
Distributed Distributionally Robust Optimization with Non-Convex Objectives

Yang Jiao
Tongji University
yangjiao@tongji.edu.cn

Kai Yang*
Tongji University
kaiyang@tongji.edu.cn

Dongjin Song
University of Connecticut
dongjin.song@uconn.edu

Abstract

Distributionally Robust Optimization (DRO), which aims to find an optimal decision that minimizes the worst case cost over the ambiguity set of probability distribution, has been widely applied in diverse applications, *e.g.*, network behavior analysis, risk management, *etc.* However, existing DRO techniques face three key challenges: 1) how to deal with the asynchronous updating in a distributed environment; 2) how to leverage the prior distribution effectively; 3) how to properly adjust the degree of robustness according to different scenarios. To this end, we propose an asynchronous distributed algorithm, named **Asynchronous Single-loop Alternative Gradient Projection (ASPIRE)** algorithm with the **Iterative Active Set (EASE)** method to tackle the distributed distributionally robust optimization (DDRO) problem. Furthermore, a new uncertainty set, *i.e.*, constrained D -norm uncertainty set, is developed to effectively leverage the prior distribution and flexibly control the degree of robustness. Finally, our theoretical analysis elucidates that the proposed algorithm is guaranteed to converge and the iteration complexity is also analyzed. Extensive empirical studies on real-world datasets demonstrate that the proposed method can not only achieve fast convergence, and remain robust against data heterogeneity as well as malicious attacks, but also tradeoff robustness with performance.

1 Introduction

The past decade has witnessed the proliferation of smartphones and Internet of Things (IoT) devices, which generate a plethora of data everyday. Centralized machine learning requires gathering the data to a particular server to train models which incurs high communication overhead [46] and suffers privacy risks [43]. As a remedy, distributed machine learning methods have been proposed. Considering a distributed system composed of N workers (devices), we denote the dataset of these workers as $\{D_1, \dots, D_N\}$. For the j^{th} ($1 \leq j \leq N$) worker, the labeled dataset is given as $D_j = \{\mathbf{x}_j^i, y_j^i\}$, where $\mathbf{x}_j^i \in \mathbb{R}^d$ and $y_j^i \in \{1, \dots, c\}$ denote the i^{th} data sample and the corresponding label, respectively. The distributed learning tasks can be formulated as the following optimization problem,

$$\min_{\mathbf{w} \in \mathcal{W}} F(\mathbf{w}) \quad \text{with} \quad F(\mathbf{w}) := \sum_j f_j(\mathbf{w}), \quad (1)$$

where $\mathbf{w} \in \mathbb{R}^p$ is the model parameter to be learned and $\mathcal{W} \subseteq \mathbb{R}^p$ is a nonempty closed convex set, $f_j(\cdot)$ is the empirical risk over the j^{th} worker involving only the local data:

$$f_j(\mathbf{w}) = \sum_{i: \mathbf{x}_j^i \in D_j} \frac{1}{|D_j|} \mathcal{L}_j(\mathbf{x}_j^i, y_j^i; \mathbf{w}), \quad (2)$$

*Corresponding author.

where \mathcal{L}_j is the local objective function over the j^{th} worker. Problem in Eq. (1) arises in numerous areas, such as distributed signal processing [19], multi-agent optimization [36], *etc.* However, such problem does not consider the data heterogeneity [57, 40, 39, 30] among different workers (*i.e.*, data distribution of workers could be substantially different from each other [44]). Indeed, it has been shown that traditional federated approaches, such as FedAvg [33], built for independent and identically distributed (IID) data may perform poorly when applied to Non-IID data [27]. This issue can be mitigated via learning a robust model that aims to achieve uniformly good performance over all workers by solving the following distributionally robust optimization (DRO) problem in a distributed manner:

$$\min_{\mathbf{w} \in \mathcal{W}} \max_{\mathbf{p} \in \Omega \subseteq \Delta_N} F(\mathbf{w}, \mathbf{p}) := \sum_j p_j f_j(\mathbf{w}), \quad (3)$$

