Learning Robust to Distributional Uncertainties and Adversarial Data

Alireza Sadeghi¹⁰, *Member, IEEE*, Gang Wang¹⁰, *Senior Member, IEEE*, and Georgios B. Giannakis¹⁰, *Fellow, IEEE*

Abstract-Successful training of data-intensive deep neural networks critically rely on vast, clean, and high-quality datasets. In practice however, their reliability diminishes, particularly with noisy, outlier-corrupted data samples encountered in testing. This challenge intensifies when dealing with anonymized, heterogeneous data sets stored across geographically distinct locations due to, e.g., privacy concerns. This present paper introduces robust learning frameworks tailored for centralized and federated learning scenarios. Our goal is to fortify model resilience with a focus that lies in (i) addressing distribution shifts from training to inference time; and, (ii) ensuring testtime robustness, when a trained model may encounter outliers or adversarially contaminated test data samples. To this aim, we start with a centralized setting where the true data distribution is considered unknown, but residing within a Wasserstein ball centered at the empirical distribution. We obtain robust models by minimizing the worst-case expected loss within this ball, yielding an intractable infinite-dimensional optimization problem. Upon leverage the strong duality condition, we arrive at a tractable surrogate learning problem. We develop two stochastic primal-dual algorithms to solve the resultant problem: one for ϵ -accurate convex sub-problems and another for a single gradient ascent step. We further develop a distributionally robust federated learning framework to learn robust model using heterogeneous data sets stored at distinct locations by solving per-learner's sub-problems locally, offering robustness with modest computational overhead and considering data distribution. Numerical tests corroborate merits of our training algorithms against distributional uncertainties and adversarially corrupted test data samples.

Index Terms—Wasserstein distance, distributionally robust optimization, minimax, primal-dual, federated learning.

I. INTRODUCTION

THE RELIABLE performance of contemporary machine learning models largely hinges on their training algorithms. In practice, their generalization is often compromised during inference, especially when confronted with noisy data tainted by outliers, as cited in previous research [1]. These issues primarily stem from imperfect data acquisition and distribution shifts from training to inference phases. Furthermore,

Manuscript received 25 October 2023; revised 3 February 2024; accepted 19 March 2024. Date of publication 26 March 2024; date of current version 16 April 2024. The work of Alireza Sadeghi and Georgios B. Giannakis was supported in part by NSF under Grant 1901134, Grant 2128593, Grant 2126052, Grant 2212318, Grant 2220292, Grant 2312547, Grant 2103256, and Grant 2212318. (Corresponding author: Alireza Sadeghi.)

The authors are with the Department of ECE, University of Minnesota, Minneapolis, MN 55455 USA (e-mail: sadeghi@umn.edu; gangwang@bit.edu.cn; georgios@umn.edu).

Digital Object Identifier 10.1109/JSAIT.2024.3381869

massive datasets are typically not confined to a single device, instead they are generated and distributed across multiple locations (referred to as workers) with limited computing and storage capabilities. Distribution shifts among workers, along with privacy and confidentiality concerns, pose additional challenges in training reliable models from distributed datasets [2], [3].

Recent research efforts have been devoted to providing reliable models with enhanced generalization guarantees, particularly by addressing the susceptibility of models to adversarial data samples. These strategies fall into two categories: detection approaches, which determine if the model input is anomalous [4], [5], and robust training methods that endow the model with resilience against attacks [6], [7].

While detection-based approaches have shown some success in ensuring trustworthy inference [8], [9], [10], adversarial training remains a promising approach for robust model training [11], [12]. By carefully maximizing the model's loss, adversarial training methods introduce imperceptible noise to clean input training data to enhance robustness [13]. This approach typically involves solving an optimization problem and relies on the trained model or a pre-trained one [14]. However, solving the resultant optimization problem is often challenging, especially when tuning hyperparameters, which can be cumbersome in certain cases [15], [16]. Moreover, these methods often focus on per-datum processing to generate adversarial noise, simply disregarding the underlying datagenerating process or distribution. Additionally, quantifying the robustness gained through such training methods remains a challenge [17], and is less explored in research.

Additional challenges arise from the prevalent assumption in many methods that data is confined to a single device. In reality, data is often generated and stored across distributed locations, each having subsets of data with potentially different distributions. While data localization is essential for privacy among other reasons, the federated learning (FL) paradigm aims to train a global model by coordinating multiple devices through a central parameter server [2], [18]. Existing FL approaches have primarily focused on the trade-offs between communication and computation when aggregating model updates from learners [19], [20], [21], with limited attention to robust FL methods. Among those addressing robust FL, some focus on robust aggregation strategies [22], [23], or a blend of periodic averaging with adaptive sampling [24], [25], while others target robust models and assume affine distribution shifts across workers' data see, e.g., [26]. However, these methods have limited scopes and do not fully account for underlying data distributions when designing training algorithms. In a recent effort to deal with discrepancies in data distributions across workers, [27] has extended standard results of optimal transport theory to the multi-marginal optimal transport problem and studied personalized federated learning. Albeit interesting, the proposed method requires an additional model per worker and does not target robustness at all. From a theoretical perspective, the optimal transport theory has further been investigated to quantify adversarial risk of classifiers on adversarially perturbed data [28], [29].

Our work leverages the distributionally robust optimization framework to develop procedures for training robust models against distributional uncertainties and adversarial data. It builds on the implicit assumption that training data samples are clean or nominal, but obtained from an "unknown" data distribution. The objective is to target distribution shifts of data from training to inference, and ensure robustness to adversarial inputs during the inference phase. Our algorithms enable learning robust models from multiple heterogeneous datasets with potentially different distributions stored at different locations. Unlike most existing works, our proposed algorithms account for the data distribution when creating adversarial perturbations. We obtain robust models by minimizing the worst-case expected loss over all distributions residing within a Wasserstein ball centered at the empirical data distribution. The resulting formulation leads to a challenging and intractable infinite-dimensional functional optimization problem, which we address by exploiting the strong duality to arrive at a tractable, and under certain conditions, an equivalent unconstrained and affordable minimization problem. To solve this robust surrogate problem, we develop a stochastic proximal gradient descent (SPGD) algorithm based on an ϵ -accurate oracle, along with its lightweight stochastic proximal gradient descent-ascent (SPGDA) iterations. The former algorithm uses an oracle to solve the involved convex sub-problems to ϵ accuracy, while the latter approximates the solution through a single gradient ascent step. To further accommodate learning from distributed datasets, we develop a distributionally robust federated learning (DRFL) framework. In summary, the main contributions of this paper are as follows.

- A generic distributionally robust learning framework, designed to endow machine learning models with resilience against distributional uncertainties and adversarial input data encountered during inference.
- Two scalable distributionally robust learning algorithms, with rigorous theoretical convergence guarantees.
- A DRFL framework, enabling the training of robust models from multiple sources with potentially different underlying data distributions.

The rest of this paper is structured as follows. Problem formulation and its robust surrogate are the subjects of Section II. The proposed SPGD with ϵ -accurate oracle and SPGDA algorithms with their convergence analyses are presented in Sections III and IV, respectively. The DRFL implementation is discussed in Section V. Numerical tests are given in Section VI with conclusions drawn in Section VII. Technical proofs are deferred to the Appendix.

II. PROBLEM STATEMENT

Consider the standard regularized statistical learning task

$$\min_{\theta \in \Theta} \ \mathbb{E}_{z \sim P_0} [\ell(\theta; z)] + r(\theta) \tag{1}$$

where $\ell(\theta;z)$ denotes the loss of a model parameterized by the unknown parameter vector θ on a datum $z=(x,y)\sim P_0$, with feature x and label y, drawn from some nominal distribution P_0 . Here, Θ denotes the feasible set for model parameters. To prevent over fitting or incorporate prior information, regularization term $r(\theta)$ is oftentimes added to the expected loss. Popular regularizers include $r(\theta) \coloneqq \beta \|\theta\|_1^2$ or $\beta \|\theta\|_2^2$, where $\beta \ge 0$ is a hyper-parameter controlling the importance of the regularization term relative to the expected loss.

In practice, the nominal distribution P_0 is typically unknown. Instead, we are given with independently and identically distributed training data samples $\{z_n\}_{n=1}^N \sim P_0$. Leveraging these samples, we form the empirical data distribution $\widehat{P}_0^{(N)}$ to replace P_0 in (1) in *empirical* loss minimization

$$\min_{\boldsymbol{\theta} \in \Theta} \ \bar{\mathbb{E}}_{z \sim \widehat{P}_0^{(N)}} [\ell(\boldsymbol{\theta}; z)] + r(\boldsymbol{\theta}) \tag{2}$$

where $\overline{\mathbb{E}}_{z\sim\widehat{P}_0^{(N)}}[\ell(\theta;z)]=N^{-1}\sum_{n=1}^N\ell(\theta;z_n)$. Indeed, a variety of machine learning tasks can be cast as (2), including, e.g., ridge and Lasso regression, logistic regression, and reinforcement learning. The resultant models obtained by solving (2) however, have been shown vulnerable to adversarially corrupted data in $\widehat{P}_0^{(N)}$. Further, the testing data distribution often deviates from the available $\widehat{P}_0^{(N)}$. For this reason, targeting an adversarially robust model against a set of distributions corresponding to perturbations of the underlying data distribution, has led to the formulation

$$\min_{\theta \in \Theta} \sup_{P \in \mathcal{P}} \mathbb{E}_{z \sim P}[\ell(\theta; z)] + r(\theta)$$
 (3)

where \mathcal{P} represents a set of distributions centered around the data generating distribution $\widehat{\mathcal{P}}_0^{(N)}$. Compared with (1), the worst-case formulation (3), yields models ensuring reasonable performance across a continuum of distributions characterized by \mathcal{P} . In practice, different types of ambiguity sets \mathcal{P} can be considered, and they lead to different robustness guarantees and computational requirements. Popular choices of \mathcal{P} include momentum [30], [31], KL divergence [32], statistical test [33], and Wasserstein distance-based ambiguity sets [33], [34]; see, e.g., [35] for a recent overview. Among all choices, it has been shown that the Wasserstein ambiguity set \mathcal{P} results in a tractable realization of (3), thanks to the strong duality result of [33] and [34], which also motivates this work.

To formalize this, consider two probability measures P and Q supported on set Z, and let $\Pi(P,Q)$ be the set of all joint measures supported on Z^2 , with marginals P and Q. Let $c: Z \times Z \to [0, \infty)$ measure the cost of transporting a unit of mass from z in P to another element z' in Q. The celebrated optimal transport problem is given by [36, p. 111]

$$W_c(P,Q) := \inf_{\pi \in \Pi} \mathbb{E}_{\pi} [c(z,z')]. \tag{4}$$

Remark 1: If $c(\cdot, \cdot)$ satisfies the axioms of distance, then W_c defines a distance on the space of probability measures. For

instance, if P and Q are defined over a Polish space equipped with metric d, then choosing $c(z,z')=d^p(z,z')$ for some $p\in [1,\infty)$ asserts that $W_c^{1/p}(P,Q)$ is the well-known Wasserstein distance of order p between probability measures P and Q [36, Definition 6.1].

For a given empirical distribution $\widehat{P}_0^{(N)}$, define the uncertainty set $\mathcal{P} := \{P|W_c(P,\widehat{P}_0^{(N)}) \leq \rho\}$ to include all probability distributions having at most ρ -distance from $P_0^{(N)}$. Incorporating this ambiguity set into (3), yields the following reformulation

$$\min_{\theta \in \Theta} \sup_{P} \mathbb{E}_{z \sim P}[\ell(\theta; z)] + r(\theta)$$
 (5a)

s.t.
$$W_c(P, \widehat{P}_0^{(N)}) \le \rho$$
. (5b)

Observe that the inner supremum in (5a) runs over all joint probability measures π on \mathbb{Z}^2 implicitly characterized by (5b). Intuitively, directly solving this optimization over the infinite-dimensional space of distribution functions is challenging, if not impossible. Fortunately, for a broad range of losses as well as transport costs, it has been shown that the inner maximization satisfies a strong duality condition [35]; that is, the optimal objective of this inner maximization and its Lagrangian dual optimal objective, are equal. In addition, the dual problem involves optimization over a one-dimensional dual variable. These two observations make it possible to solve (3) in the dual domain. To formally obtain a tractable surrogate to (5), we make the following assumptions.

Assumption 1: The transportation cost function $c: \mathbb{Z} \times \mathbb{Z} \to [0, \infty)$, is a lower semi-continuous function that to satisfy well-posedness of the transportation problem and the existence of optimal transport maps it satisfies the following conditions:

- 1) Non-Negativity of transportation $c(z, z') \ge 0 \quad \forall z, z' \in \mathcal{Z} \times \mathcal{Z}'$, and c(z, z) = 0, $\forall z \in \mathcal{Z} \times \mathcal{Z}'$.
- Boundedness to ensuring there is a lower bound on the transportation cost, i.e.,

$$0 < \inf_{z,z'} c(z,z') \le \sup_{z,z'} c(z,z') < \infty, \quad \forall z,z',z \ne z'.$$

- Measurability of transportation cost in the product measure P × Q on the product space Z × Z', ensuring involved integrals in (4) are well-defined.
- Strong convexity: Transportation cost is strongly convex, i.e., having a positive definite Hessian

$$H^{1} := \begin{bmatrix} \frac{\partial^{2} c}{\partial z^{2}} & \frac{\partial^{2} c}{\partial z \partial z^{\prime}} \\ \frac{\partial^{2} c}{\partial z^{\prime} \partial z} & \frac{\partial^{2} c}{\partial z^{\prime}^{2}} \end{bmatrix} > \mathbf{0}.$$
 (6)

Remark 2: Unless otherwise noted, throughout this paper, a generic transportation cost $c(z,z'): \mathcal{Z} \times \mathcal{Z}' \to \mathbb{R}^+$ satisfying Assumption 1 is considered. Therefore, the Wasserstein distance (4) is for such generic cost accordingly. No additional details are considered for this distance except having power

of the Wasserstein distance to be equal to 1 for the sake of simplicity of derivatives.

Assumption 2: The loss function $\ell: \Theta \times \mathcal{Z} \to [0, \infty)$, is upper semi-continuous, and integrable.

The following proposition provides a tractable surrogate for (5), whose proof can be found in [35, Th. 1].

Proposition 1: Let $\ell: \Theta \times \mathcal{Z} \to [0, \infty)$, and $c: \mathcal{Z} \times \mathcal{Z} \to [0, \infty)$ satisfy Assumptions 1 and 2, respectively. Then, for any given $\widehat{P}_0^{(N)}$, and $\rho > 0$, it holds that

$$\sup_{\boldsymbol{e} \in \mathcal{P}} \mathbb{E}_{\boldsymbol{z} \sim P}[\ell(\boldsymbol{\theta}; \boldsymbol{z})]$$

$$= \inf_{\boldsymbol{\gamma} \geq 0} \left\{ \mathbb{E}_{\boldsymbol{z} \sim \widehat{P}_0^{(N)}} \left[\sup_{\boldsymbol{\zeta} \in \mathcal{Z}} \{\ell(\boldsymbol{\theta}; \boldsymbol{\zeta}) - \gamma(\boldsymbol{c}(\boldsymbol{z}, \boldsymbol{\zeta}) - \rho)\} \right] \right\}$$
(7)

where $\mathcal{P} := \{P|W_c(P, \widehat{P}_0^{(N)}) \leq \rho\}.$

Remark 3: Thanks to strong duality, the right-hand side in (7) simply is a uni-variate dual reformulation of the primal problem represented in the left-hand side. In sharp contrast with the primal formulation, the expectation in the dual domain is taken only over the empirical $\widehat{P}_0^{(N)}$ rather than any $P \in \mathcal{P}$. In addition, since this reformulation circumvents the need for finding the optimal $\pi \in \Pi$ to form \mathcal{P} , and characterizing the primal objective $\forall P \in \mathcal{P}$, it is practically more convenient.

Upon relying on Proposition 1, the following distributionally robust surrogate is obtained

$$\min_{\theta \in \Theta} \inf_{\gamma \ge 0} \left\{ \bar{\mathbb{E}}_{z \sim \widehat{P}_0^{(N)}} \left[\sup_{\zeta \in \mathcal{Z}} \{ \ell(\theta; \zeta) + \gamma(\rho - c(z, \zeta)) \} + r(\theta) \right] \right\}.$$
(8)

Remark 4: The robust surrogate in (3) is reminiscent of minimax (saddle-point) optimization problems. Solvers for such problems have been recently investigated in several contributions see, e.g., [37] and references therein. While such solvers address the standard minimax problems, our formulated problem under (3), has its own specific challenges, posed by the structure of (8) to promote robustness, as well as taking into the account the data distribution when training the model. It is important to note that [37] and references therein offer valuable insights for conducting convergence analysis of primal-dual type solvers, as has been leveraged in this work.

It is worth noting that a relaxed (hence suboptimal) version of (8) with a fixed γ value has recently been studied in [34]. Unfortunately, one has to select an appropriate γ value using cross validation over a grid search that is also application dependent. Heuristically choosing a γ does not guarantee optimality in solving the distributionally robust surrogate (8). Clearly, the effect of heuristically selecting γ is more pronounced when training deep neural networks. Instead, we advocate algorithms that optimize γ and θ simultaneously.

