A,S} [R_{\text{gen}}(A(S))] &\leq \mathbb{E}_{A,S} [\|\nabla F(A(S)) - \nabla F_S(A(S))\|_F] \\
&\leq 4\epsilon_F + \sqrt{n^{-1} \mathbb{E}_S [\mathbb{V}_{z \sim \mathcal{D}} (\nabla F_z(A(S)))]} \\
&= \mathcal{O}\left(T^{\frac{1}{2}} n^{-\frac{1}{2}}\right)
\end{aligned}$$

The proof is complete.

D.2 EXPANSIVENESS AND BOUNDEDNESS OF SZMOD UPDATE

Lemma 12. [Boundedness of update function of SZMOD] *Let ℓ_f be a positive constant. If $\|\hat{\nabla} F_z(x)\lambda\| \leq \ell_f$ for all $\lambda \in \Delta^M$, $z \in S$ and $x \in \{x_t\}_{t=1}^T$ generated by the SZMOD algorithm with step size $\alpha_t \leq \alpha$, then $G_{\lambda,z}(x)$ is $(\alpha\ell_f)$ -bounded on the trajectory of SZMOD, i.e.,*

$$\sup_{x \in \{x_t\}_{t=1}^T} \|G_{\lambda,z}(x) - x\| \leq \alpha\ell_f$$

Proof. For all $x \in \{x_t\}_{t=1}^T$, $\lambda \in \Delta^M$, and $z \in S$, since $\|\hat{\nabla} F_z(x)\lambda\| \leq \ell_f$, we have

$$\|G_{\lambda,z}(x) - x\| \leq \|\alpha \hat{\nabla} F_z(x)\lambda\| \leq \alpha \ell_f$$

which proves the boundedness.

Lemma 13 (Properties of update function of SZMOD in convex case). *Suppose Assumptions 2, 3 hold. Let ℓ_f be a positive constant. If for all $\lambda, \lambda' \in \Delta^M$, $z \in S$, and $x \in \{x_t\}_{t=1}^T$, $x' \in \{x'_t\}_{t=1}^T$ generated by the SZMOD algorithm on datasets S and S' , respectively, we have $\|\nabla F_z(x)\lambda\| \leq \ell_f$, $\|\nabla F_z(x')\lambda'\| \leq \ell_f$, and $\|\nabla F_z(x)\| \leq \ell_F$, $\|\nabla F_z(x')\| \leq \ell_F$, we have $\|\hat{\nabla} F_z(x)\lambda\| \leq \ell_{f,d}$, $\|\hat{\nabla} F_z(x')\lambda'\| \leq \ell_{f,d}$, and $\|\hat{\nabla} F_z(x)\| \leq \ell_{F,d}$, $\|\hat{\nabla} F_z(x')\| \leq \ell_{F,d}$, and step sizes of SZMOD satisfy $\alpha_t \leq \alpha$, $\gamma_t \leq \gamma$, it holds that*

$$\begin{aligned} \|G_{\lambda,z}(x) - G_{\lambda',z}(x')\|^2 &\leq (1 - 2\alpha\mu + 2\alpha^2\ell_{f,1}^2 + \alpha\ell_{f,1}vd) \|x - x'\|^2 \\ &\quad + 2\alpha\ell_{F,d} \|x - x'\| \|\lambda' - \lambda\| + 2\alpha^2\ell_{F,d}^2 \|\lambda' - \lambda\|^2 \\ &\quad + \alpha\ell_{f,1}vd + 2\alpha^2\ell_{f,1}^2v^2d^2 \\ \|G_{x,z_1,z_2}(\lambda) - G_{x',z_1,z_2}(\lambda')\|^2 &\leq \left((1 + \ell_{F,d}^2\gamma)^2 + (1 + \ell_{F,d}^2\gamma) \ell_{g,1}\gamma + (1 + \ell_{F,d}^2\gamma) \ell_{g,1}vd\gamma \right) \|\lambda - \lambda'\|^2 \\ &\quad + \left((1 + \ell_{F,d}^2\gamma) \ell_{g,1}\gamma + \ell_{g,1}^2\gamma^2 + \ell_{g,1}^2vd\gamma \right) \|x - x'\|^2 \\ &\quad + (1 + \ell_{F,d}^2\gamma) \ell_{g,1}vd\gamma + \ell_{g,1}^2vd\gamma \end{aligned}$$

Proof. The squared norm of the difference of $G_{\lambda,z}(x)$ and $G_{\lambda',z}(x')$ can be bounded by