where $\mathbf{p} = [p_1, \dots, p_N] \in \mathbb{R}^N$ is the adversarial distribution in N workers, the j^{th} entry in this vector, *i.e.*, p_j represents the adversarial distribution value for the j^{th} worker. $\Delta_N = \{\mathbf{p} \in \mathbb{R}_+^N : \mathbf{1}^\top \mathbf{p} = 1\}$ and Ω is a subset of Δ_N . Agnostic federated learning (AFL) [35] firstly introduces the distributionally robust (agnostic) loss in federated learning and provides the convergence rate for (strongly) convex functions. However, AFL does not discuss the setting of Ω . DRFA-Prox [16] considers $\Omega = \Delta_N$ and imposes a regularizer on adversarial distribution to leverage the prior distribution. Nevertheless, three key challenges have not yet been addressed by prior works. First, whether it is possible to construct an uncertainty framework that can not only flexibly maintain the trade-off between the model robustness and performance but also effectively leverage the prior distribution? Second, how to design asynchronous algorithms with guaranteed convergence? Compared to synchronous algorithms, the master in asynchronous algorithms can update its parameters after receiving updates from only a small subset of workers [58, 10]. Asynchronous algorithms are particularly desirable in practice since they can relax strict data dependencies and ensure convergence even in the presence of device failures [58]. Finally, whether it is possible to flexibly adjust the degree of robustness? Moreover, it is necessary to provide convergence guarantee when the objectives (*i.e.*, $f_j(\mathbf{w}_j), \forall j$) are non-convex.

To this end, we propose ASPIRE-EASE to effectively address the aforementioned challenges. Firstly, different from existing works, the prior distribution is incorporated within the constraint in our formulation, which can not only leverage the prior distribution more effectively but also achieve guaranteed feasibility for any adversarial distribution within the uncertainty set. The prior distribution can be obtained from side information or uniform distribution [41], which is necessary to construct the uncertainty (ambiguity) set and obtain a more robust model [16]. Specifically, we formulate the prior distribution informed distributionally robust optimization (PD-DRO) problem as:

$$\begin{aligned} \min_{z \in \mathcal{Z}, \{\mathbf{w}_j \in \mathcal{W}\}} \max_{\mathbf{p} \in \mathcal{P}} \sum_j p_j f_j(\mathbf{w}_j) \\ \text{s.t.} \quad z = \mathbf{w}_j, \quad j = 1, \dots, N, \\ \text{var.} \quad z, \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N, \end{aligned} \quad (4)$$

where $z \in \mathbb{R}^p$ is the global consensus variable, $\mathbf{w}_j \in \mathbb{R}^p$ is the local variable (local model parameter) of j^{th} worker and $\mathcal{Z} \subseteq \mathbb{R}^p$ is a nonempty closed convex set. $\mathcal{P} \subseteq \mathbb{R}_+^N$ is the uncertainty (ambiguity) set of adversarial distribution \mathbf{p} , which is set based on the prior distribution. To solve the PD-DRO problem in an asynchronous distributed manner, we first propose **Asynchronous Single-loop Alternating Gradient Projection (ASPIRE)**, which employs *simple* gradient projection steps for the update of primal and dual variables at every iteration, thus is computationally *efficient*. Next, the **Iterative Active Set (EASE)** method is employed to replace the traditional cutting plane method to improve the computational efficiency and speed up the convergence. We further provide the convergence guarantee for the proposed algorithm. Furthermore, a new uncertainty set, *i.e.*, constrained D -norm (CD -norm), is proposed in this paper and its advantages include: 1) it can flexibly control the degree of robustness; 2) the resulting subproblem is computationally simple; 3) it can effectively leverage the prior distribution and flexibly set the bounds for every p_j .

Contributions. Our contributions can be summarized as follows:

1. We formulate a PD-DRO problem with CD -norm uncertainty set. PD-DRO incorporates the prior distribution as constraints which can leverage prior distribution more effectively and guarantee robustness. In addition, CD -norm is developed to model the ambiguity set around the prior distribution and it provides a flexible way to control the trade-off between model robustness and performance.
2. We develop a *single-loop asynchronous* algorithm, namely ASPIRE-EASE, to optimize PD-DRO in an asynchronous distributed manner. ASPIRE employs simple gradient projection steps to

update the variables at every iteration, which is computationally efficient. And EASE is proposed to replace cutting plane method to enhance the computational efficiency and speed up the convergence. We demonstrate that even if the objectives $f_j(\mathbf{w}_j), \forall j$ are non-convex, the proposed algorithm is guaranteed to converge. We also theoretically derive the iteration complexity of ASPIRE-EASE.

3. Extensive empirical studies on four different real world datasets demonstrate the superior performance of the proposed algorithm. It is seen that ASPIRE-EASE can not only ensure the model's robustness against data heterogeneity but also mitigate malicious attacks.