Our approach to addressing this, relies on the structure of (8) to *iteratively* update parameters $\bar{\theta} := [\theta^\top \ \gamma]^\top$ and ζ . To end up with a differentiable function of $\bar{\theta}$ after maximizing over ζ , Danskin's theorem requires the sup-problem to have a unique solution [38]. For this reason, we design the inner maximization to involve a strongly concave objective function through the selection of a strongly convex transportation cost, such as $c(z, z') := \|z - z'\|_p^2$ for $p \ge 1$. For the maximization over ζ to rely on a strongly concave objective, we let

¹Although not necessary in general, $c(\cdot, \cdot)$ is considered to be strongly convex. Note that all relevant derivatives discussed in subsequent sections are derived under the assumption that $c(\cdot, \cdot)$ represents a *generic transportation* cost, adhering to the specified Assumptions 1.

 $\gamma \in \Gamma := \{\gamma | \gamma > \gamma_0\}$, where γ_0 is large enough. Since γ is the dual variable corresponding to the constraint in (5), having $\gamma \in \Gamma$ is tantamount to tuning ρ , which in turn *controls* the level of *robustness*. Replacing $\gamma \geq 0$ in (8) with $\gamma \in \Gamma$, our *robust learning model* is obtained as the solution of

$$\min_{\boldsymbol{\theta} \in \Theta} \inf_{\boldsymbol{\gamma} \in \Gamma} \bar{\mathbb{E}}_{\boldsymbol{z} \sim \widehat{P}_0^{(T)}} \left[\sup_{\boldsymbol{\xi} \in \mathcal{Z}} \psi \left(\bar{\boldsymbol{\theta}}, \boldsymbol{\xi}; \boldsymbol{z} \right) \right] + r(\bar{\boldsymbol{\theta}}) \tag{9}$$

where

$$\psi(\bar{\theta}, \zeta; z) := \ell(\theta; \zeta) + \gamma(\rho - c(z, \zeta)).$$
 (10)

Intuitively, input z in (9) is pre-processed by maximizing ψ accounting for the adversarial perturbation. To iteratively solve our objective in (9), the ensuing sections provide efficient solvers under some mild conditions. Those include cases, every inner maximization (supremum) can be solved to ϵ -optimality by an oracle.

Before developing our algorithms, we start by making several standard assumptions.

Assumption 3: Function $c(z,\cdot)$ is L_c -Lipschitz and μ -strongly convex for any given $z \in \mathcal{Z}$, with respect to the norm $\|\cdot\|$.

Assumption 4: The loss function $\ell(\theta; z)$ satisfies the following Lipschitz smoothness conditions

$$\|\nabla_{\theta}\ell(\theta;z) - \nabla_{\theta}\ell(\theta';z)\|_{*} \le L_{\theta\theta}\|\theta - \theta'\|$$
 (11a)

$$\|\nabla_{\theta}\ell(\theta;z) - \nabla_{\theta}\ell(\theta;z')\|_{*} < L_{\theta\tau}\|z - z'\| \tag{11b}$$

$$\|\nabla_z \ell(\theta; z) - \nabla_z \ell(\theta; z')\|_* \le L_{zz} \|z - z'\| \tag{11c}$$

$$\|\nabla_{z}\ell(\theta;z) - \nabla_{z}\ell(\theta';z)\|_{*} \le L_{z\theta}\|\theta - \theta'\| \tag{11d}$$

and it is continuously differentiable with respect to θ .

Assumption (4) guarantees that the supremum in (8) results in a smooth function of $\bar{\theta}$; thus, one can execute gradient descent to update θ upon solving the supremum. This will further help to provide convergence analysis of our proposed algorithms. To elaborate more on this, the following lemma characterizes the smoothness and gradient Lipschitz properties obtained upon solving the maximization problem in (9).

Lemma 1: For each $z \in \mathcal{Z}$, let us define $\bar{\psi}(\bar{\theta};z) = \sup_{\zeta} \psi(\bar{\theta},\zeta;z)$ with $\zeta_*(\bar{\theta};z) := \arg\max_{\zeta \in \mathcal{Z}} \psi(\bar{\theta},\zeta;z)$. Then $\bar{\psi}(\cdot)$ is differentiable, and its gradient is $\nabla_{\bar{\theta}}\bar{\psi}(\bar{\theta};z) = \nabla_{\bar{\theta}}\psi(\bar{\theta},\zeta_*(\bar{\theta};z);z)$. Moreover, the following conditions hold

$$\left\| \zeta_* (\bar{\theta}_1; z) - \zeta_* (\bar{\theta}_2; z) \right\| \le \frac{L_{z\theta}}{\lambda} \|\theta_2 - \theta_1\| + \frac{L_c}{\lambda} \|\gamma_2 - \gamma_1\| \tag{12a}$$

and

$$\|\nabla_{\bar{\theta}}\bar{\psi}(\bar{\theta}_{1};z) - \nabla_{\bar{\theta}}\bar{\psi}(\bar{\theta}_{2};z)\| \leq \frac{L_{\theta z}L_{c} + L_{c}^{2}}{\lambda} \|\gamma_{2} - \gamma_{1}\| + \left(L_{\theta\theta} + \frac{L_{\theta z}L_{z\theta} + L_{c}L_{z\theta}}{\lambda}\right) \|\theta_{2} - \theta_{1}\|.$$
 (12b)

where $\gamma^{1,2} \in \Gamma$, and $\psi(\bar{\theta}, \cdot; z)$ is λ -strongly concave.

Proof: See Appendix-A for the proof.

Lemma 1 paves the way for iteratively solving the surrogate optimization (9), intuitively because it guarantees a differentiable and smooth objective upon solving the inner supremum to its optimum.

Remark 5: Equation (12a) is appealing in practice. Indeed, if $\bar{\theta}^t = [\theta^t, \gamma^t]$ is updated with a small enough step size, the corresponding $\zeta_*(\theta^{t+1};z)$ is close enough to $\zeta_*(\theta^t;z)$. Building on this observation, instead of using an oracle to find the optimum $\zeta_*(\theta^{t+1};z)$, an ϵ -accurate solution $\zeta_\epsilon(\theta^{t+1};z)$ suffices to obtain comparable performance. This also circumvents the need to find the optimum for the inner maximization per iteration, which could be computationally demanding.

III. STOCHASTIC PROXIMAL GRADIENT DESCENT WITH ϵ -ACCURATE ORACLE

A standard solver of regularized optimization problems is the proximal gradient algorithm. In this section, we develop a variant of it to tackle the robust surrogate (9). For convenience, let us define

$$f(\theta, \gamma) := \mathbb{E}\left[\sup_{\zeta \in \mathcal{Z}} \{\ell(\theta; \zeta) + \gamma(\rho - c(z, \zeta))\}\right]$$
(13)

and rewrite our objective as

$$\min_{\theta \in \Theta} \inf_{\gamma \in \Gamma} F(\theta, \gamma) := f(\theta, \gamma) + r(\theta)$$
 (14)

where $f(\theta, \gamma)$ is the smooth function in (13), and $r(\cdot)$ is a non-smooth and convex regularizer, such as the ℓ_1 -norm. With a slight abuse of notation, upon introducing $\bar{\theta} := [\theta \ \gamma]$, we define $f(\bar{\theta}) := f(\theta, \gamma)$ and $F(\bar{\theta}) := F(\theta, \gamma)$.

The proximal gradient algorithm the updates $\bar{\theta}^{t}$ as

$$\bar{\theta}^{t+1} = \arg\min_{\theta} \alpha_t r(\theta) + \alpha_t (\theta - \bar{\theta}^t, g(\bar{\theta}^t)) + \frac{1}{2} \|\theta - \bar{\theta}^t\|^2$$

where $g(\bar{\theta}^t) := \nabla f(\bar{\theta})|_{\bar{\theta} = \bar{\theta}^t}$, and $\alpha_t > 0$ is some step size. The last update is expressed in the compact form

$$\bar{\boldsymbol{\theta}}^{t+1} = \operatorname{prox}_{\alpha,r} \left[\bar{\boldsymbol{\theta}}^t - \alpha_t g \left(\bar{\boldsymbol{\theta}}^t \right) \right] \tag{15}$$

where the proximal gradient operator is given by

$$\operatorname{prox}_{\alpha r}[v] := \arg\min_{\theta} \alpha r(\theta) + \frac{1}{2} \|\theta - v\|^2. \tag{16}$$

The working assumption is that this optimization problem can be solved efficiently using off-the-shelf solvers.

Starting from the guess $\bar{\theta}^0$, the proposed SPGD with ϵ -accurate oracle executes two steps per iteration $t=1,2,\ldots$ First, it relies on an ϵ -accurate maximum oracle to solve the inner problem $\sup_{\zeta \in \mathcal{Z}} \{\ell(\theta^t; \zeta) - \gamma^t c(z, \zeta)\}$ for randomly drawn samples $\{z_n\}_{n=1}^N$ to yield ϵ -optimal $\xi_{\epsilon}(\bar{\theta}^t, z_n)$ with the corresponding objective values $\psi(\bar{\theta}^t, \zeta_{\epsilon}(\bar{\theta}^t, z_n); z_n)$. Next, $\bar{\theta}^t$ is updated using a stochastic proximal gradient step as

$$\bar{\theta}^{t+1} = \operatorname{prox}_{\alpha_{t}r} \left[\bar{\theta}^{t} - \frac{\alpha_{t}}{N} \sum_{n=1}^{N} \nabla_{\bar{\theta}} \psi \left(\bar{\theta}, \zeta_{\epsilon}(\bar{\theta}^{t}; z_{n}); z_{n} \right) \right].$$

For implementation, the proposed SPGD algorithm with ϵ -accurate oracle is summarized in Alg. 1. Convergence performance of this algorithm is analyzed in the ensuing subsection.

Algorithm 1: SPGD With ϵ -Accurate Oracle

Input: Initial guess $\bar{\theta}^0$, step size sequence $\{\alpha_t > 0\}_{t=0}^T$, ϵ -accurate oracle

1 for $t = 1, \ldots, T$ do

2 | Draw i.i.d samples $\{z_n\}_{n=1}^N$ 3 | Find ϵ -optimizer $\zeta_{\epsilon}(\bar{\theta}^t; z_n)$ via the oracle

4 | Update: $\bar{\theta}^{t+1} = \max_{prox_{\alpha_{t}r}} \left[\bar{\theta}^t - \frac{\alpha_t}{N} \sum_{n=1}^N \nabla_{\bar{\theta}} \psi(\bar{\theta}, \zeta_{\epsilon}(\bar{\theta}^t; z_n); z_n) \Big|_{\bar{\theta} = \bar{\theta}^t}\right]$ 5 end

A. Convergence of SPGD With ϵ -Accurate Oracle

In general, the postulated model is nonlinear, and the robust surrogate $F(\bar{\theta})$ is nonconvex. In this section, we characterize the convergence performance of Alg. 1 to a stationary point. However, lack of convexity and smoothness implies that stationary points must be understood in the sense of the Frèchet subgradient. Specifically, the Frèchet subgradient $\partial F(\check{\theta})$ for the composite optimization in (14), is the set [39]

$$\partial F\left(\check{\boldsymbol{\theta}}\right) \coloneqq \left\{ v \mid \lim_{\bar{\boldsymbol{\theta}} \to \check{\boldsymbol{\theta}}} \inf \frac{F(\bar{\boldsymbol{\theta}}) - F\left(\check{\boldsymbol{\theta}}\right) - v^{\top}\left(\bar{\boldsymbol{\theta}} - \check{\boldsymbol{\theta}}\right)\right)}{\|\bar{\boldsymbol{\theta}} - \check{\boldsymbol{\theta}}\|} \ge 0 \right\}.$$

Consequently, the distance between vector $\mathbf{0}$ and the set $\partial F(\check{\boldsymbol{\theta}})$ is a measure characterizing whether a point is stationary or not. To this end, define the distance between a vector \boldsymbol{v} and a set \mathcal{S} as $\operatorname{dist}(\boldsymbol{v},\mathcal{S}) := \min_{s \in \mathcal{S}} \|\boldsymbol{v} - s\|$, and the notion of δ -stationary points as defined next.

Definition 1: Given a small $\delta > 0$, we call vector $\check{\theta}$ a δ -stationary point if and only if $\operatorname{dist}(0, \partial F(\check{\theta})) \leq \delta$.

Since $f(\cdot)$ in (13) is smooth, we have that $\partial F(\bar{\theta}) = \nabla f(\bar{\theta}) + \partial r(\bar{\theta})$ [39]. Hence, it suffices to prove that the algorithm converges to a δ -stationary point $\check{\theta}$ satisfying

$$\operatorname{dist}(\mathbf{0}, \nabla f(\check{\boldsymbol{\theta}}) + \partial r(\check{\boldsymbol{\theta}})) \le \delta. \tag{17}$$

We further adopt the following assumption that is standard in stochastic optimization.

Assumption 5: Function f satisfies the next two conditions.

- 1) Gradient estimates are unbiased and have a bounded variance, i.e., $\mathbb{E}[g^*(\bar{\theta}^t) \nabla f(\bar{\theta}^t)] = 0$, and there is a constant $\sigma^2 < \infty$, so that $\mathbb{E}[\|\nabla f(\bar{\theta}^t) g^*(\bar{\theta}^t)\|_2^2] \le \sigma^2$.
- 2) Function $f(\bar{\theta})$ is smooth with L_f -Lipschitz continuous gradient, i.e., $\|\nabla f(\bar{\theta}_1) \nabla f(\bar{\theta}_2)\| \le L_f \|\bar{\theta}_1 \bar{\theta}_2\|$.

We are now ready to claim the convergence guarantees for Alg. 1; see Appendix-B for the proof.

Theorem 1: Let Alg. 1 run for T iterations with constant step sizes α , $\eta > 0$. Under Assumptions 1–5, Alg. 1 generates a sequence of $\{\bar{\theta}^t\}$ that satisfies

$$\mathbb{E}\left[\operatorname{dist}\left(\mathbf{0}, \partial F(\bar{\boldsymbol{\theta}}^{t'})\right)^{2}\right] \leq \left(\frac{2}{\alpha} + \beta\right) \frac{\Delta_{F}}{T} + \left(\frac{\beta}{\eta} + 2\right) \sigma^{2} + \frac{(\beta + 2)L_{\bar{\boldsymbol{\theta}}z}^{2} \epsilon}{\lambda_{0}}$$

$$(18)$$

where t' is uniformly sampled from $\{1, \ldots, T\}$; here, $\Delta_F := F(\bar{\theta}^0) - F(\bar{\theta}^{T+1})$; $L_{\bar{\theta}z}^2 := L_{\theta z}^2 + \lambda_0 L_c$, and β , $\lambda_0 > 0$ are some constants.

Theorem 1 asserts that $\{\bar{\theta}^t\}_{t=1}^T$ generated by Alg. 1 converges to a stationary point on average. The upper bound here is characterized by the initial error Δ_F , which decays at the rate of $\mathcal{O}(1/T)$; and, the constant bias terms induced by the gradient estimate variance σ^2 as well as the oracle accuracy ϵ .

Remark 6: . Note in Theorem 1 the established convergence rate of O(1/T) to a stationary point is for the squared error distance, which, upon taking the square root, yields a rate of $O(1/\sqrt{T})$, which relies on a *constant* stepsize. By carefully designing the step size, it may be possible to obtain even a tighter upper bound. In addition, it is important to note that this Theorem demonstrates the attainment of a δ -stationary point, where parameter δ depends on the accuracy of the oracle captured via ϵ , and the variance of the gradient estimates σ^2 . Please observe that σ^2 depends on the batch size M, as is evident from equation (22). While we have not explicitly stated the dependency of σ^2 on M in Assumption 5, yet it is evident from equation (22), that a larger batch size should yield to a decrease in gradient estimate variance. All in all, intuitively having a larger batch size will yield a smaller variance for gradient estimates, i.e., $\delta \propto \sigma^2 \propto \frac{1}{M}$, and it also holds that $\delta \propto \epsilon$. Thus, having an accurate oracle with small $\epsilon \to 0$ with a larger batch size $M \to \infty$ results in obtaining a stationary point (i.e., $\delta \to 0$). These details give an intuitive explanation why the convergence guarantee comes with a bias term.

The computational complexity of Alg. 1 can grow prohibitively when dealing with large-size datasets and complex models. This motivates lightweight, scalable, yet efficient methods. To this end, we introduce next a stochastic proximal gradient descent-ascent (SPGDA) algorithm.

IV. STOCHASTIC PROXIMAL GRADIENT DESCENT-ASCENT

Leveraging the strong concavity of the inner maximization problem and Lemma 1, a lightweight variant of the SPGD with ϵ -accurate oracle is developed here. Instead of optimizing the inner maximization problem to ϵ -accuracy by an oracle, we approximate its solution after only a single gradient ascent step. Specifically, for a batch of data $\{z_m^t\}_{m=1}^M$ per iteration t, our SPGDA algorithm first perturbs each datum via a gradient ascent step

$$\zeta_m^t = z_m^t + \eta_t \nabla_{\zeta} \psi \left(\bar{\theta}^t, \zeta; z_m^t \right) \Big|_{\zeta = z_m^t}, \ \forall m = 1, \dots, M \ (19)$$

and then forms

$$g^{t}(\bar{\theta}^{t}) := \frac{1}{M} \sum_{m=1}^{M} \nabla_{\bar{\theta}} \psi(\bar{\theta}, \zeta_{m}^{t}; z_{m}^{t}) \big|_{\bar{\theta} = \bar{\theta}^{t}}. \tag{20}$$

Using (20), an extra proximal gradient step is taken to obtain

$$\bar{\boldsymbol{\theta}}^{t+1} = \operatorname{prox}_{\alpha_t r} \left[\bar{\boldsymbol{\theta}}^t - \alpha_t g^t \left(\bar{\boldsymbol{\theta}}^t \right) \right]. \tag{21}$$

The SPGDA steps are summarized in Alg. 2. Besides its simplicity and scalability, SPGDA enjoys convergence to a stationary point as elaborated next.