$$\begin{aligned} &\|G_{\lambda,z}(x) - G_{\lambda',z}(x')\|^2 \\ &= \|x - x'\|^2 - 2\alpha \left\langle x - x', \hat{\nabla} F_z(x)\lambda - \hat{\nabla} F_z(x')\lambda' \right\rangle + \alpha^2 \left\| \hat{\nabla} F_z(x)\lambda - \hat{\nabla} F_z(x')\lambda' \right\|^2 \\ &\stackrel{(a)}{\leq} \|x - x'\|^2 - 2\alpha \left\langle x - x', \left(\hat{\nabla} F_z(x) - \hat{\nabla} F_z(x') \right) \lambda \right\rangle + 2\alpha^2 \left\| \left(\hat{\nabla} F_z(x) - \hat{\nabla} F_z(x') \right) \lambda \right\|^2 \\ &\quad + 2\alpha \left\langle x - x', \hat{\nabla} F_z(x')(\lambda' - \lambda) \right\rangle + 2\alpha^2 \left\| \hat{\nabla} F_z(x')(\lambda - \lambda') \right\|^2 \\ &\stackrel{(b)}{=} \|x - x'\|^2 - 2\alpha \left\langle x - x', (\nabla F_z(x) - \nabla F_z(x')) \lambda + \left(\frac{\ell_{f,1}v}{2}s(x, u) - \frac{\ell_{f,1}v}{2}s(x', u) \right) \lambda \right\rangle \\ &\quad + 2\alpha^2 \left\| (\nabla F_z(x) - \nabla F_z(x')) \lambda + \left(\frac{\ell_{f,1}v}{2}s(x, u) - \frac{\ell_{f,1}v}{2}s(x', u) \right) \lambda \right\|^2 \\ &\quad + 2\alpha \left\langle x - x', \hat{\nabla} F_z(x')(\lambda' - \lambda) \right\rangle + 2\alpha^2 \left\| \hat{\nabla} F_z(x')(\lambda - \lambda') \right\|^2 \\ &\stackrel{(c)}{\leq} \|x - x'\|^2 - 2\alpha \left\langle x - x', (\nabla F_z(x) - \nabla F_z(x')) \lambda \right\rangle + 2\alpha^2 \left\| (\nabla F_z(x) - \nabla F_z(x')) \lambda \right\|^2 \\ &\quad - 2\alpha \left\langle x - x', \left(\frac{\ell_{f,1}v}{2}s(x, u) - \frac{\ell_{f,1}v}{2}s(x', u) \right) \lambda \right\rangle + 2\alpha^2 \left\| \left(\frac{\ell_{f,1}v}{2}s(x, u) - \frac{\ell_{f,1}v}{2}s(x', u) \right) \lambda \right\|^2 \\ &\quad + 2\alpha \left\langle x - x', \hat{\nabla} F_z(x')(\lambda' - \lambda) \right\rangle + 2\alpha^2 \left\| \hat{\nabla} F_z(x')(\lambda - \lambda') \right\|^2 \\ &\stackrel{(d)}{\leq} (1 - 2\alpha\mu + 2\alpha^2\ell_{f,1}^2) \|x - x'\|^2 + 2\alpha \left\langle x - x', \hat{\nabla} F_z(x')(\lambda' - \lambda) \right\rangle + 2\alpha^2\ell_{F,d}^2 \|\lambda' - \lambda\|^2 \\ &\quad - 2\alpha \left\langle x - x', \left(\frac{\ell_{f,1}v}{2}s(x, u) - \frac{\ell_{f,1}v}{2}s(x', u) \right) \lambda \right\rangle + 2\alpha^2 \left\| \left(\frac{\ell_{f,1}v}{2}s(x, u) - \frac{\ell_{f,1}v}{2}s(x', u) \right) \lambda \right\|^2 \\ &\stackrel{(e)}{\leq} (1 - 2\alpha\mu + 2\alpha^2\ell_{f,1}^2 + \alpha\ell_{f,1}vd) \|x - x'\|^2 + 2\alpha\ell_{F,d} \|x - x'\| \|\lambda' - \lambda\| + 2\alpha^2\ell_{F,d}^2 \|\lambda' - \lambda\|^2 \\ &\quad + \alpha\ell_{f,1}vd + 2\alpha^2\ell_{f,1}^2v^2d^2 \end{aligned}$$

where (a) follows from rearranging and that $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$; (b) follows from $\hat{\nabla} F(x) = \nabla F(x) + \frac{\ell_{f,1}v}{2}s(x, u)$; (c) follows from that $\|a + b\| \leq \|a\| + \|b\|$; (d) follows from the μ -strong convexity of $F_z(x)\lambda$, $\ell_{f,1}$ -Lipschitz continuity of $\nabla F_z(x)\lambda$, and that $\|\hat{\nabla} F_z(x')\| \leq \ell_{F,d}$ for $x' \in \{x'_t\}_{t=1}^T$; and, (e) follows from Cauchy-Schwartz inequality. And $\|G_{x,z_1,z_2}(\lambda) - G_{x',z_1,z_2}(\lambda')\|$ can be bounded by

$$\begin{aligned}
& \|G_{x,z_1,z_2}(\lambda) - G_{x',z_1,z_2}(\lambda')\| \\
&= \left\| \Pi_{\Delta^M} \left(\lambda - \gamma \left(\hat{\nabla} F_{z_1}(x)^\top \hat{\nabla} F_{z_2}(x) \right) \lambda \right) - \Pi_{\Delta^M} \left(\lambda' - \gamma \left(\hat{\nabla} F_{z_1}(x')^\top \hat{\nabla} F_{z_2}(x') \right) \lambda' \right) \right\| \\
&\stackrel{(f)}{\leq} \left\| \lambda - \lambda' - \gamma \left(\hat{\nabla} F_{z_1}(x)^\top \hat{\nabla} F_{z_2}(x) \lambda - \hat{\nabla} F_{z_1}(x')^\top \hat{\nabla} F_{z_2}(x') \lambda' \right) \right\| \\
&\stackrel{(g)}{\leq} \|\lambda - \lambda'\| + \gamma \left\| \hat{\nabla} F_{z_1}(x)^\top \hat{\nabla} F_{z_2}(x) (\lambda - \lambda') \right\| + \gamma \left\| \left(\hat{\nabla} F_{z_1}(x)^\top \hat{\nabla} F_{z_2}(x) - \hat{\nabla} F_{z_1}(x')^\top \hat{\nabla} F_{z_2}(x') \right) \lambda' \right\| \\
&\stackrel{(h)}{\leq} \|\lambda - \lambda'\| + \gamma \ell_{F,d}^2 \|\lambda - \lambda'\| + \gamma \left\| \left(\hat{\nabla} F_{z_1}(x)^\top \hat{\nabla} F_{z_2}(x) - \hat{\nabla} F_{z_1}(x')^\top \hat{\nabla} F_{z_2}(x') \right) \lambda' \right\| \\
&\stackrel{(i)}{\leq} \|\lambda - \lambda'\| + \gamma \ell_{F,d}^2 \|\lambda - \lambda'\| + \gamma \left(\left\| \left(\hat{\nabla} F_{z_1}(x) - \hat{\nabla} F_{z_1}(x') \right)^\top \hat{\nabla} F_{z_2}(x) \lambda' \right\| \right. \\
&\quad \left. + \left\| \hat{\nabla} F_{z_1}(x')^\top \left(\hat{\nabla} F_{z_2}(x) - \hat{\nabla} F_{z_2}(x') \right) \lambda' \right\| \right) \\
&\leq (1 + \ell_{F,d}^2 \gamma) \|\lambda - \lambda'\| + (\ell_{f,d} \ell_{F,1} + \ell_{F,d} \ell_{f,1}) \gamma \|x - x'\| + (\ell_{f,d} \ell_{F,1} + \ell_{F,d} \ell_{f,1}) \gamma v d
\end{aligned}$$