2 Preliminaries

2.1 Distributionally Robust Optimization

Optimization problems often contain uncertain parameters. A small perturbation of the parameters could render the optimal solution of the original optimization problem infeasible or completely meaningless [5]. Distributionally robust optimization (DRO) [28, 17, 7] assumes that the probability distributions of uncertain parameters are unknown but remain in an ambiguity (uncertainty) set and aims to find a decision that minimizes the worst case expected cost over the ambiguity set, whose general form can be expressed as,

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{P \in \mathbf{P}} \mathbb{E}_P[r(\mathbf{x}, \boldsymbol{\xi})], \quad (5)$$

where $\mathbf{x} \in \mathcal{X}$ represents the decision variable, \mathbf{P} is the ambiguity set of probability distributions P of uncertain parameters $\boldsymbol{\xi}$. Existing methods for solving DRO can be broadly grouped into two widely-used categories [42]: 1) Dual methods [15, 50, 18] reformulate the primal DRO problems as deterministic optimization problems through duality theory. Ben-Tal et al. [2] reformulate the robust linear optimization (RLO) problem with an ellipsoidal uncertainty set as a second-order cone optimization problem (SOCP). 2) Cutting plane methods [34, 6] (also called adversarial approaches [21]) continuously solve an approximate problem with a finite number of constraints of the primal DRO problem, and subsequently check whether new constraints are needed to refine the feasible set. Recently, several new methods [41, 29, 23] have been developed to solve DRO, which need to solve the inner maximization problem at every iteration.

2.2 Cutting Plane Method for PD-DRO

In this section, we introduce the cutting plane method for PD-DRO in Eq. (4). We first reformulate PD-DRO by introducing an additional variable $h \in \mathcal{H}$ ($\mathcal{H} \subseteq \mathbb{R}^1$ is a nonempty closed convex set) and protection function $g(\{\mathbf{w}_j\})$ [55]. Introducing additional variable h is an epigraph reformulation [3, 56]. In this case, Eq. (4) can be reformulated as the form with uncertainty in the constraints:

$$\begin{aligned} & \min_{\mathbf{z} \in \mathcal{Z}, \{\mathbf{w}_j \in \mathcal{W}\}, h \in \mathcal{H}} h \\ \text{s.t. } & \sum_j \bar{p}_j f_j(\mathbf{w}_j) + g(\{\mathbf{w}_j\}) - h \leq 0, \\ & \mathbf{z} = \mathbf{w}_j, j=1, \dots, N, \\ \text{var. } & \mathbf{z}, \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N, h. \end{aligned} \quad (6)$$

where \bar{p} is the nominal value of the adversarial distribution for every worker and $g(\{\mathbf{w}_j\}) = \max_{P \in \mathbf{P}} \sum_j (p_j - \bar{p}) f_j(\mathbf{w}_j)$ is the protection function. Eq. (6) is a semi-infinite program (SIP) which contains infinite constraints and cannot be solved directly [42]. Denoting the set of cutting plane parameters in $(t+1)^{\text{th}}$ iteration as $\mathbf{A}^t \subseteq \mathbb{R}^N$, the following function is used to approximate $g(\{\mathbf{w}_j\})$:

$$\bar{g}(\{\mathbf{w}_j\}) = \max_{\mathbf{a}_l \in \mathbf{A}^t} \mathbf{a}_l^\top \mathbf{f}(\mathbf{w}) = \max_{\mathbf{a}_l \in \mathbf{A}^t} \sum_j a_{l,j} f_j(\mathbf{w}_j), \quad (7)$$

where $\mathbf{a}_l = [a_{l,1}, \dots, a_{l,N}] \in \mathbb{R}^N$ denotes the parameters of l^{th} cutting plane in \mathbf{A}^t and $\mathbf{f}(\mathbf{w}) = [f_1(\mathbf{w}_1), \dots, f_N(\mathbf{w}_N)] \in \mathbb{R}^N$. Substituting the protection function $g(\{\mathbf{w}_j\})$ with $\bar{g}(\{\mathbf{w}_j\})$, we can obtain the following approximate problem:

$$\begin{aligned} & \min_{\mathbf{z} \in \mathcal{Z}, \{\mathbf{w}_j \in \mathcal{W}\}, h \in \mathcal{H}} h \\ \text{s.t. } & \sum_j (\bar{p} + a_{l,j}) f_j(\mathbf{w}_j) - h \leq 0, \forall \mathbf{a}_l \in \mathbf{A}^t, \\ & \mathbf{z} = \mathbf{w}_j, j=1, \dots, N, \\ \text{var. } & \mathbf{z}, \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N, h. \end{aligned} \quad (8)$$

3 ASPIRE

Distributed optimization is an attractive approach for large-scale learning tasks [54, 8] since it does not require data aggregation, which protects data privacy while also reducing bandwidth requirements [45]. When the neural network models (*i.e.*, $f_j(\mathbf{w}_j), \forall j$ are non-convex functions) are used, solving problem in Eq. (8) in a distributed manner facing two challenges. 1) Computing the optimal solution to a non-convex subproblem requires a large number of iterations and therefore is highly computationally intensive if not impossible. Thus, the traditional Alternating Direction Method of Multipliers (ADMM) is ineffective. 2) The communication delays of workers may differ significantly [11], thus, asynchronous algorithms are strongly preferred.