Algorithm 2: SPGDA

Input: Initial guess
$$\bar{\theta}^0$$
, step size sequence $\{\alpha_t, \eta_t > 0\}_{t=0}^T$, batch size M

1 for $t = 1, \ldots, T$ do

2 | Draw a batch of i.i.d samples $\{z_m\}_{m=1}^M$

3 | Find $\{\zeta_m^t\}_{m=1}^M$ via gradient ascent: $\zeta_m^t = z_m^t + \eta_t \nabla_{\zeta} \psi(\bar{\theta}^t, \zeta; z_m^t) \big|_{\zeta = z_m^t}, \quad m = 1, \ldots, M$

Update: $\bar{\theta}^{t+1} = \operatorname{prox}_{\alpha_{tr}} \left[\bar{\theta}^t - \frac{\alpha_t}{M} \sum_{m=1}^M \nabla_{\bar{\theta}} \psi(\bar{\theta}^t, \zeta_m^t; z_m^t) \big|_{\bar{\theta} = \bar{\theta}^t}\right]$

4 end

A. Convergence of SPGDA

To prove convergence of Alg. 2, let us start by defining

$$g^*\left(\bar{\boldsymbol{\theta}}^t\right) := \frac{1}{M} \sum_{m=1}^{M} \nabla_{\bar{\boldsymbol{\theta}}} \boldsymbol{\psi}^*\left(\bar{\boldsymbol{\theta}}^t, \boldsymbol{\zeta}_m^*; \boldsymbol{z}_m^t\right). \tag{22}$$

Different from (20), the gradient here is obtained at the optimum $\zeta_m^* = \zeta_m^*(\overline{\theta}; z) := \arg \max_{\zeta \in \mathcal{Z}} \psi(\overline{\theta}, \zeta; z_m)$.

To establish convergence, one more assumption is needed. Assumption 6: Function f satisfies the following conditions.

- 1) Gradient estimates $\nabla_{\bar{\theta}} \psi^*(\bar{\theta}^t, \zeta_m^*; z_m)$ at ζ_m^* are unbiased and have bounded variance. That is, for $m = 1 \cdots M$, we have $\mathbb{E}\left[\nabla_{\bar{\theta}} \psi^*(\theta, \zeta_m^*; z_m) \nabla_{\bar{\theta}} f(\theta)\right] = \mathbf{0}$ and $\mathbb{E}\left[\|\nabla_{\bar{\theta}} \psi^*(\theta, \zeta_m^*; z_m) \nabla f(\theta)\right]^2\right] \leq \sigma^2$.
- 2) The expected norm of $g^t(\theta)$ is bounded, that is, $\mathbb{E}\|g^t(\bar{\theta})\|^2 \leq B^2$.

We can now present a theorem on the convergence of Alg. 2; see Appendix-C for the proof.

Theorem 2 (Convergence of Alg. 2): Let $\Delta_F := F(\bar{\theta}^0) - \inf_{\bar{\theta}} F(\bar{\theta})$, and D denote the diameter of the feasible set Θ . Under As. 1–4 and 6, for a constant step size $\alpha > 0$, and a fixed batch size M > 0, after T iterations, Alg. 2 satisfies

$$\mathbb{E}\left[\operatorname{dist}\left(\mathbf{0}, \partial F(\bar{\boldsymbol{\theta}}^T)\right)^2\right] \leq \frac{\upsilon}{T+1} \Delta_F + \frac{2L_{\boldsymbol{\theta}z}^2 \upsilon}{M} \left[(1-\alpha\mu)D^2 + \alpha^2 B^2 \right] + \frac{4\sigma^2}{M}$$
(23)

where v, v, and $\mu = \gamma_0 - L_{zz}$ are some positive constants.

Theorem 2 implies that the sequence $\{\bar{\theta}^I\}_{t=1}^T$ generated by Alg. 2 converges to a stationary point. The upper bound in (23) is characterized by a vanishing term induced by initial error Δ_F , and constant bias terms.

V. DISTRIBUTIONALLY ROBUST FEDERATED LEARNING

In practice, massive datasets are distributed geographically across multiple sites, where scalability, data privacy and integrity, as well as bandwidth scarcity typically discourage uploading them to a central server. This has propelled the so-called federated learning framework, where multiple workers exchange information with a server to learn a centralized

model using data locally generated and/or stored across workers [2]. Workers in this learning framework communicate *iteratively* with the server. Albeit appealing for its scalability, one needs to carefully address the bandwidth bottleneck associated with server-worker links. Furthermore, the workers' data may have (slightly) different distributions, which further challenges the learning task. To seek a model robust to distribution drifts across workers, we will adapt our novel SPGDA approach to design a privacy-respecting and robust algorithm.

To that end, consider K workers with each worker $k \in \mathcal{K}$ collecting samples $\{z_n(k)\}_{n=1}^N$. A globally shared model parameterized by θ is to be updated at the server by aggregating gradients computed locally per worker. For simplicity, we consider workers having the same number of samples N. The goal is to learn a single global model from stored data at all workers by minimizing the following objective function

$$\min_{\boldsymbol{\theta} \in \Omega} \bar{\mathbb{E}}_{z \sim \widehat{P}}[\ell(\boldsymbol{\theta}; z)] + r(\boldsymbol{\theta}) \tag{24}$$

where $\bar{\mathbb{E}}_{z\sim \widehat{\mathcal{P}}}[\ell(\theta;z)] := \frac{1}{NK} \sum_{n=1}^{N} \sum_{k=1}^{K} \ell(\theta,z_n(k))$. To endow the learned model with robustness against distributional uncertainties, our novel formulation will solve the following problem in a distributed fashion

$$\min_{\theta \in \Theta} \sup_{P \in \mathcal{P}} \mathbb{E}_{z \sim P}[\ell(\theta; z)] + r(\theta)$$
s. to.
$$\mathcal{P} := \left\{ P \middle| \sum_{k=1}^{K} W_c \Big(P, \widehat{P}^{(N)}(k) \Big) \le \rho \right\}$$
 (25)

where $W_c(P, \widehat{P}^{(N)}(k))$ denotes the Wasserstein distance between distribution P and the local $\widehat{P}^{(N)}(k)$, per worker k.

Clearly, the constraint $P \in \mathcal{P}$ couples the optimization in (25) across all workers. To offer distributed implementations, we resort to Prop. 1, to arrive at the equivalent reformulation

$$\min_{\theta \in \Theta} \inf_{\gamma \in \Gamma} \sum_{k=1}^{K} \left\{ \bar{\mathbb{E}}_{z(k) \sim \widehat{P}^{(N)}(k)} \left[\sup_{\zeta \in \mathcal{Z}} \{ \ell(\theta; \zeta) + \gamma(\rho - c(z(k), \zeta)) \} \right] \right\} + r(\theta).$$
(26)

Next, we present our communication- and computationefficient DRFL that builds on the SPGDA scheme in Section IV.

Specifically, our DRFL hinges on the fact that with fixed server parameters $\bar{\theta}^t := [\theta^{t\top}, \gamma^t]^\top$ per iteration t, the optimization problem becomes *separable* across all workers. Hence, upon receiving $\bar{\theta}^t$ from the server, each worker $k \in \mathcal{K}$: i) samples a minibatch $\mathcal{B}^t(k)$ of data from $\widehat{P}^{(N)}(k)$; ii) forms the *perturbed* loss $\psi_k(\bar{\theta}^t, \zeta; z) := \ell(\theta^t; \zeta) + \gamma^t(\rho - c(z, \zeta))$ for each $z \in \mathcal{B}^t(k)$; iii) lazily maximizes $\psi_k(\bar{\theta}^t, \zeta; z)$ over ζ using a single gradient ascent step to yield $\zeta(\bar{\theta}^t; z) = z + \eta_t \nabla_{\zeta} \psi_k(\bar{\theta}^t, \zeta; z)|_{\zeta=z}$; and, iv) sends the stochastic gradient $|\mathcal{B}^t(k)|^{-1} \sum_{z \in \mathcal{B}^t(k)} \nabla_{\bar{\theta}} \psi_k(\bar{\theta}^t, \zeta(\bar{\theta}^t; z); z)|_{\bar{\theta}=\bar{\theta}^t}$ back to the

²Further details regarding this assumption are provided in Section VII-D.

Algorithm 3: DRFL

```
Input: Initial guess \bar{\theta}^1, a set of workers \mathcal{K} with data samples \{z_n(k)\}_{n=1}^N per worker k \in \mathcal{K}, step size
                     sequence \{\alpha_t, \eta_t > 0\}_{t=1}^T
    Output: \bar{\theta}^{\dot{T}+1}
1 for t = 1, ..., T do
            Each worker:
2
            Samples a minibatch \mathcal{B}^t(k) of samples
3
            Given \bar{\theta}^t and z \in \mathcal{B}^t(k), forms local perturbed loss
                         \psi_k(\bar{\boldsymbol{\theta}}^t, \zeta; z) := \ell(\bar{\boldsymbol{\theta}}^t; \zeta) + \gamma^t(\rho - c(z, \zeta))
              Lazily maximizes \psi_k(\bar{\theta}^t, \zeta; z) over \zeta to find
                               \zeta(\bar{\theta}^t;z) = z + \eta_t \nabla_{\zeta} \psi_k(\bar{\theta}^t,\zeta;z)|_{\zeta=\zeta}
              Computes stochastic gradient
                           \frac{1}{|\mathcal{B}^{t}(k)|} \sum_{z \in \mathcal{B}(t)} \nabla_{\bar{\theta}} \psi_{k}(\bar{\theta}^{t}, \zeta(\bar{\theta}^{t}; z); z) \Big|_{\bar{\theta} = \bar{\theta}^{t}}
            and uploads to server
5
           Updates \bar{\theta}^t according to (27) Broadcasts \bar{\theta}^{t+1} to workers
s end
```

server. Upon receiving all local gradients, the server updates $\bar{\theta}^t$ using a proximal gradient descent step to find $\bar{\theta}^{t+1}$, that is

$$\bar{\theta}^{t+1} = \operatorname{prox}_{\alpha_{t}r} \left[\bar{\theta}^{t} - \frac{\alpha_{t}}{K} \sum_{k=1}^{K} \frac{1}{|\mathcal{B}^{t}(k)|} \times \sum_{z \in \mathcal{B}^{t}(k)} \nabla_{\bar{\theta}} \psi_{k} \left(\bar{\theta}^{t}, \zeta(\bar{\theta}^{t}; z); z \right) \Big|_{\bar{\theta} = \bar{\theta}^{t}} \right]$$
(27)

which is then broadcast to all workers to begin a new round of local updates. Our DRFL approach is tabulated in Alg. 3.

VI. NUMERICAL TESTS

To assess the performance in the presence of distribution drifts and adversarial perturbations, we will rely on empirical classification of standard MNIST and Fashion- (F-)MNIST datasets using standard cross-entropy loss. Specifically, we compare performance using models trained with empirical risk minimization (ERM), the fast-gradient method (FGSM) [13], its iterated variant (IFGM) [40], and the Wasserstein robust method (WRM) [34]. We further evaluate the testing performance using the projected gradient descent (PGD) attack [41]. We first test performance of SPGD with ϵ -accurate oracle, and SPGDA on these standard classification tasks.

A. SPGD With €-Accurate Oracle and SPGDA

The FGSM attack performs one step gradient update along the direction of the gradient's sign to find an adversarial sample; that is,

$$x_{\text{adv}} = \text{Clip}_{[-1,1]}\{x + \epsilon_{\text{adv}} \text{sign}(\nabla \ell_x(\theta; (x, y)))\}$$
 (28)

where $\epsilon_{\rm avd}$ controls the maximum ℓ_{∞} perturbation of adversarial samples. The element-wise ${\rm Clip}_{[a,b]}\{\}$ operator forces its input to reside in the prescribed range [-1,1]. By running $T_{\rm adv}$ iterations of (28) iterative (I) FGSM attack samples are generated [13]. Starting with an initialization $x_{\rm adv}^0 = x$, and considering the ℓ_{∞} norm, the PGD attack iterates [41]

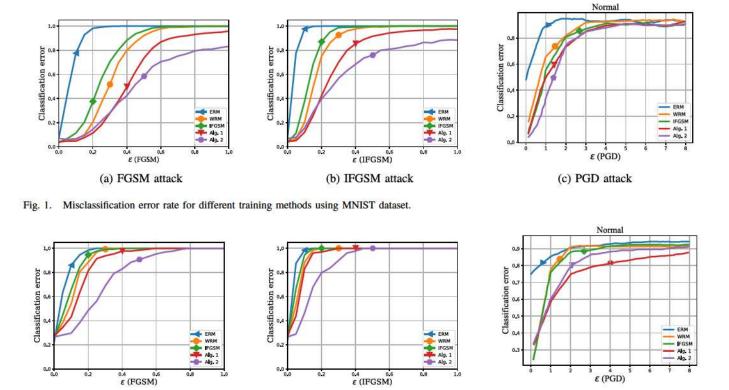
$$x_{\text{adv}}^{t+1} = \prod_{\mathcal{B}_{\epsilon}(x_{\text{adv}}^{t})} \left\{ x_{\text{adv}}^{t} + \alpha \text{sign} \left(\nabla \ell_{x}(\theta; (x_{\text{adv}}^{t}, y)) \right) \right\}$$
(29)

for T_{adv} steps, where Π denotes projection onto the ball $\mathcal{B}_{\epsilon}(x_{\text{adv}}^t) \coloneqq \{x : \|x - x_{\text{adv}}^t\|_{\infty} \le \epsilon_{\text{adv}}\}$, and $\alpha > 0$ is the stepsize set to 1 in our experiments. We use $T_{\text{adv}} = 10$ iterations for all iterative methods both in training and attack samples. The PGD can also be interpreted as an iterative algorithm that solves the optimization problem $\max_{x'} \ell(\theta; (x', y))$ subject to $\|x' - x\|_{\ell_{\infty}} \le \alpha$. The Wasserstein attack on the other hand, generates adversarial samples by solving a perturbed training loss with an ℓ_2 -based transportation cost associated with the Wasserstein distance between the training and adversarial data distributions [34].

Remark 7 (ϵ -accurate Oracle Implementation): To implement the ϵ -accurate oracle for Algorithm 1, several optimization algorithms can be employed. Here we rely on 10 iterations of gradient ascent with a constant step size $\eta=0.001$. While this method provides only an approximate solution, in theory for large enough number of iterations it can generate ϵ -accurate solution. This holds specifically because the feasible set for γ is such that the function $\psi(\cdot)$ is guaranteed to be strongly concave, and thus with at least $O(1/(\epsilon \eta))$ gradient ascent steps an ϵ -accurate solution can be obtained.

For the MNIST and F-MNIST datasets, a convolutional neural network (CNN) classifier consisting of 8×8 , 6×6 , and 5×5 filter layers with rectified linear units (ReLU) [42] and the same padding, is used. Its first, second, and third layers have 64, 128, and 128 channels, respectively, followed by a fully connected layer, and a softmax layer at the output.

CNNs with the same architecture are trained, using different adversarial samples. Specifically, to train a Wasserstein robust CNN model (WRM), $\gamma = 1$ was used to generate Wasserstein adversarial samples, ϵ_{adv} was set to 0.1 for the other two methods, and $\rho = 25$ was used to define the uncertainty set for both Algs. 1 and 2. Unless otherwise noted, we set the batch size to 128, the number of epochs to 30, the learning rates to $\alpha = 0.001$ and $\eta = 0.02$, and used the Adam optimizer [43]. Fig. 1(a) shows the classification error on the MNIST dataset. The error rates were obtained using testing samples generated according to the FGSM method with ϵ_{adv} . Clearly all training methods outperform ERM, and our proposed Algs. 1 and 2 offer improved performance over competing alternatives. The testing accuracy of all methods using samples generated according to an IFGSM attack is presented in Fig. 1(b). Likewise, Algs. 1 and 2 outperform other methods in this case. Fig. 1(c) depicts the testing accuracy of the considered methods under different levels of



(b) IFGSM attack

Fig. 2. Misclassification error rate for different training methods using F-MNIST dataset.

a PGD attack. The plots in Fig. 1 showcase the improved performance obtained by CNNs trained using Algs. 1 and 2.

(a) FGSM attack

The F-MNIST article image dataset is adopted in the second experiment. Similar to the MNIST dataset, each example in F-MNIST is also a 28×28 gray-scale image, associated with a label from 10 classes. F-MNIST is a modern replacement for the original MNIST dataset for benchmarking machine learning algorithms. Using CNNs with similar architectures as before, the classification error is depicted for different training methods in Fig. 2. Three different attacks, namely FGSM, IFGSM, and PGD are used during testing. The resulting classification error rates are reported in Figs. 2(a), 2(b), and 2(c), respectively. The proposed SPGD and SPGDA algorithms outperform the other training methods, verifying the superiority of Algs. 1 and 2 in terms of yielding robust models.

B. Extended Numerical Tests on CIFAR-10 and CIFAR-100

In this section, we present an extended evaluation of our approach on more practical tasks using large-scale pre-trained CNNs for classification. Specifically, we leverage well-established models such as ResNet-18 [44], WideResNet-50-20 [45], and VGG-11-BN [46] for both 10-class and 100-class challenges on the CIFAR dataset. These models are initially pretrained on the ImageNet dataset [47] and fine-tuned on CIFAR-10 and CIFAR-100. CIFAR-10 comprises a total of 60,000 color images of sizes 32 × 32 across 10 classes, with 50,000 training images and 10,000 for

testing. On the other hand, CIFAR-100 features 600 images per class, with 500 for training and 100 for testing per class. To ensure fair comparisons, we kept hyperparameters consistent across experiments, utilizing SGD with a fixed learning rate of 0.001 and a momentum of 0.9 for fine-tuning over 15 epochs.