where (f) follows from non-expansiveness of projection; (g) follows from triangle inequality, (h) follows from $\|\hat{\nabla} F_z(x)\| \leq \ell_{F,d}$ for $x \in \{x'_t\}_{t=1}^T$, (i) follows from triangle inequality; and (j) follows from $\ell_{F,1}$ -Lipschitz continuity of $\nabla F_z(x)\lambda'$, $\ell_{f,1}$ Lipschitz continuity of $\nabla F_z(x)$, $\|\hat{\nabla} F_z(x)\| \leq \ell_{F,d}$ for $x \in \{x'_t\}_{t=1}^T$, and $\|\nabla F_z(x)\lambda'\| \leq \ell_f$ for $x \in \{x_t\}_{t=1}^T$. Let $\ell_{g,1} = \ell_{f,d} \ell_{F,1} + \ell_{F,d} \ell_{f,1}$. Taking square on both sides of above yields

$$\begin{aligned}
& \|G_{x,z_1,z_2}(\lambda) - G_{x',z_1,z_2}(\lambda')\|^2 \\
&\leq \left((1 + \ell_{F,d}^2 \gamma) \|\lambda - \lambda'\| + \ell_{g,1} \gamma \|x - x'\| + \gamma \ell_{g,1} v d \right)^2 \\
&= (1 + \ell_{F,d}^2 \gamma)^2 \|\lambda - \lambda'\|^2 + \ell_{g,1}^2 \gamma^2 \|x - x'\|^2 + \gamma^2 \ell_{g,1}^2 v^2 d^2 + 2(1 + \ell_{F,d}^2 \gamma) \ell_{g,1} \gamma \|\lambda - \lambda'\| \|x - x'\| \\
&\quad + 2(1 + \ell_{F,d}^2 \gamma) \ell_{g,1} v d \gamma \|\lambda - \lambda'\| + 2\ell_{g,1}^2 v d \gamma \|x - x'\| \\
&\leq (1 + \ell_{F,d}^2 \gamma)^2 \|\lambda - \lambda'\|^2 + \ell_{g,1}^2 \gamma^2 \|x - x'\|^2 + (1 + \ell_{F,d}^2 \gamma) \ell_{g,1} \gamma (\|\lambda - \lambda'\|^2 + \|x - x'\|^2) \\
&\quad + (1 + \ell_{F,d}^2 \gamma) \ell_{g,1} v d \gamma (\|\lambda - \lambda'\|^2 + 1) + \ell_{g,1}^2 v d \gamma (\|x - x'\|^2 + 1) \\
&= \left((1 + \ell_{F,d}^2 \gamma)^2 + (1 + \ell_{F,d}^2 \gamma) \ell_{g,1} \gamma + (1 + \ell_{F,d}^2 \gamma) \ell_{g,1} v d \gamma \right) \|\lambda - \lambda'\|^2 \\
&\quad + \left((1 + \ell_{F,d}^2 \gamma) \ell_{g,1} \gamma + \ell_{g,1}^2 \gamma^2 + \ell_{g,1}^2 v d \gamma \right) \|x - x'\|^2 + (1 + \ell_{F,d}^2 \gamma) \ell_{g,1} v d \gamma + \ell_{g,1}^2 v d \gamma
\end{aligned}$$

The proof is complete.

D.3 GROWTH RECURSION

Lemma 14 (Growth recursion with approximate expansiveness). *Fix an arbitrary sequence of updates G_1, \dots, G_T and another sequence G'_1, \dots, G'_T . Let $x_0 = x'_0$ be a starting point in Ω and define $\delta_t = \|x'_t - x_t\|$ where x_t, x'_t are defined recursively through*

$$x_{t+1} = G_t(x_t), \quad x'_{t+1} = G'_t(x'_t) \quad (t \geq 0)$$

Let $\eta_t > 0, \nu_t \geq 0$, and $s_t \geq 0$. Then, for any $p > 0$, and $t \in [T]$, we have the recurrence relation (with $\delta_0 = 0$)

$$\delta_{t+1}^2 \leq \begin{cases} \eta_t \delta_t^2 + \nu_t, & G_t = G'_t \text{ is } (\eta_t, \nu_t) \text{-approximately expansive in square} \\ (1+p) \min \{ \eta_t \delta_t^2 + \nu_t, \delta_t^2 \} + \left(1 + \frac{1}{p}\right) 4\varsigma_t^2 & G_t \text{ and } G'_t \text{ are } \varsigma_t \text{-bounded,} \\ & G_t \text{ is } (\eta_t, \nu_t) \text{-approximately expansive in square.} \end{cases}$$

Proof. When G_t and G'_t are ς_t -bounded, we can bound δ_{t+1} by

$$\begin{aligned} \delta_{t+1} &= \|x_{t+1} - x'_{t+1}\| = \|G_t(x_t) - G'_t(x'_t)\| \\ &= \|G_t(x_t) - x_t - G'_t(x'_t) + x'_t + x_t - x'_t\| \\ &\leq \|G_t(x_t) - x_t\| + \|G'_t(x'_t) - x'_t\| + \|x_t - x'_t\| \\ &\leq 2\varsigma_t + \delta_t \end{aligned}$$

Alternatively, when G_t and G'_t are ς_t -bounded, G_t is (η_t, ν_t) -approximately expansive, we have

$$\delta_{t+1} = \|x_{t+1} - x'_{t+1}\| = \|G_t(x_t) - G'_t(x'_t)\|$$

When $G_t = G'_t$, is (η_t, ν_t) -approximately expansive in square, given $\delta_t^2, \delta_{t+1}^2$ can be bounded by

$$\delta_{t+1}^2 = \|x_{t+1} - x'_{t+1}\|^2 = \|G_t(x_t) - G_t(x'_t)\|^2 \leq \eta_t \|x_t - x'_t\|^2 + \nu_t = \eta_t \delta_t^2 + \nu_t$$