(c) PGD attack

We conducted extensive adversarial testing using various attack strategies, including variance-tuning based (VNFIGSM) [48], ensemble adversarial (RFGSM) [49], transferable adversarial (TIFGSM) [50], transferable attack with input diversity (DIFGSM) [51], Nesterov accelerated scale invariance (NIFGSM) [52], Carlini and Wagner (CW) [53], KL-divergence based (TPGD) [54], and parameter-free attack (APGD) [55]. The hyperparameters used during adversarial testing in this case are reported in Table III in the Appendix.

The results are presented in Tables I and II. The normal testing accuracies of all models for CIFAR-10 dataset were above 90% and for CIFAR-100 were above 78%, which shows that models were trained to classify normal input data accurately. The results in Tables I and II reveal that models trained using ERM demonstrate increased vulnerability to adversarial attacks compared to models trained using our robust method (Algorithm 1). This aligns with the expectation of the inherent vulnerability of conventionally trained models to adversarial perturbations. However, a notable finding is the consistent improvement in model performance when transitioning from ERM to robust models. This transition enhances performance across most adversarial testing scenarios while causing only a negligible drop in normal testing accuracy. Similar results were observed for the 100-class classification

Models	FGSM	TPGD	CW	PGD	APGD	TIFGSM	NIFGSM	DIFGSM	RFGSM	VNIFGSM
ResNet-18-ERM	37.03	27.57	0.45	0.57	0.68	2.30	0.55	1.29	0.54	0.88
ResNet-18-Robust	40.13	31.66	2.16	8.10	11.03	13.76	7.81	8.74	6.76	8.19
VGG-11-BN-ERM	28.75	24.72	0.06	0.16	0.24	1.16	0.16	0.33	0.16	0.71
VGG-11-BN-Robust	31.13	25.11	1.03	3.11	5.31	7.45	3.46	4.10	2.83	3.89
WideResNet-50-20-ERM	42.97	44.52	1.06	1.40	1.55	4.04	1.46	2.71	1.35	2.1
WideResNet-50-20-Robust	44.17	45.14	2.71	10.07	11.83	15.13	11.17	11.85	8.72	11.73

TABLE II
TESTING ACCURACY RESULTS IN % FOR THE CIFAR-100 DATASET

Models	FGSM	TPGD	CW	PGD	APGD	TIFGSM	NIFGSM	DIFGSM	RFGSM	VNIFGSM
ResNet-18-ERM	17.50	19.17	1.13	18.88	17.61	18.47	18.16	20.09	18.86	19.36
ResNet-18-Robust	21.13	24.34	4.54	22.97	21.69	21.84	21.53	24.88	21.93	24.52
VGG-11-BN-ERM	17.47	17.82	2.52	16.81	12.15	20.37	16.95	19.27	16.81	17.89
VGG-11-BN-Robust	21.02	21.92	6.26	18.84	15.31	20.70	19.65	21.70	19.79	21.79
WideResNet-50-20-ERM	20.05	21.44	5.87	21.02	17.78	22.87	21.24	22.69	21.00	21.66
WideResNet-50-20-Robust	23.54	25.06	9.45	23.83	18.36	26.71	24.67	26.52	23.63	23.92

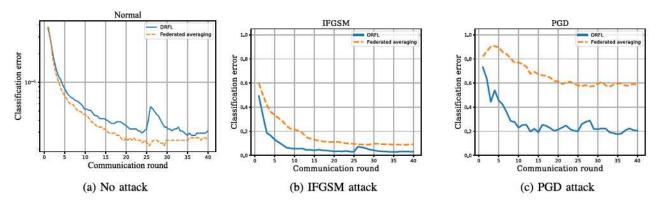


Fig. 3. Distributionally robust federated learning for image classification using the non-i.i.d. F-MNIST dataset.

problem, which are reported in Table II. The hyperparameters used for testing are reported in Table IV in the Appendix.

These observations underscore the practical significance of designing and implementing distributionally robust algorithms. Furthermore, the consistency of these improvements across a range of models, including ResNet, VGG, and WideResNet, across different classification tasks, highlights the effectiveness and generalizability of our robust training algorithms in enhancing model reliability.

C. Distributionally Robust Federated Learning

To validate the performance of our DRFL algorithm, we considered an FL environment consisting of a server and 10 workers, with local batch size 64, and assigned to every worker an equal-sized subset of training data containing i.i.d. samples from 10 different classes. All workers participated in each communication round. To benchmark the DRFL, we simulated the federated averaging method [56]. The testing accuracy on the MNIST dataset per communication round using clean (normal) images is depicted in Fig. 4(a). Clearly, both DRFL and federated averaging algorithms exhibit reasonable performance when the data is not corrupted. The performance is further tested against IFGSM and PGD

attacks with a fixed $\epsilon_{adv} = 0.1$ during each communication round, and the corresponding misclassification error rates are shown in Figs. 4(b) and 4(c), respectively. The classification performance using federated averaging does not improve in Fig. 4(b), whereas the DRFL performance keeps improving across communication rounds. This is a direct consequence of accounting for the data uncertainties during the learning process. Moreover, Fig. 4(c) showcases that the federated averaging becomes even worse as the model gets progressively trained under the PGD attack. This indeed motivates our DRFL approach when data are from untrusted entities with possibly adversarial input perturbations. Similarly, Fig. 5 depicts the misclassification rate of the proposed DRFL method compared with federated averaging, when using the F-MNIST dataset.

As the distribution of data across devices may influence performance, we further considered a biased local data setting. In particular, each worker k = 1, ..., 10 has data from only one class, so the distributions at workers are highly perturbed, and data stored across workers are thus non-i.i.d. The testing error rate for normal inputs is reported in Fig. 3(a), while the test error against adversarial attacks is depicted in Figs. 3(b) and 3(c). This additional set of tests shows that having distributional shifts across workers can indeed enhance testing performance when the samples are adversarially manipulated.

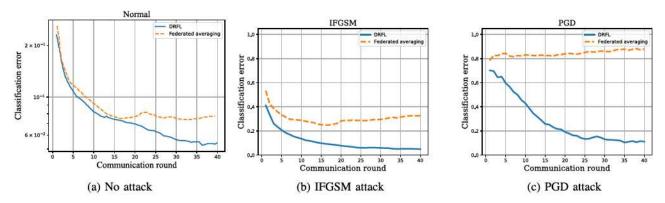


Fig. 4. Federated learning for image classification using the MNIST dataset.

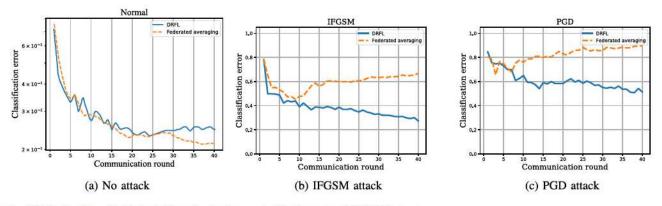


Fig. 5. Distributionally robust federated learning for image classification using F-MNIST dataset.

VII. CONCLUSION

A framework to robustify parametric machine learning models against distributional uncertainties was put forth in this paper. The learning task was cast as a distributionally robust optimization problem, for which two scalable stochastic optimization algorithms were developed. The first algorithm relies on an ϵ -accurate maximum-oracle to solve the inner convex subproblem, while the second approximates its solution via a single gradient ascent step. Convergence guarantees for both algorithms to a stationary point were obtained. The upshot of the proposed approach is that it is amenable to federated learning from unreliable datasets across multiple workers. The novel DRFL algorithm ensures data privacy and integrity, while offering robustness with minimal computational and communication overhead. Numerical tests for classifying standard real images showcased the merits of the proposed algorithms against distributional uncertainties and adversaries. This study uncovers promising avenues for future research, particularly in the domain of personalization. A key focus is on exploring the integration of distributionally robust learning with personalization. In this context, there is potential to extend this approach to heterogeneous datasets by incorporating the concept of multi-distribution Wasserstein distance [27]. This would involve utilizing two distinct models—one globally shared and another dedicated to aligning local data distributions within a common probability space. Further advancements in this research direction can be achieved by leveraging Bayesian and ensemble learning techniques to enhance both the robustness and personalization of federated learning methods.

APPENDIX

A. Proof of Lemma 1

Since function $\zeta \mapsto \psi(\bar{\theta}, \zeta; z)$ is λ -strongly concave, then $\zeta_*(\bar{\theta}) = \sup_{\zeta \in \mathcal{Z}} \psi(\bar{\theta}, \zeta; z)$ is unique. In addition, the first-order optimality condition gives $\langle \nabla_{\zeta} \psi(\bar{\theta}, \zeta_*\bar{\theta}); z), \zeta - \zeta_*(\bar{\theta}) \rangle \leq 0$. Let us define $\zeta_*^1 = \zeta_*(\bar{\theta}_1), \zeta_*^2 = \zeta_*(\bar{\theta}_2)$, and use the strong concavity for any $\bar{\theta}_1$ and $\bar{\theta}_2$, to write

$$\psi(\bar{\theta}_{2}, \zeta_{*}^{2}; z) \leq \psi(\bar{\theta}_{2}, \zeta_{*}^{1}; z) + \langle \nabla_{\zeta} \psi(\bar{\theta}_{2}, \zeta_{*}^{1}; z), \zeta_{*}^{2} - \zeta_{*}^{1} \rangle
- \frac{\lambda}{2} \|\zeta_{*}^{2} - \zeta_{*}^{1}\|^{2}$$
(30)

and

$$\psi(\bar{\theta}_{2}, \zeta_{*}^{1}; z) \leq \psi(\bar{\theta}_{2}, \zeta_{*}^{2}; z) + \langle \nabla_{\zeta} \psi(\bar{\theta}_{2}, \zeta_{*}^{2}; z), \zeta_{*}^{1} - \zeta_{*}^{2} \rangle
- \frac{\lambda}{2} \|\zeta_{*}^{2} - \zeta_{*}^{1}\|^{2}
\leq \psi(\bar{\theta}_{2}, \zeta_{*}^{2}; z) - \frac{\lambda}{2} \|\zeta_{*}^{2} - \zeta_{*}^{1}\|^{2}$$
(31)

where the last inequality is a consequence of the firstorder optimality condition. Summing (30) and (31), we find that

$$\begin{split} &\lambda \|\boldsymbol{\zeta}_{*}^{2} - \boldsymbol{\zeta}_{*}^{1}\|^{2} \leq \left\langle \nabla_{\zeta} \psi \left(\bar{\boldsymbol{\theta}}_{2}, \boldsymbol{\zeta}_{*}^{1}; \boldsymbol{z}\right), \boldsymbol{\zeta}_{*}^{2} - \boldsymbol{\zeta}_{*}^{1} \right\rangle \\ &\leq \left\langle \nabla_{\zeta} \psi \left(\bar{\boldsymbol{\theta}}_{2}, \boldsymbol{\zeta}_{*}^{1}; \boldsymbol{z}\right), \boldsymbol{\zeta}_{*}^{2} - \boldsymbol{\zeta}_{*}^{1} \right\rangle - \left\langle \nabla_{\zeta} \psi \left(\bar{\boldsymbol{\theta}}_{1}, \boldsymbol{\zeta}_{*}^{1}; \boldsymbol{z}\right), \boldsymbol{\zeta}_{*}^{2} - \boldsymbol{\zeta}_{*}^{1} \right\rangle \\ &= \left\langle \nabla_{\zeta} \psi \left(\bar{\boldsymbol{\theta}}_{2}, \boldsymbol{\zeta}_{*}^{1}; \boldsymbol{z}\right) - \nabla_{\zeta} \psi \left(\bar{\boldsymbol{\theta}}_{1}, \boldsymbol{\zeta}_{*}^{1}; \boldsymbol{z}\right), \boldsymbol{\zeta}_{*}^{2} - \boldsymbol{\zeta}_{*}^{1} \right\rangle. \end{split} \tag{32}$$

And using Hölder's inequality, we obtain that

$$\lambda \| \xi_{*}^{2} - \xi_{*}^{1} \|^{2}$$

$$\leq \| \nabla_{\zeta} \psi (\bar{\theta}_{2}, \xi_{*}^{1}; z) - \nabla_{\zeta} \psi (\bar{\theta}_{1}, \xi_{*}^{1}; z) \|_{\perp} \| \xi_{*}^{2} - \xi_{*}^{1} \|$$
(33)

from which we deduce

$$\|\xi_*^2 - \xi_*^1\| \le \frac{1}{\lambda} \|\nabla_{\xi} \psi(\bar{\theta}_2, \xi_*^1; z) - \nabla_{\xi} \psi(\bar{\theta}_1, \xi_*^1; z)\|_{\star}.$$
 (34)

Using $\psi(\bar{\theta}, \zeta; z) := \ell(\theta; \zeta) + \gamma(\rho - c(z, \zeta))$, we have that

$$\begin{split} \left\| \nabla_{\zeta} \psi \left(\bar{\theta}_{2}, \zeta_{*}^{1}; z \right) - \nabla_{\zeta} \psi \left(\bar{\theta}_{1}, \zeta_{*}^{1}; z \right) \right\|_{\star} \\ &= \left\| \nabla_{\zeta} \ell \left(\theta_{2}; \zeta_{*}^{1} \right) - \nabla_{\zeta} \ell \left(\theta_{1}; \zeta_{*}^{1} \right) \right. \\ &+ \gamma_{1} \nabla_{\zeta} c \left(z, \zeta_{*}^{1} \right) - \gamma_{2} \nabla_{\zeta} c \left(z, \zeta_{*}^{1} \right) \right\|_{\star} \\ &\leq \left\| \nabla_{\zeta} \ell \left(\theta_{2}; \zeta_{*}^{1} \right) - \nabla_{\zeta} \ell \left(\theta_{1}; \zeta_{*}^{1} \right) \right\|_{\star} \\ &+ \left\| \gamma_{1} \nabla_{\zeta} c \left(z, \zeta_{*}^{1} \right) - \gamma_{2} \nabla_{\zeta} c \left(z, \zeta_{*}^{1} \right) \right\|_{\star} \\ &\leq L_{z\theta} \left\| \theta_{2} - \theta_{1} \right\| + \left\| \nabla_{\zeta} c \left(z, \zeta_{*}^{1} \right) \right\|_{\star} \left\| \gamma_{2} - \gamma_{1} \right\|. \end{split}$$
(35)

Substituting (35) into (34), yields

$$\begin{aligned} \|\zeta_{*}^{2} - \zeta_{*}^{1}\| &\leq \frac{L_{z\theta}}{\lambda} \|\theta_{2} - \theta_{1}\| + \frac{1}{\lambda} \|\nabla_{\zeta} c(z, \zeta_{*}^{1})\|_{*} \|\gamma_{2} - \gamma_{1}\| \\ &\leq \frac{L_{z\theta}}{\lambda} \|\theta_{2} - \theta_{1}\| + \frac{L_{c}}{\lambda} \|\gamma_{2} - \gamma_{1}\| \end{aligned} (36)$$

where the last inequality holds because $\zeta \mapsto c(z, \zeta)$ is L_c -Lipschitz as per Assumption 3.

To obtain (12b), we first suppose without loss of generality that only a single datum z is given, and in order to prove existence of the gradient of $\bar{\psi}(\bar{\theta},z)$ with respect to $\bar{\theta}$, we resort to the Danskin's theorem as follows.

Danskin's Theorem [57]: Consider the following minimax optimization problem

$$\min_{\theta \in \Theta} \max_{t \in \mathcal{X}} f(\theta, \zeta) \tag{37}$$

where \mathcal{X} is a nonempty compact set, and $f: \Theta \times \mathcal{X} \to [0, \infty)$ is such that $f(\cdot, \zeta)$ is differentiable for any $\zeta \in \mathcal{X}$, and $\nabla_{\theta} f(\theta, \zeta)$ is continuous on $\Theta \times \mathcal{X}$. With $\mathcal{S}(\theta) := \{\zeta_* | \zeta_* = \arg \max_{\zeta} f(\theta, \zeta)\}$, the function

$$\bar{f}(\theta) \coloneqq \max_{\zeta \in \mathcal{Z}} f(\theta, \zeta)$$

is locally Lipschitz and directionally differentiable, where the directional derivatives satisfy

$$\bar{f}(\theta, d) = \sup_{\zeta \in \mathcal{S}(\theta)} \langle d, \nabla_{\theta} f(\theta, \zeta) \rangle. \tag{38}$$

For a given θ , if the set $S(\theta)$ is a singleton, then the function $\bar{f}(\theta)$ is differentiable at θ with gradient

$$\nabla_{\theta} \bar{f}(\theta) = \nabla_{\theta} f(\theta, \zeta_{*}(\theta)). \tag{39}$$

Given θ , and the μ -strongly convex $c(z,\cdot)$, function $\psi(\bar{\theta},\cdot;z)$ is concave if $L_{zz} - \gamma \mu < 0$, which holds true for $\gamma_0 > L_{zz}/\mu$. Replacing $\bar{f}(\theta,\zeta)$ with $\psi(\bar{\theta},\zeta;z)$, and given the concavity of

 $\zeta\mapsto \psi(\bar{\theta},\zeta;z),$ we have that $\bar{\psi}(\bar{\theta};z)$ is a continuous function with gradient

$$\nabla_{\bar{\theta}} \bar{\psi}(\bar{\theta}; z) = \nabla_{\bar{\theta}} \bar{\psi}(\bar{\theta}, \zeta_*(\bar{\theta}; z); z). \tag{40}$$

We can then obtain the second inequality, as

$$\begin{split} & \| \nabla_{\bar{\theta}} \psi \left(\bar{\theta}_{1}, \zeta_{*}^{1}; z \right) - \nabla_{\bar{\theta}} \psi \left(\bar{\theta}_{2}, \zeta_{*}^{2}; z \right) \| \\ & \leq \left\| \nabla_{\bar{\theta}} \psi \left(\bar{\theta}_{1}, \zeta_{*}^{1}; z \right) - \nabla_{\bar{\theta}} \psi \left(\bar{\theta}_{1}, \zeta_{*}^{2}; z \right) \right\| \\ & + \left\| \nabla_{\bar{\theta}} \psi \left(\bar{\theta}_{1}, \zeta_{*}^{2}; z \right) - \nabla_{\bar{\theta}} \psi \left(\bar{\theta}_{2}, \zeta_{*}^{2}; z \right) \right\| \\ & \leq \left\| \begin{bmatrix} \nabla_{\theta} \ell(\theta_{1}, \zeta_{*}^{1}) - \nabla_{\theta} \ell(\theta_{1}, \zeta_{*}^{2}) \\ c(z, \zeta_{*}^{2}) - c(z, \zeta_{*}^{1}) \end{bmatrix} \right\| \\ & + \left\| \begin{bmatrix} \nabla_{\theta} \ell(\theta_{1}, \zeta_{*}^{2}) - \nabla_{\theta} \ell(\theta_{2}, \zeta_{*}^{2}) \\ 0 \end{bmatrix} \right\| \\ & \leq L_{\theta z} \| \zeta_{*}^{1} - \zeta_{*}^{2} \| + L_{c} \| \zeta_{*}^{1} - \zeta_{*}^{2} \| + L_{\theta \theta} \| \theta_{1} - \theta_{2} \| \\ & \leq \left(L_{\theta \theta} + \frac{L_{\theta z} L_{z \theta} + L_{c} L_{z \theta}}{\lambda} \right) \| \theta_{2} - \theta_{1} \| \\ & + \frac{L_{\theta z} L_{c} + L_{c}^{2}}{\lambda} \| \gamma_{2} - \gamma_{1} \| \end{split}$$

$$(41)$$

where we again used inequality (36). As a technical note, if the considered model is a neural network with a non-smooth activation function, the loss will not be continuously differentiable. However, we will not encounter this challenge often in practice.

B. Proof of Theorem 1

With slight abuse of notation, define for convenience $F(\theta, \gamma) := f(\theta, \gamma) + r(\theta) + h(\gamma)$, where $h(\gamma)$ is the indicator function

$$h(\gamma) = \begin{cases} 0, & \text{if } \gamma \in \Gamma \\ \infty, & \text{if } \gamma \notin \Gamma \end{cases}$$
 (42)

with $\Gamma := \{\gamma | \gamma \ge \gamma_0\}$, and for ease of representation we use $\bar{r}(\bar{\theta}) := r(\theta) + h(\gamma)$. Having an L_f -smooth function f, yields

$$f(\bar{\boldsymbol{\theta}}^{t+1}) \leq f(\bar{\boldsymbol{\theta}}^{t}) + \left\langle \nabla f(\bar{\boldsymbol{\theta}}^{t}), \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\rangle + \frac{L_f}{2} \|\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t}\|^{2}.$$

$$(43)$$

For a given z^t , the gradients are

$$\begin{split} \boldsymbol{g}^* \Big(\bar{\boldsymbol{\theta}}^t \Big) &\coloneqq \begin{bmatrix} \nabla_{\boldsymbol{\theta}} \boldsymbol{\psi} \Big(\boldsymbol{\theta}^t, \boldsymbol{\gamma}, \boldsymbol{\zeta}_* (\bar{\boldsymbol{\theta}}^t; \boldsymbol{z}^t); \boldsymbol{z}^t \Big) \\ \partial_{\boldsymbol{\gamma}} \boldsymbol{\psi} \Big(\boldsymbol{\theta}^t, \boldsymbol{\gamma}, \boldsymbol{\zeta}_* (\bar{\boldsymbol{\theta}}^t; \boldsymbol{z}^t); \boldsymbol{z}^t \Big) \end{bmatrix} \\ &= \begin{bmatrix} \nabla_{\boldsymbol{\theta}} \boldsymbol{\psi} \Big(\boldsymbol{\theta}^t, \boldsymbol{\gamma}, \boldsymbol{\zeta}_* (\bar{\boldsymbol{\theta}}^t; \boldsymbol{z}^t); \boldsymbol{z}^t \Big) \\ \rho - c(\boldsymbol{z}^t, \boldsymbol{\zeta}_* (\bar{\boldsymbol{\theta}}^t; \boldsymbol{z}^t) \end{bmatrix}. \end{split}$$

and

$$g^{\epsilon}(\bar{\theta}^{t}) := \begin{bmatrix} \nabla_{\theta} \psi(\theta^{t}, \gamma, \zeta_{\epsilon}(\bar{\theta}^{t}; z^{t}); z^{t}) \\ \partial_{\gamma} \psi(\theta^{t}, \gamma, \zeta_{\epsilon}(\bar{\theta}^{t}; z^{t}); z^{t}) \end{bmatrix}$$
$$= \begin{bmatrix} \nabla_{\theta} \psi(\theta^{t}, \gamma, \zeta_{\epsilon}(\bar{\theta}^{t}; z^{t}); z^{t}) \\ \rho - c(z^{t}, \zeta_{\epsilon}(\bar{\theta}^{t}; z^{t}) \end{bmatrix}$$

obtained by an oracle at the optimal ζ_* and the ϵ -optimal ζ_ϵ solvers, respectively. Now, we define the error vector $\delta(\bar{\boldsymbol{\theta}}^t) := \nabla f(\bar{\boldsymbol{\theta}}^t) - g^{\epsilon}(\bar{\boldsymbol{\theta}}^t)$, and replace this into (43), to obtain

$$f(\bar{\boldsymbol{\theta}}^{t+1}) \leq f(\bar{\boldsymbol{\theta}}^{t}) + \left\langle g^{\epsilon}(\bar{\boldsymbol{\theta}}^{t}) + \delta(\bar{\boldsymbol{\theta}}^{t}), \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\rangle + \frac{L_{f}}{2} \|\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t}\|^{2}. \tag{44}$$

The following properties hold equivalently for the proximal operator, and for any x, y

$$u = \operatorname{prox}_{\alpha r}(x) \iff (x - u)^{T}(y - u) \le \alpha r(y) - \alpha r(u).$$
 (45)

With $u = \bar{\theta}^{t+1}$ and $x = \bar{\theta}^t - \alpha_t g^{\epsilon}(\bar{\theta}^t)$ in (45), it holds that

$$\left\langle \bar{\theta}^{t} - \alpha_{t} g^{\epsilon} \left(\bar{\theta}^{t} \right) - \bar{\theta}^{t+1}, \, \bar{\theta}^{t} - \bar{\theta}^{t+1} \right\rangle \leq \alpha_{t} \bar{r} \left(\bar{\theta}^{t} \right) - \alpha_{t} \bar{r} \left(\bar{\theta}^{t+1} \right)$$

and upon rearranging, we obtain

$$\left\langle g^{\epsilon}\left(\bar{\boldsymbol{\theta}}^{t}\right), \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t}\right\rangle \leq \bar{r}\left(\bar{\boldsymbol{\theta}}^{t}\right) - \bar{r}\left(\bar{\boldsymbol{\theta}}^{t+1}\right) - \frac{1}{\alpha_{t}} \|\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t}\|^{2}.$$
(46)

Adding inequalities in (46) and (44) gives

$$f(\bar{\boldsymbol{\theta}}^{t+1}) \leq f(\bar{\boldsymbol{\theta}}^{t}) + \left\langle \delta(\bar{\boldsymbol{\theta}}^{t}), \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\rangle + \frac{L_{f}}{2} \|\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t}\|^{2} + \bar{r}(\bar{\boldsymbol{\theta}}^{t}) - \bar{r}(\bar{\boldsymbol{\theta}}^{t+1}) - \frac{1}{\alpha_{t}} \|\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t}\|^{2}$$

and with $F(\bar{\theta}) := f(\bar{\theta}) + \bar{r}(\bar{\theta})$, we can write

$$F(\bar{\boldsymbol{\theta}}^{t+1}) - F(\bar{\boldsymbol{\theta}}^{t}) \le \left\langle \delta(\bar{\boldsymbol{\theta}}^{t}), \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\rangle + \left(\frac{L_f}{2} - \frac{1}{\alpha_t} \right) \|\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t}\|^2. \tag{47}$$

Using Young's inequality for any $\eta>0$ gives $\langle\delta(\bar{\theta}^t),\bar{\theta}^{t+1}-\bar{\theta}^t\rangle\leq \frac{\eta}{2}\|\bar{\theta}^{t+1}-\bar{\theta}^t\|^2+\frac{1}{2\eta}\|\delta(\bar{\theta}^t)\|^2$, and hence

$$F(\bar{\boldsymbol{\theta}}^{t+1}) - F(\bar{\boldsymbol{\theta}}^t) \le \left(\frac{L_f + \eta}{2} - \frac{1}{\alpha_t}\right) \|\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^t\|^2 + \frac{\|\delta(\bar{\boldsymbol{\theta}}^t)\|^2}{2\eta}. \tag{48}$$

Next, we will bound $\delta(\bar{\theta}^t) := \nabla f(\bar{\theta}^t) - g^{\epsilon}(\bar{\theta}^t)$. By adding and subtracting $g^*(\bar{\theta}^t)$ to the right hand side, we find

$$\left\|\delta\left(\bar{\boldsymbol{\theta}}^{t}\right)\right\|^{2} \leq 2\left\|\nabla f\left(\bar{\boldsymbol{\theta}}^{t}\right) - g^{*}\left(\bar{\boldsymbol{\theta}}^{t}\right)\right\|^{2} + 2\left\|g^{*}\left(\bar{\boldsymbol{\theta}}^{t}\right) - g^{\epsilon}\left(\bar{\boldsymbol{\theta}}^{t}\right)\right\|^{2}.\tag{49}$$

The Lipschitz smoothness of the gradient, implies that

$$\begin{aligned} & \left\| g^* \left(\bar{\theta}^t \right) - g^{\epsilon} \left(\bar{\theta}^t \right) \right\|^2 \\ & = \left\| \begin{bmatrix} \nabla_{\theta} \psi \left(\theta^t, \gamma, \zeta_* (\bar{\theta}^t; z^t); z^t \right) \\ \rho - c z^t, \zeta_* (\bar{\theta}^t; z^t) \end{bmatrix} \right] \end{aligned}$$

$$-\left[\begin{array}{c} \nabla_{\theta}\psi(\theta^{t},\gamma,\zeta_{\epsilon}(\bar{\theta}^{t};z^{t});z^{t}) \\ \rho - c(z^{t},\zeta_{\epsilon}(\bar{\theta}^{t};z^{t})) \end{array}\right]^{2}$$

$$= \left\|\nabla_{\theta}\psi(\theta^{t},\gamma,\zeta_{*}(\bar{\theta}^{t};z^{t});z^{t}) - \nabla_{\theta}\psi(\theta^{t},\gamma,\zeta_{\epsilon}(\bar{\theta}^{t};z^{t});z^{t})\right\|^{2}$$

$$+ \left\|c(z^{t},\zeta_{*}^{t}) - c(z^{t},\zeta_{\epsilon}^{t})\right\|^{2}$$

$$\stackrel{\text{(a)}}{\leq} \left(\frac{L_{\theta z}^{2}}{\lambda^{t}} + L_{c}\right) \|\zeta_{*}^{t} - \zeta_{\epsilon}^{t}\|^{2}$$

$$\stackrel{\text{(b)}}{\leq} \left(\frac{L_{\theta z}^{2}}{\lambda^{t}} + L_{c}\right) \epsilon$$

$$\leq \left(\frac{L_{\theta z}^{2}}{\lambda_{0}} + L_{c}\right) \epsilon$$

$$(50)$$

where (a) uses the $\lambda^t = \mu \gamma^t - L_{zz}$ strong-concavity of $\zeta \mapsto \psi(\bar{\theta}, \gamma, \zeta; z)$, and the second term is bounded by $L_c \|\zeta_*^t - \zeta_\epsilon^t\|^2$ according to Assumption 3. The last inequality holds for $\lambda_0 := \mu \gamma_0 - L_{zz}$, where we used (42) to bound $\gamma^t \ge \gamma_0 > L_{zz}$. So far, we have established that

$$\left\| g^* \left(\bar{\theta}^t \right) - g^{\epsilon} \left(\bar{\theta}^t \right) \right\|^2 \le \frac{L_{\bar{\theta}z}^2 \epsilon}{\lambda_0} \tag{51}$$

where for notational convenience we let $L_{\bar{\theta}_z}^2 := L_{\theta_z}^2 + \lambda_0 L_c$. Substituting (51) into (49), the error can be bounded as

$$\left\|\delta\left(\bar{\boldsymbol{\theta}}^{t}\right)\right\|^{2} \leq 2\left\|\nabla f\left(\bar{\boldsymbol{\theta}}^{t}\right) - g^{*}\left(\bar{\boldsymbol{\theta}}^{t}\right)\right\|^{2} + \frac{2L_{\bar{\boldsymbol{\theta}}z}^{2}\epsilon}{\lambda_{0}}.$$
 (52)

Combining (48) and (50) yields

$$F(\bar{\boldsymbol{\theta}}^{t+1}) - F(\bar{\boldsymbol{\theta}}^{t}) \le \left(\frac{L_f + \eta}{2} - \frac{1}{\alpha_t}\right) \left\|\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t}\right\|^2 + \frac{1}{\eta} \left\|\nabla f(\bar{\boldsymbol{\theta}}^{t}) - g^*(\bar{\boldsymbol{\theta}}^{t})\right\|^2 + \frac{L_{\bar{\boldsymbol{\theta}}z}^2 \epsilon}{\eta \lambda_0}. \tag{53}$$

Considering a constant step size α and summing these inequalities over t = 1, ..., T yields

$$\left(\frac{1}{\alpha} - \frac{L_f + \eta}{2}\right) \sum_{t=0}^{T} \left\| \bar{\theta}^{t+1} - \bar{\theta}^t \right\|^2 \le F\left(\bar{\theta}^0\right) - F\left(\bar{\theta}^T\right) \\
+ \frac{1}{\eta} \sum_{t=0}^{T} \left\| \nabla f\left(\bar{\theta}^t\right) - g^*\left(\bar{\theta}^t\right) \right\|^2 + \frac{(T+1)L_{\bar{\theta}z}^2 \epsilon}{\lambda_0}. \tag{54}$$

From the proximal gradient update

$$\bar{\boldsymbol{\theta}}^{t+1} = \arg\min_{\boldsymbol{\theta}} \, \alpha \bar{r}(\boldsymbol{\theta}) + \alpha \langle \boldsymbol{\theta} - \bar{\boldsymbol{\theta}}^t, g^{\epsilon} \left(\bar{\boldsymbol{\theta}}^t \right) \rangle + \frac{1}{2} \| \boldsymbol{\theta} - \bar{\boldsymbol{\theta}}^t \|^2$$
(55)

the optimality of $\bar{\theta}^{t+1}$ in (55), implies that

$$\bar{r}(\bar{\boldsymbol{\theta}}^{t+1}) + \langle \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t}, g^{\epsilon}(\bar{\boldsymbol{\theta}}^{t}) \rangle + \frac{1}{2\alpha} \left\| \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\|^{2} \leq \bar{r}(\bar{\boldsymbol{\theta}}^{t})$$

which combined with the smoothness of f (c.f. (43)) yields

$$\left\langle \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t}, \boldsymbol{g}^{\epsilon} \left(\bar{\boldsymbol{\theta}}^{t} \right) - \nabla f \left(\bar{\boldsymbol{\theta}}^{t} \right) \right\rangle + \left(\frac{1}{2\alpha} - \frac{L_{f}}{2} \right) \|\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \|^{2} \\
\leq F \left(\bar{\boldsymbol{\theta}}^{t} \right) - F \left(\bar{\boldsymbol{\theta}}^{t+1} \right) \tag{56}$$

Subtracting $\langle \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^t, \nabla f(\bar{\boldsymbol{\theta}}^{t+1}) \rangle$ from both sides gives $\left\langle \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^t, g^{\epsilon}(\bar{\boldsymbol{\theta}}^t) - \nabla f(\bar{\boldsymbol{\theta}}^{t+1}) \right\rangle + \left(\frac{1}{2\alpha} - \frac{L_f}{2} \right) \left\| \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^t \right\|^2$ $\leq F(\bar{\boldsymbol{\theta}}^t) - F(\bar{\boldsymbol{\theta}}^{t+1}) - \left\langle \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^t, \nabla f(\bar{\boldsymbol{\theta}}^{t+1}) - \nabla f(\bar{\boldsymbol{\theta}}^t) \right\rangle.$