When G_t and G'_t are ς_t -bounded, applying (B.57), we can bound δ_{t+1}^2 by

$$\delta_{t+1}^2 \leq (\delta_t + 2\varsigma_t)^2 \leq (1+p)\delta_t^2 + (1+1/p)4\varsigma_t^2$$

where $p > 0$ and the last inequality follows from $(a+b)^2 \leq (1+p)a^2 + (1+1/p)b^2$. Alternatively, when G_t and G'_t are ς_t -bounded, G_t is (η_t, ν_t) -approximately expansive in square, the following holds

$$\begin{aligned} \delta_{t+1}^2 &= \|x_{t+1} - x'_{t+1}\|^2 = \|G_t(x_t) - G'_t(x'_t)\|^2 \\ &= \|G_t(x_t) - G_t(x'_t) + G_t(x'_t) - G'_t(x'_t)\|^2 \\ &\leq (1+p) \|G_t(x_t) - G_t(x'_t)\|^2 + (1+1/p) \|G_t(x'_t) - G'_t(x'_t)\|^2 \\ &\leq (1+p) (\eta_t \delta_t^2 + \nu_t) + (1+1/p) \|G_t(x'_t) - x'_t - G'_t(x'_t) + x'_t\|^2 \\ &\leq (1+p) (\eta_t \delta_t^2 + \nu_t) + 2(1+1/p) \left(\|G_t(x'_t) - x'_t\|^2 + \|G'_t(x'_t) - x'_t\|^2 \right) \\ &\leq (1+p) (\eta_t \delta_t^2 + \nu_t) + (1+1/p) 4\varsigma_t^2 \end{aligned}$$

The proof is complete.

D.4 UPPER BOUND OF SZMOD UNIFORM STABILITY

In Theorem 4 we bound the argument stability, which is then used to derive the MOL uniform stability and PS generalization error in Theorem 4.

Theorem 7 (Argument stability bound in strongly convex case). *Suppose Assumptions 2, 3, hold. Let A be the SZMOD algorithm in Algorithm 3. Choose the step sizes $\alpha_t \leq \alpha \leq \min \left\{ 1/(2\ell_{f,1}), \mu/(2\ell_{f,1}^2) \right\}$, and $\gamma_t \leq \gamma \leq \min \left\{ \frac{\mu^2}{484\ell_{f,d}^2\ell_{g,1}}, \frac{1}{8(3\ell_f^2+2\ell_{g,1})} \right\} / T$. Then it holds that*

$$\mathbb{E}_A \left[\|A(S) - A(S')\|^2 \right] \leq \frac{48}{\mu n} \ell_{f,d}^2 \left(\alpha + \frac{12 + 4M\ell_{f,d}^2}{\mu n} + \frac{10M\ell_f^4 \gamma}{\mu} \right) + \frac{4}{\mu n} \left(\frac{10\alpha M\ell_{f,d}^2 \gamma + \mu \gamma}{\mu \alpha} + \alpha \ell_{f,1} + \frac{2\alpha \ell_{f,1}^2}{n} \right)$$

Proof of Theorem 7. Under Assumptions 2, 3, Lemma 6 implies that for $\{x_t\}$ generated by the

$$\|\hat{\nabla} F_z(x_t)\| \leq \ell_{f,d}. \quad \text{and} \quad \|\hat{\nabla} F_z(x_t)\| \leq \ell_{F,d}$$

For notation simplicity, denote $\delta_t = \|x_t - x'_t\|$, $\zeta_t = \|\lambda_t - \lambda'_t\|$, $x_T = A_T(S)$ and $x'_T = A_T(S')$. Denote the index of the different sample in S and S' as i^* , and the set of indices selected at the t -th iteration as I_t , i.e., $I_t = \{i_{t,s}\}_{s=1}^3$. When $i^* \notin I_t$, for any $c_1 > 0$, based on Lemma 13, we have

$$\begin{aligned} \delta_{t+1}^2 &\leq (1 - 2\alpha_t \mu + 2\alpha_t^2 \ell_{f,1}^2 + \alpha_t \ell_{f,1} v d) \delta_t^2 + 2\alpha_t \ell_{F,d} \delta_t \zeta_{t+1} + 2\alpha_t^2 \ell_{F,d}^2 \zeta_{t+1}^2 + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^d \\ &\leq (1 - 2\alpha_t \mu + 2\alpha_t^2 \ell_{f,1}^2 + \alpha_t \ell_{f,1} v d) \delta_t^2 + \alpha_t \ell_{F,d} (c_1 \delta_t^2 + c_1^{-1} \zeta_{t+1}^2) + 2\alpha_t^2 \ell_{F,d}^2 \zeta_{t+1}^2 + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^d \\ &\leq (1 - \alpha_t \mu + \alpha_t \ell_{f,1} v d) \delta_t^2 + \alpha_t \ell_{F,d} (c_1 \delta_t^2 + c_1^{-1} \zeta_{t+1}^2) + 2\alpha_t^2 \ell_{F,d}^2 \zeta_{t+1}^2 + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^d \end{aligned}$$

where the second last inequality is due to Young's inequality; the last inequality is due to choosing $\alpha_t \leq \mu / (2\ell_{f,1}^2)$. When $i^* \in I_t$, from Lemma 12, the $(\alpha_t \ell_{f,d})$ -boundedness of the update at t -th iteration, and Lemma 9, the growth recursion, for a given constant $p > 0$, we have

$$\delta_{t+1}^2 \leq (1 + p) \delta_t^2 + (1 + 1/p) 4\alpha_t^2 \ell_{f,d}^2$$

Taking expectation of δ_{t+1}^2 over I_t , we have

$$\begin{aligned} \mathbb{E}_{I_t} [\delta_{t+1}^2] &\leq \mathbb{P}(i^* \notin I_t) ((1 - \alpha_t \mu + \alpha_t \ell_{f,1} v d) \delta_t^2 + \alpha_t \ell_{F,d} c_1 \delta_t^2 + (\alpha_t \ell_{F,d} c_1^{-1} + 2\alpha_t^2 \ell_{F,d}^2) \mathbb{E}_{I_t} [\zeta_{t+1}^2 | i^* \notin I_t]) \\ &\quad + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^2 + \mathbb{P}(i^* \in I_t) ((1 + p) \delta_t^2 + (1 + 1/p) 4\alpha_t^2 \ell_{f,d}^2) \\ &\leq (1 - \alpha_t (\mu - \ell_{F,d} c_1 - \ell_{f,1} v d)) \mathbb{P}(i^* \notin I_t) + p \mathbb{P}(i^* \in I_t) \delta_t^2 \\ &\quad + \alpha_t \underbrace{(\ell_{F,d} c_1^{-1} + 2\alpha_t \ell_{F,d}^2)}_{c_2} \mathbb{E}_{I_t} [\zeta_{t+1}^2 | i^* \notin I_t] \mathbb{P}(i^* \notin I_t) + \left(1 + \frac{1}{p}\right) \mathbb{P}(i^* \in I_t) 4\alpha_t^2 \ell_{f,d}^2 \\ &\quad + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^2 \end{aligned} \tag{27}$$

At each iteration of SZMOD, we randomly select three independent samples (instead of one) from the training set S . Then the probability of selecting the different sample from S and S' at the t -th iteration, $\mathbb{P}(i^* \in I_t)$ in the above equation, can be computed as follows

$$\mathbb{P}(i^* \in I_t) = 1 - \left(\frac{n-1}{n} \right)^3 \leq \frac{3}{n}$$

Consequently, the probability of selecting the same sample from S and S' at the t -th iteration is $\mathbb{P}(i^* \notin I_t) = 1 - \mathbb{P}(i^* \in I_t)$. Let $\ell_{g,1} = \ell_{f,d} \ell_{F,1} + \ell_{F,d} \ell_{f,1}$. Recalling when $i^* \notin I_t$, $\zeta_{t+1} \leq (1 + \ell_{F,d}^2 \gamma) \zeta_t + \ell_{g,1} \gamma \delta_t + \gamma \ell_{g,1} v d$ from Lemma 8, it follows that

$$\begin{aligned} \zeta_{t+1}^2 &\leq \left((1 + \ell_{F,d}^2 \gamma)^2 + (1 + \ell_{F,d}^2 \gamma) \ell_{g,1} \gamma + (1 + \ell_{F,d}^2 \gamma) \ell_{g,1} v d \gamma \right) \zeta_t^2 \\ &\quad + ((1 + \ell_{F,d}^2 \gamma) \ell_{g,1} \gamma + \ell_{g,1}^2 \gamma^2 + \ell_{g,1}^2 v d \gamma) \delta_t^2 + (1 + \ell_{F,d}^2 \gamma) \ell_{g,1} v d \gamma + \ell_{g,1}^2 v d \gamma \\ &\leq (1 + \underbrace{(3\ell_{F,d}^2 + 4\ell_{g,1})}_{c_3} \gamma) \zeta_t^2 + 4\ell_{g,1} \gamma \delta_t^2 + (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma \end{aligned}$$

where the last inequality follows from $\ell_{g,1}\gamma_t \leq 1$, and $\ell_{F,d}^2\gamma_t \leq 1$. And since ζ_t and δ_t are independent of I_t , it follows that

$$\mathbb{E}_{I_t} [\zeta_{t+1}^2 \mid i^* \notin I_t] \leq (1 + c_3\gamma_t) \zeta_t^2 + 4\ell_{g,1}\gamma_t\delta_t^2 + (2\ell_{g,1} + \ell_{g,1}^2)vd\gamma_t \quad (28)$$

Combining equation 27 and equation 28, we have

$$\begin{aligned} \mathbb{E}_{I_t} [\delta_{t+1}^2] &\leq (1 - \alpha_t (\mu - \ell_{F,d}c_1 - \ell_{f,1}vd)) \mathbb{P}(i^* \notin I_t) + p\mathbb{P}(i^* \in I_t) \delta_t^2 + \left(1 + \frac{1}{p}\right) \mathbb{P}(i^* \in I_t) 4\alpha_t^2\ell_{f,d}^2 \\ &\quad + \alpha_t c_2 ((1 + c_3\gamma_t) \zeta_t^2 + 4\ell_{g,1}\gamma_t\delta_t^2 + (2\ell_{g,1} + \ell_{g,1}^2)vd\gamma_t) \mathbb{P}(i^* \notin I_t) + \alpha_t\ell_{f,1}vd + 2\alpha_t^2\ell_{f,1}^2v^2d^2 \\ &= (1 - \alpha_t (\mu - \ell_{F,d}c_1 - \ell_{f,1}vd - 4c_2\ell_{g,1}\gamma_t)) \mathbb{P}(i^* \notin I_t) + p\mathbb{P}(i^* \in I_t) \delta_t^2 + \left(1 + \frac{1}{p}\right) \mathbb{P}(i^* \in I_t) 4\alpha_t^2\ell_{f,d}^2 \\ &\quad + \alpha_t c_2 ((1 + c_3\gamma_t) \zeta_t^2 + (2\ell_{g,1} + \ell_{g,1}^2)vd\gamma_t) \mathbb{P}(i^* \notin I_t) + \alpha_t\ell_{f,1}vd + 2\alpha_t^2\ell_{f,1}^2v^2d^2 \\ &= (\eta_t \mathbb{P}(i^* \notin I_t) + p\mathbb{P}(i^* \in I_t)) \delta_t^2 + \left(1 + \frac{1}{p}\right) \mathbb{P}(i^* \in I_t) 4\alpha_t^2\ell_{f,d}^2 \\ &\quad + \alpha_t c_2 ((1 + c_3\gamma_t) \zeta_t^2 + (2\ell_{g,1} + \ell_{g,1}^2)vd\gamma_t) \mathbb{P}(i^* \notin I_t) + \alpha_t\ell_{f,1}vd + 2\alpha_t^2\ell_{f,1}^2v^2d^2 \end{aligned} \quad (29)$$

where we define $\eta_t = 1 - \alpha(\mu - \ell_{F,d}c_1 - \ell_{f,1}vd - 4c_2\ell_{g,1}\gamma_t)$