Considering $\|g^{\epsilon}(\bar{\boldsymbol{\theta}}^t) - \nabla f(\bar{\boldsymbol{\theta}}^{t+1}) + \frac{1}{\alpha}(\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^t)\|^2$ on the left hand side, and adding relevant terms to the right hand side, we arrive at

$$\begin{aligned} & \left\| g^{\epsilon} \left(\bar{\theta}^{t} \right) - \nabla f \left(\bar{\theta}^{t+1} \right) + \frac{1}{\alpha} (\bar{\theta}^{t+1} - \bar{\theta}^{t}) \right\|^{2} \\ & \leq \left\| g^{\epsilon} \left(\bar{\theta}^{t} \right) - \nabla f \left(\bar{\theta}^{t+1} \right) \right\|^{2} + \frac{1}{\alpha^{2}} \left\| \bar{\theta}^{t+1} - \bar{\theta}^{t} \right\|^{2} \\ & + \left(\frac{L_{f}}{\alpha} - \frac{1}{\alpha^{2}} \right) \left\| \bar{\theta}^{t+1} - \bar{\theta}^{t} \right\|^{2} + \frac{2}{\alpha} \left(F \left(\bar{\theta}^{t} \right) - F \left(\bar{\theta}^{t+1} \right) \right) \\ & - \frac{2}{\alpha} \left\langle \bar{\theta}^{t+1} - \bar{\theta}^{t}, \nabla f \left(\bar{\theta}^{t+1} \right) - \nabla f \left(\bar{\theta}^{t} \right) \right\rangle \\ & \leq \left\| g^{\epsilon} \left(\bar{\theta}^{t} \right) - \nabla f \left(\bar{\theta}^{t} \right) \right\|^{2} + \frac{1}{\alpha^{2}} \left\| \bar{\theta}^{t+1} - \bar{\theta}^{t} \right\|^{2} \\ & + \left(\frac{L_{f}}{\alpha} - \frac{1}{\alpha^{2}} \right) \left\| \bar{\theta}^{t+1} - \bar{\theta}^{t} \right\|^{2} + \frac{2}{\alpha} \left(F \left(\bar{\theta}^{t} \right) - F \left(\bar{\theta}^{t+1} \right) \right) \\ & - \frac{2}{\alpha} \left\langle \bar{\theta}^{t+1} - \bar{\theta}^{t} \right\rangle, \nabla f \left(\bar{\theta}^{t+1} \right) - \nabla f \left(\bar{\theta}^{t} \right) \right\rangle \\ & \leq \left\| g^{\epsilon} \left(\bar{\theta}^{t} \right) - \nabla f \left(\bar{\theta}^{t} \right) \right\|^{2} + \frac{1}{\alpha^{2}} \left\| \bar{\theta}^{t+1} - \bar{\theta}^{t} \right\|^{2} \\ & + \left(\frac{L_{f}}{\alpha} - \frac{1}{\alpha^{2}} \right) \left\| \bar{\theta}^{t+1} - \bar{\theta}^{t} \right\|^{2} + \frac{2}{\alpha} \left(F \left(\bar{\theta}^{t} \right) - F \left(\bar{\theta}^{t+1} \right) \right) \\ & + \frac{\eta}{\alpha} \left\| \bar{\theta}^{t+1} - \bar{\theta}^{t} \right\|^{2} + \frac{L_{f}^{2}}{\eta} \left\| \bar{\theta}^{t+1} - \bar{\theta}^{t} \right\|^{2} \end{aligned}$$

$$(57)$$

where the last inequality is obtained by applying Young's inequality, and then using the L_f -Lipschitz continuity of $f(\cdot)$. By simplifying the last inequality, we obtain

$$\left\| g^{\epsilon} \left(\bar{\boldsymbol{\theta}}^{t} \right) - \nabla f \left(\bar{\boldsymbol{\theta}}^{t+1} \right) + \frac{1}{\alpha} \left(\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right) \right\|^{2}$$

$$\leq \left\| g^{\epsilon} \left(\bar{\boldsymbol{\theta}}^{t} \right) - \nabla f \left(\bar{\boldsymbol{\theta}}^{t} \right) \right\|^{2} + \frac{2}{\alpha} \left(F \left(\bar{\boldsymbol{\theta}}^{t} \right) - F \left(\bar{\boldsymbol{\theta}}^{t+1} \right) \right)$$

$$+ \left(\frac{L_{f}^{2}}{\eta} + \frac{L_{f} + \eta}{\alpha} \right) \left\| \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\|^{2}. \tag{58}$$

The first term in the right hand side can be bounded by adding and subtracting $g^*(\bar{\theta}^t)$ and using (51), to arrive at

$$\left\| g^{\epsilon} \left(\bar{\boldsymbol{\theta}}^{t} \right) - \nabla f \left(\bar{\boldsymbol{\theta}}^{t+1} \right) + \frac{1}{\alpha} \left(\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right) \right\|^{2}$$

$$\leq 2 \left\| \nabla f \left(\bar{\boldsymbol{\theta}}^{t} \right) - g^{*} \left(\bar{\boldsymbol{\theta}}^{t} \right) \right\|^{2} \frac{2L_{\bar{\boldsymbol{\theta}}}^{2} \epsilon}{\lambda_{0}} + \frac{2}{\alpha} \left(F \left(\bar{\boldsymbol{\theta}}^{t} \right) - F \left(\bar{\boldsymbol{\theta}}^{t+1} \right) \right)$$

$$+ \left(\frac{L_{f}^{2}}{\eta} + \frac{L_{f} + \eta}{\alpha} \right) \left\| \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\|^{2}. \tag{59}$$

Summing these inequalities over t = 1, ..., T, we find

$$\sum_{t=0}^{T} \left\| g^{\epsilon} \left(\bar{\boldsymbol{\theta}}^{t} \right) - \nabla f \left(\bar{\boldsymbol{\theta}}^{t+1} \right) + \frac{1}{\alpha} \left(\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right) \right\|^{2}$$

$$\leq 2\sum_{t=0}^{T} \left\| \nabla f(\bar{\boldsymbol{\theta}}^{t}) - g^{*}(\bar{\boldsymbol{\theta}}^{t}) \right\|^{2} + \frac{2(T+1)L_{\bar{\boldsymbol{\theta}}}^{2}\epsilon}{\lambda_{0}} + \frac{2}{\alpha} \left[F(\bar{\boldsymbol{\theta}}^{0}) - F(\bar{\boldsymbol{\theta}}^{T}) \right] + \left(\frac{L_{f}^{2}}{\eta} + \frac{L_{f} + \eta}{\alpha} \right) \sum_{t=0}^{T} \left\| \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\|^{2}.$$
 (60)

Using (54) to bound the last term yields

$$\sum_{t=0}^{T} \left\| g^{\epsilon} \left(\bar{\boldsymbol{\theta}}^{t} \right) - \nabla f \left(\bar{\boldsymbol{\theta}}^{t+1} \right) + \frac{1}{\alpha} \left(\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right) \right\|^{2} \\
\leq 2 \sum_{t=0}^{T} \left\| \nabla f \left(\bar{\boldsymbol{\theta}}^{t} \right) - g^{*} \left(\bar{\boldsymbol{\theta}}^{t} \right) \right\|^{2} + \frac{2(T+1)L_{\bar{\boldsymbol{\theta}}}^{2} \epsilon}{\lambda_{0}} \\
+ \frac{2}{\alpha} \Delta_{F} + \beta \Delta_{F} + \frac{\beta}{\eta} \sum_{t=0}^{T} \left\| \nabla f \left(\bar{\boldsymbol{\theta}}^{t} \right) - g^{*} \left(\bar{\boldsymbol{\theta}}^{t} \right) \right\|^{2} \\
+ \frac{\beta(T+1)L_{\bar{\boldsymbol{\theta}}z}^{2} \epsilon}{\lambda_{0}} \tag{61}$$

where $\beta = (\frac{L_f^2}{\eta} + \frac{L_f + \eta}{\alpha}) \frac{2\alpha}{2 - (L_f + \eta)\alpha}$. By taking expectation of both sides of this inequality, we obtain

$$\frac{1}{T+1} \mathbb{E} \left[\sum_{t=0}^{T} \left\| g^{\epsilon} \left(\bar{\boldsymbol{\theta}}^{t} \right) - \nabla f \left(\bar{\boldsymbol{\theta}}^{t+1} \right) + \frac{1}{\alpha} \left(\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right) \right\|^{2} \right] \\
\leq \left(\frac{2}{\alpha} + \beta \right) \frac{\Delta_{F}}{T+1} + \left(\frac{\beta}{\eta} + 2 \right) \sigma^{2} + \frac{(\beta+2)L_{\tilde{\boldsymbol{\theta}}}^{2} \epsilon}{\lambda_{0}} \tag{62}$$

where we have used $\mathbb{E}[\|\nabla f(\bar{\theta}^t) - g^*(\bar{\theta}^t)\|_2^2] \le \sigma^2$, which holds according to Assumption 5. By [39, Th. 10] and [58], we know that

$$-g^{\epsilon}(\bar{\theta}^{t}) - \frac{1}{\alpha}(\bar{\theta}^{t+1} - \bar{\theta}^{t}) \in \partial \bar{r}(\bar{\theta}^{t+1})$$
(63)

which gives

$$\nabla f(\bar{\boldsymbol{\theta}}^{t+1}) - g^{\epsilon}(\bar{\boldsymbol{\theta}}^{t}) - \frac{1}{\alpha}(\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t})$$

$$\in \nabla f(\bar{\boldsymbol{\theta}}^{t+1}) + \partial \bar{r}(\bar{\boldsymbol{\theta}}^{t+1})$$

$$= \partial F(\bar{\boldsymbol{\theta}}^{t+1}).$$

Upon replacing the latter in the left hand side of (62), and recalling the definition of distance, we deduce that

$$\mathbb{E}\left[\operatorname{dist}(\mathbf{0}, \partial \hat{F}(\bar{\boldsymbol{\theta}}^{t'}))\right] \leq \left(\frac{2}{\alpha} + \beta\right) \frac{\Delta_F}{T} + \left(\frac{\beta}{\eta} + 2\right) \sigma^2 + \frac{(\beta + 2)L_{\bar{\theta}z}^2 \epsilon}{\lambda_0}$$

where t' is randomly drawn from $t' \in \{1, 2, ..., T+1\}$, which concludes the proof.

C. Proof of Theorem 2

Instead of resorting to an oracle to obtain an ϵ -optimal solver for the surrogate loss, here we utilize a single step

stochastic gradient ascent with mini-batch size M to solve the maximization step. Consequently, the updates become

$$\bar{\boldsymbol{\theta}}^{t+1} = \operatorname{prox}_{\alpha_t r} \left(\bar{\boldsymbol{\theta}}^t - \alpha_t g^t (\bar{\boldsymbol{\theta}}^t) \right) \tag{64}$$

where $g^t(\bar{\theta}^t) := \frac{1}{M} \sum_{m=1}^M g(\bar{\theta}^t, \zeta_m^t; z_m)$. Letting $\delta(\bar{\theta}^t) := \nabla f(\bar{\theta}^t) - g^t(\bar{\theta}^t)$, and using the L_f -smoothness of $f(\bar{\theta})$, we obtain

$$f(\bar{\boldsymbol{\theta}}^{t+1}) \leq f(\bar{\boldsymbol{\theta}}^{t}) + \left\langle \nabla f(\bar{\boldsymbol{\theta}}^{t}), \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\rangle + \frac{L_{f}}{2} \left\| \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\|^{2}$$

$$\leq f(\bar{\boldsymbol{\theta}}^{t}) + \left\langle g^{t}(\bar{\boldsymbol{\theta}}^{t}) + \delta(\bar{\boldsymbol{\theta}}^{t}), \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\rangle$$

$$+ \frac{L_{f}}{2} \left\| \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\|^{2}. \tag{65}$$

Next, we substitute $\bar{\theta}^{t+1} \to u$, $\theta^t \to y$, and $\bar{\theta}^t - \alpha_t g^t(\bar{\theta}^t) \to x$ in (45), to arrive at

$$\left\langle \bar{\theta}^{t} - \alpha_{t} g^{t} \left(\bar{\theta}^{t} \right) - \bar{\theta}^{t+1}, \, \bar{\theta}^{t} - \bar{\theta}^{t+1} \right\rangle \leq \alpha_{t} \bar{r} \left(\bar{\theta}^{t} \right) - \alpha_{t} \bar{r} \left(\bar{\theta}^{t+1} \right)$$

which leads to

$$\left\langle g^{t}\left(\bar{\theta}^{t}\right), \bar{\theta}^{t+1} - \bar{\theta}^{t}\right\rangle \leq \bar{r}\left(\bar{\theta}^{t}\right) - \bar{r}\left(\bar{\theta}^{t+1}\right) - \frac{1}{\alpha_{t}}\left\|\bar{\theta}^{t+1} - \bar{\theta}^{t}\right\|^{2}.$$

Substituting the latter into (65), gives

$$f(\bar{\boldsymbol{\theta}}^{t+1}) \leq f(\bar{\boldsymbol{\theta}}^{t}) + \left\langle \delta(\bar{\boldsymbol{\theta}}^{t}), \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\rangle + \frac{L_{f}}{2} \|\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t}\|^{2} + \bar{r}(\bar{\boldsymbol{\theta}}^{t}) - \bar{r}(\bar{\boldsymbol{\theta}}^{t+1}) - \frac{1}{\alpha_{t}} \|\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t}\|^{2}$$

and with $F(\theta) := f(\theta) + \bar{r}(\theta)$, we have

$$F(\bar{\boldsymbol{\theta}}^{t+1}) - F(\bar{\boldsymbol{\theta}}^{t}) \le \left\langle \delta(\bar{\boldsymbol{\theta}}^{t}), \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\rangle + \left(\frac{L_f}{2} - \frac{1}{\alpha_t} \right) \left\| \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\|^2. \tag{66}$$

Using Young's inequality $\langle \delta(\bar{\theta}^t), \bar{\theta}^{t+1} - \bar{\theta}^t \rangle \leq \frac{1}{2} \|\delta(\bar{\theta}^t)\|^2 + \frac{1}{2} \|\bar{\theta}^{t+1} - \bar{\theta}^t\|^2$ implies that

$$F(\bar{\boldsymbol{\theta}}^{t+1}) - F(\bar{\boldsymbol{\theta}}^t) \le \left(\frac{L_f + 1}{2} - \frac{1}{\alpha_t}\right) \left\|\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^t\right\|^2 + \frac{\|\delta(\bar{\boldsymbol{\theta}}^t)\|^2}{2} \tag{67}$$

and after adding the term $\langle \bar{\theta}^{t+1} - \bar{\theta}^t, \nabla f(\bar{\theta}^{t+1}) \rangle$ to both sides in (67), and simplifying terms, yields

$$\left\langle \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t}, g^{t} \left(\bar{\boldsymbol{\theta}}^{t} \right) - \nabla f \left(\bar{\boldsymbol{\theta}}^{t+1} \right) \right\rangle \\
\leq - \left(\frac{1}{2\alpha_{t}} - \frac{L_{f}}{2} \right) \left\| \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\|^{2} + F \left(\bar{\boldsymbol{\theta}}^{t} \right) - F \left(\bar{\boldsymbol{\theta}}^{t+1} \right) \\
- \left\langle \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t}, \nabla f \left(\bar{\boldsymbol{\theta}}^{t+1} \right) - \nabla f \left(\bar{\boldsymbol{\theta}}^{t} \right) \right\rangle. \tag{68}$$

Completing the square yields

$$\begin{split} & \left\| g^{t} \left(\bar{\boldsymbol{\theta}}^{t} \right) - \nabla f \left(\bar{\boldsymbol{\theta}}^{t+1} \right) + \frac{1}{\alpha_{t}} \left(\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right) \right\|^{2} \\ & \leq \left\| g^{t} \left(\bar{\boldsymbol{\theta}}^{t} \right) - \nabla f \left(\bar{\boldsymbol{\theta}}^{t+1} \right) \right\|^{2} + \frac{1}{\alpha_{t}^{2}} \left\| \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\|^{2} \end{split}$$

$$+\left(\frac{L_{f}}{\alpha_{t}}-\frac{1}{\alpha_{t}^{2}}\right)\left\|\bar{\theta}^{t+1}-\bar{\theta}^{t}\right\|^{2}+\frac{2\left(F\left(\bar{\theta}^{t}\right)-F\left(\bar{\theta}^{t+1}\right)\right)}{\alpha_{t}}$$

$$-\frac{2}{\alpha_{t}}\left(\bar{\theta}^{t+1}-\bar{\theta}^{t},\nabla f\left(\bar{\theta}^{t+1}\right)-\nabla f\left(\bar{\theta}^{t}\right)\right)$$

$$\leq 2\left\|g^{t}\left(\bar{\theta}^{t}\right)-\nabla f\left(\bar{\theta}^{t}\right)\right\|^{2}+2\left\|\nabla f\left(\bar{\theta}^{t}\right)-\nabla f\left(\bar{\theta}^{t+1}\right)\right\|^{2}$$

$$+\frac{1}{\alpha_{t}^{2}}\left\|\bar{\theta}^{t+1}-\bar{\theta}^{t}\right\|^{2}+\left(\frac{L_{f}}{\alpha_{t}}-\frac{1}{\alpha_{t}^{2}}\right)\left\|\bar{\theta}^{t+1}-\bar{\theta}^{t}\right\|^{2}$$

$$+\frac{2\left(F\left(\bar{\theta}^{t}\right)-F\left(\bar{\theta}^{t+1}\right)\right)}{\alpha_{t}}$$

$$-\frac{2}{\alpha_{t}}\left(\bar{\theta}^{t+1}-\bar{\theta}^{t},\nabla f\left(\bar{\theta}^{t+1}\right)-\nabla f\left(\bar{\theta}^{t}\right)\right)$$

$$\leq 2\left\|g^{t}\left(\bar{\theta}^{t}\right)-\nabla f\left(\bar{\theta}^{t}\right)\right\|^{2}+2L_{f}^{2}\left\|\bar{\theta}^{t+1}-\bar{\theta}^{t}\right\|^{2}$$

$$+\frac{1}{\alpha_{t}^{2}}\left\|\bar{\theta}^{t+1}-\bar{\theta}^{t}\right\|^{2}+\left(\frac{L_{f}}{\alpha_{t}}-\frac{1}{\alpha_{t}^{2}}\right)\left\|\bar{\theta}^{t+1}-\bar{\theta}^{t}\right\|^{2}$$

$$+\frac{2\left(F\left(\bar{\theta}^{t}\right)-F\left(\bar{\theta}^{t+1}\right)\right)}{\alpha_{t}}+\frac{2L_{f}}{\alpha_{t}}\left\|\bar{\theta}^{t+1}-\bar{\theta}^{t}\right\|^{2}$$

$$\leq 2\left\|g^{t}\left(\bar{\theta}^{t}\right)-\nabla f\left(\bar{\theta}^{t}\right)\right\|^{2}+\frac{2\left(F\left(\bar{\theta}^{t}\right)-F\left(\bar{\theta}^{t+1}\right)\right)}{\alpha_{t}}$$

$$+\frac{3L_{f}+2L_{f}^{2}\alpha_{t}}{\alpha_{t}}\left\|\bar{\theta}^{t+1}-\bar{\theta}^{t}\right\|^{2}.$$
(69)