While when $i^* \in I_t$, for a given constant $p_2 > 0$, we have

$$\begin{aligned} \zeta_{t+1} &= \left\| \Pi_{\Delta^M} \left(\lambda_t - \gamma_t \hat{\nabla} F_{t,1}(x_t)^\top \hat{\nabla} F_{t,2}(x_t) \lambda_t \right) - \Pi_{\Delta^M} \left(\lambda'_t - \gamma_t \hat{\nabla} F_{t,1}(x'_t)^\top \hat{\nabla} F_{t,2}(x'_t) \lambda'_t \right) \right\| \\ &\leq \left\| \lambda_t - \lambda'_t - \gamma_t \left(\hat{\nabla} F_{t,1}(x_t)^\top \hat{\nabla} F_{t,2}(x_t) \lambda_t - \hat{\nabla} F'_{t,1}(x'_t)^\top \hat{\nabla} F'_{t,2}(x'_t) \lambda'_t \right) \right\| \\ &\leq \|\lambda_t - \lambda'_t\| + 2\gamma_t \ell_{F,d} \ell_{f,d} \leq \zeta_t + 2\gamma_t \sqrt{M} \ell_{f,d}^2 \\ \zeta_{t+1}^2 &\leq (1 + p_2) \zeta_t^2 + (1 + 1/p_2) 4\gamma_t^2 M \ell_{f,d}^4 \end{aligned}$$

Taking expectation of ζ_{t+1}^2 over I_t gives

$$\begin{aligned} \mathbb{E}_{I_t} [\zeta_{t+1}^2] &= \mathbb{E}_{I_t} [\zeta_{t+1}^2 \mid i^* \in I_t] \mathbb{P}(i^* \in I_t) + \mathbb{E}_{I_t} [\zeta_{t+1}^2 \mid i^* \notin I_t] \mathbb{P}(i^* \notin I_t) \\ &\leq ((1 + p_2) \zeta_t^2 + (1 + 1/p_2) 4\gamma_t^2 M \ell_{f,d}^4) \mathbb{P}(i^* \in I_t) + ((1 + c_3\gamma_t) \zeta_t^2 + 3\ell_{g,1}\gamma_t\delta_t^2) \mathbb{P}(i^* \notin I_t) \\ &\leq \left(1 + c_3\gamma_t + \frac{3}{n}p_2\right) \zeta_t^2 + \left(1 + \frac{1}{p_2}\right) 4\gamma_t^2 M \ell_f^4 \frac{3}{n} + 4\ell_{g,1}\gamma_t\delta_t^2 + (2\ell_{g,1} + \ell_{g,1}^2) vd\gamma_t. \end{aligned} \quad (30)$$

Based on linearity of expectation and applying equation 30 recursively yields

$$\begin{aligned}
\mathbb{E} [\zeta_{t+1}^2] &\leq \sum_{t'=0}^t \left(\left(1 + \frac{1}{p_2}\right) 4\gamma^2 M \ell_f^4 \frac{3}{n} + 4\ell_{g,1} \gamma \mathbb{E} [\delta_{t'}^2] + (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma_t \right) \left(\prod_{k=t'+1}^t \left(1 + c_3 \gamma + \frac{3}{n} p_2\right) \right) \\
&= \sum_{t'=0}^t \left(\left(1 + \frac{1}{p_2}\right) 4\gamma^2 M \ell_f^4 \frac{3}{n} + 4\ell_{g,1} \gamma \mathbb{E} [\delta_{t'}^2] + (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma_t \right) \left(1 + c_3 \gamma + \frac{3}{n} p_2\right)^{t-t'} \\
&\stackrel{(a)}{\leq} \sum_{t'=0}^t \left(\left(1 + \frac{8T}{n}\right) 4\gamma^2 M \ell_f^4 \frac{3}{n} + 4\ell_{g,1} \gamma \mathbb{E} [\delta_{t'}^2] + (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma_t \right) \left(1 + \frac{1}{2T}\right)^{t-t'} \\
&\stackrel{(b)}{\leq} \sum_{t'=0}^t \left(\left(1 + \frac{8T}{n}\right) 4\gamma^2 M \ell_f^4 \frac{3}{n} + 4\ell_{g,1} \gamma \mathbb{E} [\delta_{t'}^2] + (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma_t \right) e^{\frac{1}{2}} \\
&\stackrel{(c)}{\leq} 2\gamma \sum_{t'=0}^t \left(\left(1 + \frac{8T}{n}\right) 4\gamma M \ell_f^4 \frac{3}{n} + 4\ell_{g,1} \mathbb{E} [\delta_{t'}^2] + (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma_t \right)
\end{aligned} \tag{31}$$

where (a) follows from choosing $\gamma_t \leq \gamma \leq 1/(8c_3T)$, $p_2 = n/(8T)$, (b) follows from $t - t' \leq T$, and $(1 + \frac{a}{T})^T \leq e^a$, and the inequality (c) follows from $e^{\frac{1}{2}} < 2$. Note that $\delta_0 = 0$, $\zeta_1 = 0$. Applying equation 27 at $t = 0$ gives

$$\mathbb{E} [\delta_1^2] \leq \frac{3}{n} \left(1 + \frac{1}{p}\right) 4\alpha^2 \ell_{f,d}^2 + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^2 + \alpha_t c_2 (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma_t$$

which together with equation 30 gives

$$\mathbb{E} [\zeta_2^2] \leq 4\ell_{g,1} \gamma_1 \delta_1^2 + \left(1 + \frac{1}{p_2}\right) 4\gamma_1^2 M \ell_{f,d}^4 \frac{3}{n} + (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma_t$$

Therefore, for $0 \leq t \leq 1$, it satisfies that

$$\begin{aligned}
\mathbb{E} [\delta_t^2] &\leq \left(\frac{3}{n} \left(1 + \frac{1}{p}\right) 4\alpha^2 \ell_{f,d}^2 + 24M \ell_{f,d}^4 c_2 \left(\frac{8\gamma T}{n} + \gamma \right) \frac{\alpha}{n} + (2\alpha c_2 + 1) (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^2 \right) \\
&\quad \underbrace{\left(\sum_{t'=0}^{t-1} \left(1 - \frac{1}{2} \alpha \mu + \frac{3p}{n}\right)^{t-t'-1} \right)}_{\beta_t} \\
&= \left(\frac{3}{n} \left(1 + \frac{1}{p}\right) 4\alpha^2 \ell_{f,d}^2 + 24M \ell_{f,d}^4 c_2 \left(\frac{8\gamma T}{n} + \gamma \right) \frac{\alpha}{n} + (2\alpha c_2 + 1) (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^2 \right)
\end{aligned} \tag{32}$$