Recalling that $\delta(\bar{\theta}^t) := \nabla f(\bar{\theta}^t) - g^t(\bar{\theta}^t)$, we can bound the first term as

$$\mathbb{E}\left[\left\|g^{t}(\bar{\boldsymbol{\theta}}^{t}) - \nabla f(\bar{\boldsymbol{\theta}}^{t})\right\|^{2} | \boldsymbol{\theta}^{t}\right]$$

$$= \mathbb{E}\left[\left\|g^{*}(\bar{\boldsymbol{\theta}}^{t}) - \nabla f(\bar{\boldsymbol{\theta}}^{t}) + \delta^{t}\right\|^{2} | \boldsymbol{\theta}^{t}\right]$$

$$= \left\|g^{*}(\bar{\boldsymbol{\theta}}^{t}) - \nabla f(\bar{\boldsymbol{\theta}}^{t})\right\|^{2} + \left\|\delta^{t}\right\|^{2}$$

$$+ 2\mathbb{E}\left[\left\langle g^{*}(\bar{\boldsymbol{\theta}}^{t}) - \nabla f(\bar{\boldsymbol{\theta}}^{t}), \delta^{t}\right\rangle | \boldsymbol{\theta}^{t}\right]$$
(70)

where the third equality is obtained by expanding the square term, and using $\mathbb{E}[\langle g^*(\bar{\theta}^t) - \nabla f(\bar{\theta}^t), \delta^t \rangle | \bar{\theta}^t] = \mathbf{0}$. We will further bound the right hand side here as follows. Recalling that $\delta^t = \frac{1}{M} \sum_{m=1}^M g(\bar{\theta}^t, \zeta_m^t; z_m) - g^*(\bar{\theta}^t)$, where $g^*(\theta^t) := \frac{1}{M} \sum_{m=1}^M \nabla_{\bar{\theta}} \psi(\bar{\theta}^t, \zeta_m^{*t}; z_m)$, it holds that

$$\mathbb{E}\left[\|\delta^{t}\|^{2}\left|\bar{\theta}^{t},\zeta_{m}^{t}\right]\right] \\
= \mathbb{E}\left[\left\|\frac{1}{M}\sum_{m=1}^{M}\left[g(\bar{\theta}^{t},\zeta_{m}^{t};z_{m})-g^{*}(\bar{\theta}^{t})\right]\right\|^{2}\left|\bar{\theta}^{t},\zeta_{m}^{t}\right]\right] \\
= \frac{1}{M^{2}}\sum_{m=1}^{M}\mathbb{E}\left[\left\|\nabla_{\bar{\theta}}\psi(\bar{\theta}^{t},\zeta_{m}^{t};z_{m})-\nabla_{\bar{\theta}}\psi(\bar{\theta}^{t},\zeta_{m}^{*t};z_{m})\right\|^{2}\left|\bar{\theta}^{t},\zeta_{m}^{t}\right]\right] \\
\leq \frac{L_{\theta z}^{2}}{M^{2}}\sum_{m=1}^{M}\left\|\zeta_{m}^{t}-\zeta_{m}^{*t}\right\|^{2} \tag{71}$$

where the second equality is because the samples $\{z_m\}_{m=1}^M$ are i.i.d., and last inequality holds due to the Lipschitz smoothness

of $\psi(\cdot)$. Since ζ_m^t is obtained by a single gradient ascent update over a μ -strongly concave function, we have that

$$\frac{L_{\theta z}^{2}}{M^{2}} \sum_{m=1}^{M} \left\| \xi_{m}^{t} - \xi_{m}^{*t} \right\|^{2} \leq \frac{L_{\theta z}^{2}}{M} \left[(1 - \alpha_{t} \mu) D^{2} + \alpha_{t}^{2} B^{2} \right]$$
 (72)

where D is the diameter of the feasible set, and $\alpha_t > 0$ is the step size. The following holds for the expected error term

$$\mathbb{E}\Big[\|\delta^t\|^2\big|\bar{\theta}^t,\zeta_m^t\Big] \le \frac{L_{\theta z}^2}{M}\Big[(1-\alpha_t\mu)D^2 + \alpha_t^2B^2\Big] \tag{73}$$

and using it in (70), we arrive at

$$\mathbb{E}\Big[\left\|\mathbf{g}^{t}\left(\bar{\boldsymbol{\theta}}^{t}\right) - \nabla f\left(\bar{\boldsymbol{\theta}}^{t}\right)\right\|^{2} \middle|\boldsymbol{\theta}^{t}\Big] \leq 2\left\|\mathbf{g}^{*}\left(\bar{\boldsymbol{\theta}}^{t}\right) - \nabla f\left(\bar{\boldsymbol{\theta}}^{t}\right)\right\|^{2} + \frac{L_{\bar{\boldsymbol{\theta}}z}^{2}}{M}\Big[(1 - \alpha_{t}\mu)D^{2} + \alpha_{t}^{2}B^{2}\Big]. \tag{74}$$

Substituting the last inequality into (69) boils down to

$$\mathbb{E}\left[\left\|g^{t}\left(\bar{\theta}^{t}\right) - \nabla f\left(\bar{\theta}^{t+1}\right) + \frac{1}{\alpha_{t}}\left(\bar{\theta}^{t+1} - \bar{\theta}^{t}\right)\right\|^{2} |\bar{\theta}^{t}] \\
\leq 4\left\|g^{*}\left(\bar{\theta}^{t}\right) - \nabla f\left(\bar{\theta}^{t}\right)\right\|^{2} \\
+ \frac{3L_{f} + 2L_{f}^{2}\alpha_{t}}{\alpha_{t}} \mathbb{E}\left[\left\|\bar{\theta}^{t+1} - \bar{\theta}^{t}\right\|^{2} |\bar{\theta}^{t}\right] \\
+ \frac{2F\left(\bar{\theta}^{t}\right) - 2\mathbb{E}\left[F\left(\bar{\theta}^{t+1}\right)|\bar{\theta}^{t}\right]}{\alpha_{t}} \\
+ \frac{L_{\frac{\bar{\theta}z}}^{2}}{M}\left[(1 - \alpha_{t}\mu)D^{2} + \alpha_{t}^{2}B^{2}\right]. \tag{75}$$

Taking again expectation over $\bar{\theta}^t$ on both sides, yields

$$\mathbb{E} \left\| g^{t} \left(\bar{\theta}^{t} \right) - \nabla f \left(\bar{\theta}^{t+1} \right) + \frac{1}{\alpha_{t}} \left(\bar{\theta}^{t+1} - \bar{\theta}^{t} \right) \right\|^{2} \\
\leq 4 \mathbb{E} \left[\left\| g^{*} \left(\bar{\theta}^{t} \right) - \nabla f \left(\bar{\theta}^{t} \right) \right\|^{2} \right] + \frac{L_{\bar{\theta}z}^{2}}{M} \left[(1 - \alpha_{t}\mu) D^{2} + \alpha_{t}^{2} B^{2} \right] \\
+ \mathbb{E} \left[\frac{2F \left(\bar{\theta}^{t} \right) - 2F \left(\bar{\theta}^{t+1} \right)}{\alpha_{t}} + \frac{3L_{f} + 2L_{f}^{2} \alpha_{t}}{\alpha_{t}} \left\| \bar{\theta}^{t+1} - \bar{\theta}^{t} \right\|^{2} \right]. \tag{76}$$

Recalling that $\mathbb{E}[\|\psi^*(\bar{\theta}^t, \zeta_m^t; z_m) - \nabla f(\bar{\theta}^t)\|^2] \leq \sigma^2$, and that $g^*(\bar{\theta}^t) = \frac{1}{M} \sum_{m=1}^M \psi(\bar{\theta}^t, \zeta_m^{*t}; z_m)$, the first term on the right hand side can be bounded by $\frac{4\sigma^2}{M}$. For a fixed learning rate $\alpha > 0$, summing inequalities (76) from $t = 0, \ldots, T$, yields

$$\begin{split} &\frac{1}{T+1} \mathbb{E} \bigg[\sum_{t=0}^{T} \left\| g^{t} \Big(\bar{\boldsymbol{\theta}}^{t} \Big) - \nabla f \Big(\bar{\boldsymbol{\theta}}^{t+1} \Big) + \frac{1}{\alpha_{t}} \Big(\bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \Big) \right\|^{2} \bigg] \\ &\leq \frac{2}{\alpha(T+1)} \bigg(F \Big(\boldsymbol{\theta}^{0} \Big) - \mathbb{E} \big[F \Big(\boldsymbol{\theta}^{T} \Big) \big] \bigg) \\ &\quad + \frac{3L_{f} + 2L_{f}^{2} \alpha}{\alpha} \frac{1}{T+1} \mathbb{E} \bigg[\sum_{t=0}^{T} \left\| \bar{\boldsymbol{\theta}}^{t+1} - \bar{\boldsymbol{\theta}}^{t} \right\|^{2} \bigg] + \frac{4\sigma^{2}}{M} \\ &\quad + \frac{2L_{\bar{\boldsymbol{\theta}}z}^{2}}{M} \bigg[(1 - \alpha\mu)D^{2} + \alpha^{2}B^{2} \bigg] \end{split}$$

$$\leq \frac{1}{T+1} \left\{ \frac{2}{\alpha} + \frac{6L_f + 4L_f^2 \alpha}{\left[2 - \alpha \left(L_f + \beta\right)\right]} \right\} \left(F\left(\bar{\theta}^0\right) - \mathbb{E}\left[F\left(\bar{\theta}^T\right)\right] \right) \\
+ \frac{2L_{\bar{\theta}z}^2}{M} \left\{ 1 + \frac{3L_f + 2L_f^2 \alpha}{2\left(2 - \alpha \left(L_f + \beta\right)\right)} \right\} \left[(1 - \alpha \mu)D^2 + \alpha^2 B^2 \right] \\
+ \frac{4\sigma^2}{M}.$$
(77)

Consider now replacing $F(\bar{\theta}^0) - F(\bar{\theta}^T)$ with $\Delta_F = F(\bar{\theta}^0) - \inf_{\bar{\theta}} F(\bar{\theta})$, and note that $g^t(\bar{\theta}^t) - \nabla f(\bar{\theta}^{t+1}) + \frac{1}{\alpha_t}(\bar{\theta}^{t+1} - \bar{\theta}^t) \in \partial F(\bar{\theta}^{t+1})$, where ∂F denotes the set of subgradients of F. It then becomes clear that

$$\mathbb{E}\left[\operatorname{dist}(0, \partial F)^{2}\right] \leq \frac{1}{T+1} \mathbb{E}\left[\sum_{t=0}^{T} \left\|g^{t}\left(\bar{\theta}^{t}\right) - \nabla f\left(\bar{\theta}^{t+1}\right) + \frac{1}{\alpha_{t}}\left(\bar{\theta}^{t+1} - \bar{\theta}^{t}\right)\right\|^{2}\right] \\ \leq \frac{\zeta}{T+1} \Delta_{F} + \frac{2L_{\bar{\theta}z}^{2} \nu}{N} \left[(1-\alpha\mu)D^{2} + \alpha^{2}B^{2}\right] + \frac{4\sigma^{2}}{M}$$

where $\zeta = \frac{2}{\alpha} + \frac{6L_f + 4L_f^2\alpha}{(2-\alpha(L_f + \beta))}$ and $\nu = 1 + \frac{3L_f + 2L_f^2\alpha}{2(2-\alpha(L_f + \beta))}$, which concludes the proof.

D. Bounded gradient norm assumption.

Regarding a bounded gradient norm assumption, this indeed a consequence of the other assumption already made. Specifically, the condition $\mathbb{E}[\|\nabla_{\overline{\theta}}\psi^*(\theta,\zeta_m^*;z_m) - \nabla f(\theta)\|^2] \le \sigma^2$ to hold, requires that the gradient of the loss be bounded, that is $\|\nabla_{\theta}\ell(\theta,z)\|_2 \le B_1 \ \forall \theta,z$, where B_1 is a constant dependent on σ . A similar argument is presented in [34]). To formalize this, recall the definition

$$g^{t}\left(\overline{\theta}^{t}\right) := \frac{1}{M} \sum_{m=1}^{M} \nabla_{\overline{\theta}} \psi\left(\overline{\theta}, \zeta_{m}^{t}; z_{m}^{t}\right) \bigg|_{\overline{\theta} = \overline{\theta}^{t}}.$$
 (78)

Using this, the gradient norm can be bounded as

$$\|g^{t}(\overline{\theta}^{t})\| = \left\| \frac{1}{M} \sum_{m=1}^{M} \nabla_{\overline{\theta}} \psi(\overline{\theta}, \zeta_{m}^{t}; z_{m}^{t}) \right\|$$

$$\stackrel{(a)}{\leq} \frac{1}{M} \sum_{m=1}^{M} \left\| \nabla_{\overline{\theta}} \psi(\overline{\theta}, \zeta_{m}^{t}; z_{m}^{t}) \right\|$$

$$\stackrel{(b)}{\leq} \frac{1}{M} \sum_{m=1}^{M} \max_{m} \left\| \nabla_{\overline{\theta}} \psi(\overline{\theta}, \zeta_{m}^{t}; z_{m}^{t}) \right\|$$

$$= \max_{m} \left\| \nabla_{\overline{\theta}} \psi(\overline{\theta}, \zeta_{m}^{t}; z_{m}^{t}) \right\|$$

where (a) is due to norm properties, and (b) is simply since each term is replaced by the maximum value it can take. For brevity, with a slight abuse of notation let us define

$$\left\| \nabla_{\overline{\theta}} \psi^{\star} \left(\overline{\theta}, \zeta; z \right) \right\| := \max_{m} \left\| \nabla_{\overline{\theta}} \psi \left(\overline{\theta}, \zeta_{m}^{t}; z_{m}^{t} \right) \right\|$$
(80)

Upon invoking the definition of $\psi(\overline{\theta}, \zeta; z) := \ell(\theta; \zeta) + \gamma(\rho - c(z, \zeta))$ (c.f., (10)), it holds that

$$\left\| \nabla_{\overline{\theta}} \psi^{\star} \left(\overline{\theta}, \zeta; z \right) \right\|^{2} \coloneqq \left\| \left[\begin{array}{c} \nabla_{\theta} \ell(\theta; \zeta) \\ \nabla_{\gamma} (\gamma(\rho - c(z, \zeta))) \end{array} \right] \right\|^{2}$$

Testing Method	FGSM	TPGD	CW	PGD	DIFGSM
Hyperparameters	$\epsilon=0.03$	$\begin{aligned} \epsilon &= 0.001 \\ \alpha &= 0.0007 \\ \text{steps} &= 5 \end{aligned}$	$c = 1$ $\kappa = 0.01$ $steps = 50$ $lr = 0.01$	$\epsilon = 0.03$ $\alpha = 0.008$ steps = 10 random start = True	$\epsilon = 0.03$ $\alpha = 0.007$ steps = 10 decay = 0.0
Testing Method	TIFGSM	NIFGSM	RFGSM	VNIFGSM	APGD
Hyperparameters	$\epsilon=0.03$ $\alpha=0.007$ $\mathrm{nsig}=3$ $\mathrm{steps}=10$ $\mathrm{decay}=0.0$ $\mathrm{len\ kernel}=15$ $\mathrm{resize\ rate}=0.9$ $\mathrm{kernel}=\mathrm{gaussian}$ $\mathrm{diversity\ prob}=0.5$ $\mathrm{random\ start}=\mathrm{False}$	$\begin{array}{c} \epsilon = 0.03 \\ \alpha = 0.007 \\ \text{steps} = 10 \\ \text{decay} = 0.0 \\ \text{resize rate} = 0.9 \\ \text{diversity prob} = 0.5 \\ \text{random start} = \text{False} \end{array}$	$\begin{aligned} \epsilon &= 0.03 \\ \alpha &= 0.007 \\ \text{steps} &= 10 \end{aligned}$	$\epsilon = 0.06$ $\alpha = 0.007$ steps = 10 decay = 1.0 $N = 5, \beta = 1.5$	$\epsilon = 0.03$ $\rho = 0.75$ $\text{steps} = 10$ $\text{norm} = \ell_{\infty}$