Next, we will prove by induction that equation 32 also holds for $t > 1$. Assuming that equation 32 holds for all $0 \leq t \leq k \leq T - 1$, we apply equation 29 to the case where $t = k$ to obtain

$$\begin{aligned}
\mathbb{E} [\delta_{k+1}^2] &\leq \left(\eta_k + \frac{3p}{n} \right) \mathbb{E} [\delta_k^2] + \alpha_k c_2 (1 + c_3 \gamma_k) \mathbb{E} [\zeta_k^2] \mathbb{P}(i^* \notin I_t) + \frac{3}{n} \left(1 + \frac{1}{p} \right) 4\alpha_k^2 \ell_{f,d}^2 \\
&\quad + (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma_t \mathbb{P}(i^* \notin I_t) + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^2 \\
&\stackrel{(a)}{\leq} \left(\eta_k + \frac{3p}{n} \right) \mathbb{E} [\delta_k^2] \\
&\quad + 2\alpha_k c_2 \gamma \left(\sum_{t'=0}^k \left(\left(1 + \frac{8T}{n} \right) 4\gamma M \ell_f^4 \frac{3}{n} + 4\ell_{g,1} \mathbb{E} [\delta_{t'}^2] + (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma_t \right) \right) \mathbb{P}(i^* \notin I_t) \\
&\quad + \frac{3}{n} \left(1 + \frac{1}{p} \right) 4\alpha_k^2 \ell_f^2 + (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma_t \mathbb{P}(i^* \notin I_t) + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^2 \\
&\stackrel{(b)}{\leq} \underbrace{\left(\left(\eta_k + \frac{3p}{n} \right) \beta_k + 1 + 8\alpha_k c_2 (1 + c_3 \gamma_t) \ell_{g,1} \gamma \left(\sum_{t'=1}^k \beta_{t'} \right) \mathbb{P}(i^* \notin I_t) \right)}_{J_1} \\
&\quad \times \left(\frac{3}{n} \left(1 + \frac{1}{p} \right) 4\alpha^2 \ell_{f,d}^2 + 24M \ell_{f,d}^4 c_2 \left(\frac{8\gamma T}{n} + \gamma \right) \frac{\alpha}{n} + (2\alpha c_2 + 1) (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^2 \right).
\end{aligned} \tag{33}$$

where (a) follows from equation 31, and (b) follows from equation 32 for $0 \leq t \leq k$ and that $\gamma k \leq \gamma T \leq 1$. The coefficient J_1 in equation 33 can be further bounded by

$$\begin{aligned}
J_1 &= \left(\eta_k + \frac{3p}{n} \right) \beta_k + 1 + 8\alpha_k c_2 (1 + c_3 \gamma) \ell_{g,1} \gamma \left(\sum_{t'=1}^k c_{t'} \right) \mathbb{P}(i^* \notin I_t) \\
&\stackrel{(c)}{\leq} \left(\eta_k + \frac{3p}{n} \right) \beta_k + 1 + 16\alpha_k c_2 \ell_{g,1} k \gamma \beta_k \mathbb{P}(i^* \notin I_t) \\
&\stackrel{(d)}{\leq} \left(1 - \alpha_k (\mu - \ell_{F,d} c_1 - \ell_{f,1} v d - 4c_2 \ell_{g,1} \gamma - 16c_2 \ell_{g,1} k \gamma) \mathbb{P}(i^* \notin I_t) + \frac{3p}{n} \right) \beta_k + 1 \\
&\stackrel{(e)}{\leq} \left(1 - \frac{1}{2} \alpha \mu + \frac{3p}{n} \right) \beta_k + 1
\end{aligned} \tag{34}$$

where (c) follows from $\beta_t \leq \beta_{t+1}, \gamma_t \leq \gamma$ for all $t = 0, \dots, T$; (d) follows from the definition of η_k ; (e) is because $\gamma \leq \mu^2 / (120 \ell_{F,d}^2 T)$, $\alpha \leq 1 / (2\ell_{f,1}) \leq 1 / (2\mu)$ and choosing $c_1 = \mu / (4\ell_{F,d})$ leads to

$$\begin{aligned}
\ell_{F,d} c_1 + \ell_{f,1} v d + 4c_2 \ell_{g,1} \gamma + 16c_2 \ell_{g,1} k \gamma &\leq \ell_{F,d} c_1 + 21 (\ell_{F,d} c_1^{-1} + 2\alpha \ell_{F,d}^2) \ell_{g,1} (k+1) \gamma \\
&\leq \frac{1}{4} \mu + 21 (4\mu^{-1} + 2\alpha) \ell_{F,d}^2 \ell_{g,1} \frac{k+1}{T} \frac{\mu^2}{484 \ell_{F,d}^2 \ell_{g,1}} \leq \frac{1}{2} \mu.
\end{aligned}$$

Combining equation 33 and equation 34 implies

$$\begin{aligned}
\mathbb{E} [\delta_{k+1}^2] &\leq \left(\left(1 - \frac{1}{2} \alpha \mu + \frac{3p}{n} \right) \beta_k + 1 \right) \\
&\quad \left(\frac{3}{n} \left(1 + \frac{1}{p} \right) 4\alpha^2 \ell_{f,d}^2 + 24M \ell_{f,d}^4 c_2 \left(\frac{8\gamma T}{n} + \gamma \right) \frac{\alpha}{n} + (2\alpha c_2 + 1) (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^2 \right) \\
&= \left(\frac{3}{n} \left(1 + \frac{1}{p} \right) 4\alpha^2 \ell_{f,d}^2 + 24M \ell_{f,d}^4 c_2 \left(\frac{8\gamma T}{n} + \gamma \right) \frac{\alpha}{n} + (2\alpha c_2 + 1) (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^2 \right) \\
&\quad c_{k+1}
\end{aligned} \tag{35}$$

where the equality follows by the definition of β_t given in equation 32. The above statements from equation 33-equation 35 show that if equation 32 holds for all t such that $0 \leq t \leq k \leq T-1$, it also holds for $t = k+1$.

Therefore, we can conclude that for $T \geq 0$, it follows

$$\begin{aligned}
& \mathbb{E} [\delta_T^2] \\
& \leq \left(\frac{3}{n} \left(1 + \frac{1}{p} \right) 4\alpha^2 \ell_{f,d}^2 + 24M\ell_{f,d}^4 c_2 \left(\frac{8\gamma T}{n} + \gamma \right) \frac{\alpha}{n} + (2\alpha c_2 + 1) (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^2 \right) \\
& \beta_T \\
& = \left(\frac{3}{n} \left(1 + \frac{1}{p} \right) 4\alpha^2 \ell_{f,d}^2 + 24M\ell_{f,d}^4 c_2 \left(\frac{8\gamma T}{n} + \gamma \right) \frac{\alpha}{n} + (2\alpha c_2 + 1) (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^2 \right) \\
& \left(\sum_{k=0}^{T-1} \left(1 - \frac{1}{2} \alpha \mu + \frac{3p}{n} \right)^{T-k-1} \right) \\
& = \left(\frac{3}{n} \left(1 + \frac{1}{p} \right) 4\alpha^2 \ell_{f,d}^2 + 24M\ell_{f,d}^4 c_2 \left(\frac{8\gamma T}{n} + \gamma \right) \frac{\alpha}{n} + (2\alpha c_2 + 1) (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^2 \right) \\
& \left(\frac{1}{4} \alpha \mu \right)^{-1} \left(1 - \left(1 - \frac{1}{4} \alpha \mu \right)^T \right),
\end{aligned}$$

where the last equality follows from taking $p = \alpha \mu n / 12$, and compute the sum of geometric series.

By plugging in $c_1 = \mu / (4\ell_{F,d})$, $c_2 = \ell_{F,d} c_1^{-1} + 2\alpha \ell_{F,d}^2$, $c_3 = 3\ell_{F,d}^2 + 2\ell_{g,1}$, we have that

$$\begin{aligned}
\mathbb{E} [\delta_T^2] & \leq \left(\frac{3}{n} \left(1 + \frac{1}{p} \right) 4\alpha^2 \ell_{f,d}^2 + 24M\ell_{f,d}^4 c_2 \left(\frac{8\gamma T}{n} + \gamma \right) \frac{\alpha}{n} + (2\alpha c_2 + 1) (2\ell_{g,1} + \ell_{g,1}^2) v d \gamma \right. \\
& \quad \left. + \alpha_t \ell_{f,1} v d + 2\alpha_t^2 \ell_{f,1}^2 v^2 d^2 \right) \left(\frac{1}{4} \alpha \mu \right)^{-1} \\
& \leq \frac{48}{\mu n} \ell_{f,d}^2 \left(\alpha + \frac{12}{\mu n} + \frac{2M\ell_{f,d}^2 c_2 c_3^{-1}}{n} + 2M\ell_{f,d}^2 c_2 \gamma \right) + \frac{4}{\mu n} \left(\frac{10\alpha M\ell_{f,d}^2 \gamma + \mu \gamma}{\mu \alpha} + \alpha \ell_{f,1} + \frac{2\alpha \ell_{f,1}^2}{n} \right) \\
& \leq \frac{48}{\mu n} \ell_{f,d}^2 \left(\alpha + \frac{12 + 4M\ell_{f,d}^2}{\mu n} + \frac{10M\ell_{f,d}^4 \gamma}{\mu} \right) + \frac{4}{\mu n} \left(\frac{10\alpha M\ell_{f,d}^2 \gamma + \mu \gamma}{\mu \alpha} + \alpha \ell_{f,1} + \frac{2\alpha \ell_{f,1}^2}{n} \right)
\end{aligned}$$

where the second inequality follows from $v \leq \min \left\{ \frac{1}{nd}, \frac{1}{nd(2\ell_{g,1} + \ell_{g,1}^2)} \right\}$, and the last inequality

follows from $c_2 = \ell_{F,d}^2 (4\mu^{-1} + 2\alpha) \leq 5M\ell_{f,d}^2 \mu^{-1}$, and $c_2 c_3^{-1} \leq 5\ell_{F,d}^2 \mu^{-1} / (3\ell_{F,d}^2) \leq 2\mu^{-1}$.

D.5 PROOF OF THEOREM 4

Proof of Theorem 2. Combining the argument stability in Theorem 7, and Assumption 2, the MOL uniform stability can be bounded by

$$\begin{aligned}
& \sup_z \mathbb{E}_A \left[\|\nabla F_z(A(S)) - \nabla F_z(A(S'))\|_F^2 \right] \\
& \leq \mathbb{E}_A \left[\ell_{F,1}^2 \|A(S) - A(S')\|^2 \right] \\
& \leq \frac{48}{\mu n} \ell_{f,d}^2 \ell_{F,1}^2 \left(\alpha + \frac{12 + 4M\ell_{f,d}^2}{\mu n} + \frac{10M\ell_{f,d}^4 \gamma}{\mu} \right) + \frac{4}{\mu n} \ell_{F,1}^2 \left(\frac{10\alpha M\ell_{f,d}^2 \gamma + \mu \gamma}{\mu \alpha} + \alpha \ell_{f,1} + \frac{2\alpha \ell_{f,1}^2}{n} \right)
\end{aligned}$$

by Assumption 2

Finally, based on Propositions 2-3, we have

$$\begin{aligned}\mathbb{E}_{A,S} [R_{\text{gen}}(A(S))] &\leq \mathbb{E}_{A,S} [\|\nabla F(A(S)) - \nabla F_S(A(S))\|_{\text{F}}] \\ &\leq 4\epsilon_{\text{F}} + \sqrt{n^{-1} \mathbb{E}_S [\mathbb{V}_{z \sim \mathcal{D}} (\nabla F_z(A(S)))]} \\ &= \mathcal{O}\left(n^{-\frac{1}{2}}\right).\end{aligned}$$

The proof is completed.