TABLE III
HYPERPARAMETERS USED FOR TESTING ON CIFAR-10

TABLE IV
HYPER-PARAMETERS USED FOR TESTING ON CIFAR-100

Testing Method	FGSM	TPGD	CW	PGD	DIFGSM
Hyperparameters	$\epsilon = 0.03$	$\begin{array}{c} \epsilon = 0.001 \\ \alpha = 0.005 \\ \mathrm{steps} = 2 \end{array}$	$c = 1$ $\kappa = 0.01$ $steps = 5$ $lr = 0.01$	$\begin{aligned} \epsilon &= 0.001 \\ \alpha &= 0.05 \\ \text{steps} &= 10 \\ \text{random start} &= \text{True} \end{aligned}$	$\begin{aligned} \epsilon &= 0.001 \\ \alpha &= 0.01 \\ \text{steps} &= 5 \\ \text{decay} &= 0.001 \end{aligned}$
Testing Method	TIFGSM	NIFGSM	RFGSM	VNIFGSM	APGD
Hyperparameters	$\epsilon=0.001$ $\alpha=0.01$ $\mathrm{nsig}=3$ $\mathrm{steps}=10$ $\mathrm{decay}=0.0$ $\mathrm{len\ kernel}=15$ $\mathrm{resize\ rate}=0.9$ $\mathrm{kernel}=\mathrm{gaussian}$ $\mathrm{diversity\ prob}=0.5$ $\mathrm{random\ start}=\mathrm{False}$	$\begin{array}{c} \epsilon = 0.03 \\ \alpha = 0.007 \\ \text{steps} = 10 \\ \text{decay} = 0.0 \\ \text{resize rate} = 0.9 \\ \text{diversity prob} = 0.5 \\ \text{random start} = \text{False} \end{array}$	$\begin{aligned} \epsilon &= 0.001 \\ \alpha &= 0.01 \\ \text{steps} &= 5 \end{aligned}$	$\epsilon=0.001$ $\alpha=0.01$ steps = 10 decay = 1.0 $N=5, \beta=1.5$	$\epsilon = 0.01$ $\rho = 0.8$ steps = 10 norm = ℓ_{∞}

$$\leq \|\nabla_{\theta} \ell(\theta; \zeta)\|^{2} + \|\nabla_{\gamma} (\gamma(\rho - c(z, \zeta)))\|^{2}$$

$$= \|\nabla_{\theta} \ell(\theta; \zeta)\|^{2} + \|\rho - c(z, \zeta)\|^{2}$$

$$\leq B_{1}^{2} + B_{2}^{2}$$

$$= B^{2}$$
(81)

where $\|\nabla_{\theta}\ell(\theta;\zeta)\|^2 \leq B_1^2$ follows from discussion above, and $\|\rho-c(z,\zeta)\|^2 \leq B_2^2$ is a direct consequence of boundedness constraint under the Assumption 1. Hence it is clear that $\mathbb{E}\|g^t(\overline{\theta}^t)\|^2 \leq \mathbb{E}\|\nabla_{\overline{\theta}}\psi^\star(\overline{\theta},\zeta;z)\|^2 \leq B^2$.

E. Hyperparameters for CIFAR-10 and CIFAR-100 Tests

The numerical test results presented in Table I and Table II correspond to testing against adversarial samples generated using the trained model parameters and specific algorithms with preselected hyperparameters. For reference, Table III details the hyperparameters utilized for the CIFAR-10 classification problem. To address the observed drop in performance in the 100-class classification scenario, we modified the hyperparameters, which are outlined in Table IV. These adjustments were necessary to ensure consistent and reliable testing across both classification tasks.

REFERENCES

- B. Li et al., "Trustworthy AI: From principles to practices," ACM Comput. Surveys, vol. 55, no. 9, pp. 1–46, 2023.
- [2] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 50–60, May 2020.
- [3] W. Zhang, G. Wang, J. Sun, Y. Yuan, and G. Huang, "STORM: Efficient stochastic transformer based world models for reinforcement learning," in *Proc. Conf. Neural Inf. Process. Syst.*, 2023, pp. 1–20.
- [4] A. Abusnaina et al., "Adversarial example detection using latent neighborhood graph," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 87–96.
- [5] Y.-C. Lin and F. Yu, "DeepSHAP summary for adversarial example detection," in *Proc. IEEE/ACM Int. Workshop Deep Learn. Test.*, 2023, pp. 17–24.
- [6] J. Cohen, E. Rosenfeld, and Z. Kolter, "Certified adversarial robustness via randomized smoothing," in *Proc. Int. Conf. Mach. Learn.*, Jun. 2019, pp. 310–320.
- [7] C. Guo, M. Rana, M. Cisse, and L. Van Der Maaten, "Countering adversarial images using input transformations," in *Proc. Int. Conf. Learn. Rep.*, Apr. 2018, pp. 1–8.
- [8] T. Pang et al., "Two coupled rejection metrics can tell adversarial examples apart," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2022, pp. 23–33.
- [9] Z. Deng, X. Yang, S. Xu, H. Su, and J. Zhu, "LiBRe: A practical Bayesian approach to adversarial detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 72–82.
- [10] F. Sheikholeslami, S. Jain, and G. B. Giannakis, "Minimum uncertainty based detection of adversaries in deep neural networks," in *Proc. Inf. Theory Appl. Workshop*, 2020, pp. 1–16.

- [11] A. Liu, X. Liu, H. Yu, C. Zhang, Q. Liu, and D. Tao, "Training robust deep neural networks via adversarial noise propagation," *IEEE Trans. Image Process.*, vol. 30, pp. 69–81, 2021.
- [12] Y. Wang, X. Ma, J. Bailey, J. Yi, B. Zhou, and Q. Gu, "On the convergence and robustness of adversarial training," in *Proc. Mach. Learn. Res.*, Jun. 2021, pp. 6586–6595.
- [13] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in *Proc. Int. Conf. Learn. Rep.*, Dec. 2015, pp. 1–9.
- [14] C. Zhang, P. Benz, C. Lin, A. Karjauv, J. Wu, and I. S. Kweon, "A survey on universal adversarial attack," Jan. 2022, arXiv:2103.01498.
- [15] X. Yang, J. Lin, H. Zhang, X. Yang, and P. Zhao, "Improving the transferability of adversarial examples via direction tuning," Aug. 2023, arXiv:2303.15109.
- [16] G. Wang, H. Yan, Y. Guo, and X. Wei, "Improving adversarial transferability with gradient refining," Apr. 2022, arXiv:2105.04834.
- [17] E. Dobriban, H. Hassani, D. Hong, and A. Robey, "Provable tradeoffs in adversarially robust classification," *IEEE Trans. Inf. Theory*, vol. 69, no. 12, pp. 7793–7822, Dec. 2023.
- [18] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated optimization in heterogeneous networks," in *Proc. Mach. Learn. Syst.*, vol. 2, 2020, pp. 429–450.
- [19] R. Song, L. Zhou, L. Lyu, A. Festag, and A. Knoll, "ResFed: Communication efficient federated learning with deep compressed residuals," *IEEE Internet Things J.*, vol. 11, no. 6, pp. 9458–9472, Mar. 2024.
- [20] X. Liu and T. Ratnarajah, "Computation and communication efficient federated learning over wireless networks," Sep. 2023, arXiv:2309.01816.
- [21] B. Weng, J. Sun, G. Huang, F. Deng, G. Wang, and J. Chen, "Competitive meta-learning," *IEEE/CAA J. Automatica Sinica*, vol. 10, no. 9, pp. 1902–1904, Sep. 2023.
- [22] K. Pillutla, S. M. Kakade, and Z. Harchaoui, "Robust aggregation for federated learning," *IEEE Trans. Signal Process.*, vol. 70, pp. 42–54, 2022
- [23] D. Data and S. N. Diggavi, "Byzantine-resilient high-dimensional federated learning," *IEEE Trans. Inf. Theory*, vol. 69, no. 10, pp. 6639–6670, Oct. 2023.
- [24] Y. Deng, M. M. Kamani, and M. Mahdavi, "Distributionally robust federated averaging," in *Proc. NeurIPS*, vol. 33, 2020, pp. 15111–15122.
- [25] K. Pillutla, Y. Laguel, J. Malick, and Z. Harchaoui, "Federated learning with superquantile aggregation for heterogeneous data," *Mach. Learn.*, pp. 1–68, May 2023.
- [26] A. Reisizadeh, F. Farnia, R. Pedarsani, and A. Jadbabaie, "Robust federated learning: The case of affine distribution shifts," in *Proc. Adv. Neural Inf. Process. Syst.*, 2020, pp. 1–6.
- [27] F. Farnia, A. Reisizadeh, R. Pedarsani, and A. Jadbabaie, "An optimal transport approach to personalized federated learning," *IEEE J. Sel. Areas Inf. Theory*, vol. 3, no. 2, pp. 162–171, Jun. 2022.
- [28] M. S. Pydi and V. Jog, "Adversarial risk via optimal transport and optimal couplings," in *Proc. Int. Conf. Mach. Learn*, 2020, pp. 7814–7823.
- [29] M. S. Pydi and V. Jog, "The many faces of adversarial risk: An expanded study," *IEEE Trans. Inf. Theory*, vol. 70, no. 1, pp. 550–570, Jan. 2024.
- [30] E. Delage and Y. Ye, "Distributionally robust optimization under moment uncertainty with application to data-driven problems," *Oper. Res.*, vol. 58, no. 3, pp. 595-612, 2010.
- [31] W. Wiesemann, D. Kuhn, and M. Sim, "Distributionally robust convex optimization," Oper. Res., vol. 62, no. 6, pp. 1358–1376, 2014.
- [32] Z. Hu and L. J. Hong, "Kullback-Leibler divergence constrained distributionally robust optimization," Optimization, to be published.
- [33] C. Bandi and D. Bertsimas, "Robust option pricing," Eur. J. Oper. Res., vol. 239, no. 3, pp. 842–853, 2014.
- [34] A. Sinha, H. Namkoong, and J. Duchi, "Certifying some distributional robustness with principled adversarial training," in *Proc. Int. Conf. Learn. Rep.*, Apr. 2018, pp. 1–34.
- [35] J. Blanchet and K. Murthy, "Quantifying distributional model risk via optimal transport," Math. Oper. Res., vol. 44, no. 2, pp. 565-600, 2019.
- [36] C. Villani, Optimal Transport: Old and New, vol. 338. New York, NY, USA: Springer, 2008.
- [37] T. Lin, C. Jin, and M. Jordan, "On gradient descent ascent for nonconvex-concave minimax problems," in *Proc. Int. Conf. Mach. Learn.*, May 2020, pp. 6083–6093.
- [38] D. P. Bertsekas, "Nonlinear programming," J. Oper. Res. Soc., vol. 48, no. 3, pp. 334–334, 1997.
- [39] R. T. Rockafellar and R. J.-B. Wets, Variational Analysis, vol. 317. New York, NY, USA: Springer, 2009.

- [40] A. Kurakin, I. J. Goodfellow, and S. Bengio, "Adversarial machine learning at scale," in *Proc. Int. Conf. Learn. Rep.*, Apr. 2017, pp. 1–9.
- [41] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," in *Proc. Int. Conf. Learn. Rep.*, Apr. 2018, pp. 1–9.
- [42] G. Wang, G. B. Giannakis, and J. Chen, "Learning ReLU networks on linearly separable data: Algorithm, optimality, and generalization," *IEEE Trans. Signal Process.*, vol. 67, no. 9, pp. 2357–2370, May 2019.
- [43] D. P. Kingma and J. Ba, "ADAM: A method for stochastic optimization," in *Proc. Int. Conf. Learn. Rep.*, May 2015, pp. 1–8.
- [44] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [45] S. Zagoruyko and N. Komodakis, "Wide residual networks," in Proc. Brit. Mach. Vis. Conf., Jun. 2016, pp. 1-9.
- [46] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. Intl. Conf. Learn. Rep.*, Apr. 2014, pp. 1–6.
- [47] O. Russakovsky et al., "ImageNet large scale visual recognition challenge," Int. J. Comput. Vis., vol. 115, pp. 211–252, Apr. 2015.
- [48] X. Wang and K. He, "Enhancing the transferability of adversarial attacks through variance tuning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 24–33.
- [49] F. Tramèr, A. Kurakin, N. Papernot, I. J. Goodfellow, D. Boneh, and P. McDaniel, "Ensemble adversarial training: Attacks and defenses," in Proc. Int. Conf. Learn. Rep., Apr. 2018, pp. 1–9.
- [50] Y. Dong, T. Pang, H. Su, and J. Zhu, "Evading defenses to transferable adversarial examples by translation-invariant attacks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 12–21.
- [51] C. Xie et al., "Improving transferability of adversarial examples with input diversity," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, May 2019, pp. 30–39.
- [52] J. Lin, C. Song, K. He, L. Wang, and J. E. Hopcroft, "Nesterov accelerated gradient and scale invariance for adversarial attacks," in *Proc. Int. Conf. Learn. Rep.*, May 2019, pp. 1–6.
- [53] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in Proc. IEEE Symp. Security Privacy, 2017, pp. 39–57.
- [54] H. Zhang, Y. Yu, J. Jiao, E. Xing, L. E. Ghaoui, and M. Jordan, "Theoretically principled trade-off between robustness and accuracy," in Proc. Int. Conf. Mach. Learn., May 2019, pp. 72–82.
- [55] F. Croce and M. Hein, "Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks," in *Proc. Int. Conf. Mach. Learn.*, May 2020, pp. 2206–2216.
- [56] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. Y. Arcas, "Communication-efficientlearning of deep networks from decentralized data," in *Proc. Int. Conf. Artif. Intell. Stat.*, vol. 54, Apr. 2017, pp. 1273–1282.
- [57] J. M. Danskin, "The theory of max-min, with applications," SIAM J. Appl. Math., vol. 14, no. 4, pp. 641-664, Jul. 1966.
- [58] Y. Xu, R. Jin, and T. Yang, "Non-asymptotic analysis of stochastic methods for non-smooth non-convex regularized problems," in *Proc.* Adv. Neural Inf. Process. Syst., Dec. 2019, pp. 2626–2636.



Alireza Sadeghi (Member, IEEE) received the B.Sc. degree (Hons.) in electrical engineering from the Iran University of Science and Technology, Tehran, Iran, in 2012, the M.Sc. degree in electrical engineering from the University of Tehran in 2015, and the Ph.D. degree in electrical engineering from the University of Minnesota, Minneapolis, MN, USA, in 2021.

He is currently a Postdoctoral Associate with the Department of Electrical and Computer Engineering, University of Minnesota. His research interests

include machine learning, optimization, and signal processing with applications to data science, networking, and cyber-physical systems. He was a recipient of ADC Fellowship and DDF from the University of Minnesota, and the Student Travel Awards from the IEEE Communications Society and the National Science Foundation.



Gang Wang (Senior Member, IEEE) received the B.Eng. degree in automatic control and the Ph.D. degree in control science and engineering from the Beijing Institute of Technology, Beijing, China, in 2011 and 2018, respectively, and the Ph.D. degree in electrical and computer engineering from the University of Minnesota, Minneapolis, MN, USA, in 2018.

Since July 2018, he has been a Postdoctoral Researcher with the University of Minnesota. Since August 2020, he has been a Professor with the

School of Automation, Beijing Institute of Technology. His research interests include signal processing, control, and reinforcement learning with applications to cyber-physical systems and multiagent systems.

Dr. Wang was a recipient of the Best Paper Award from the FITEE in 2021, the Excellent Doctoral Dissertation Award from the Chinese Association of Automation in 2019, the Eusipco Best Student Paper Award in 2017, and the Best Conference Paper at the 2019 IEEE Power and Energy Society General Meeting. He is currently an Associate Editor of Signal Processing and the IEEE TRANSACTIONS ON SIGNAL AND INFORMATION PROCESSING OVER NETWORKS, and an Early Career Advisory Board Member of the IEEE/CAA JOURNAL OF AUTOMATICA SINICA.



Georgios B. Giannakis (Fellow, IEEE) received the Diploma degree in electrical engineering from the National Technical University of Athens, Greece, in 1981, and the first M.Sc. degree in electrical engineering, the second M.Sc. degree in mathematics, and the Ph.D. degree in electrical engineering from the University of the Southern California in 1983, 1986, and 1986, respectively.

He was with the University of Virginia from 1987 to 1998, and since 1999, he has been with the University of Minnesota (UMN), where he held an

Endowed Chair of Telecommunications, served as the Director of the Digital Technology Center from 2008 to 2021, and since 2016, he has been a UMN Presidential Chair in ECE. He is the (co-) inventor of 36 issued patents. His interests span the areas of statistical learning, communications, and networking subjects on which he has published more than 490 journal papers, 800 conference papers, 26 book chapters, two edited books and two research monographs. His current research focuses on data science with applications to IoT and power networks with renewables. He is the (co-)recipient of ten Best Journal Paper Awards from the IEEE Signal Processing and Communications Societies, including the G. Marconi Prize. He also received the IEEE-SPS 'Nobert Wiener' Society Award in 2019; the EURASIP's 'A. Papoulis' Society Award in 2020; the Tech. Achievement Awards from the IEEE-SPS in 2000 and from EURASIP in 2005; the IEEE ComSoc Education Award in 2019; and the IEEE Fourier Technical Field Award in 2015. He is a member of the Academia Europaea, the Academy of Athens, Greece, the Royal Academy of Engineering U.K., and a Fellow of the U.S. National Academy of Inventors, the European Academy of Sciences, and EURASIP. He has served the IEEE in a number of posts, including that of a Distinguished Lecturer for the IEEE-SPS.