To this end, we propose the **A**synchronous **S**ingle-loop **P**rojection **A**lternative **g**radient **p**rojEction (ASPIRE). The advantages of the proposed algorithm include: 1) ASPIRE uses simple gradient projection steps to update variables in each iteration and therefore it is computationally more efficient than the traditional ADMM method, which seeks to find the optimal solution in non-convex (for $\mathbf{w}_j, \forall j$) and convex (for \mathbf{z} and h) optimization subproblems every iteration, 2) the proposed asynchronous algorithm does not need strict synchronization among different workers. Therefore, ASPIRE remains resilient against communication delays and potential hardware failures from workers. Details of the algorithm are given below. Firstly, we define the node as master which is responsible for updating the global variable \mathbf{z} , and we define the node which is responsible for updating the local variable \mathbf{w}_j as worker j . In each iteration, the master updates its variables once it receives updates from at least S workers, *i.e.*, active workers, satisfying $1 \leq S \leq N$. \mathbf{Q}^{t+1} denotes the index subset of workers from which the master receives updates during $(t+1)^{\text{th}}$ iteration. We also assume the master will receive updated variables from every worker at least once for each τ iterations. The augmented Lagrangian function of Eq. (8) can be written as:

$$L_p = h + \sum_l \lambda_l (\sum_j (\bar{p} + a_{l,j}) f_j(\mathbf{w}_j) - h) + \sum_j \phi_j^\top (\mathbf{z} - \mathbf{w}_j) + \sum_j \frac{\kappa_1}{2} \|\mathbf{z} - \mathbf{w}_j\|^2, \quad (9)$$

where $L_p = L_p(\{\mathbf{w}_j\}, \mathbf{z}, h, \{\lambda_l\}, \{\phi_j\})$, $\lambda_l \in \mathbf{\Lambda}, \forall l$ and $\phi_j \in \mathbf{\Phi}, \forall j$ represent the dual variables of inequality and equality constraints in Eq. (8), respectively. $\mathbf{\Lambda} \subseteq \mathbb{R}^1$ and $\mathbf{\Phi} \subseteq \mathbb{R}^p$ are nonempty closed convex sets, constant $\kappa_1 > 0$ is a penalty parameter. Note that Eq. (9) does not consider the second-order penalty term for inequality constraint since it will invalidate the distributed optimization. Following [52], the regularized version of Eq. (9) is employed to update all variables as follows,

$$\tilde{L}_p(\{\mathbf{w}_j\}, \mathbf{z}, h, \{\lambda_l\}, \{\phi_j\}) = L_p - \sum_l \frac{c_1^t}{2} \|\lambda_l\|^2 - \sum_j \frac{c_2^t}{2} \|\phi_j\|^2, \quad (10)$$

where c_1^t and c_2^t denote the regularization terms in $(t+1)^{\text{th}}$ iteration. To avoid enumerating the whole dataset, the mini-batch loss could be used. A batch of instances with size m can be randomly sampled from each worker during each iteration. The loss function of these instances from j^{th} worker is given by $\hat{f}_j(\mathbf{w}_j) = \sum_{i=1}^m \frac{1}{m} \mathcal{L}_j(\mathbf{x}_j^i, \mathbf{y}_j^i; \mathbf{w}_j)$. It is evident that $\mathbb{E}[\hat{f}_j(\mathbf{w}_j)] = f_j(\mathbf{w}_j)$ and $\mathbb{E}[\nabla \hat{f}_j(\mathbf{w}_j)] = \nabla f_j(\mathbf{w}_j)$. In $(t+1)^{\text{th}}$ master iteration, the proposed algorithm proceeds as follows.

1) Active workers update the local variables \mathbf{w}_j as follows,

$$\mathbf{w}_j^{t+1} = \begin{cases} \mathcal{P}_{\mathcal{W}}(\mathbf{w}_j^t - \eta_w \nabla_{\mathbf{w}_j} \tilde{L}_p(\{\mathbf{w}_j^{\tilde{t}_j}\}, \mathbf{z}^{\tilde{t}_j}, h^{\tilde{t}_j}, \{\lambda_l^{\tilde{t}_j}\}, \{\phi_j^{\tilde{t}_j}\})), \forall j \in \mathbf{Q}^{t+1}, \\ \mathbf{w}_j^t, \forall j \notin \mathbf{Q}^{t+1}, \end{cases} \quad (11)$$

where η_w represents the step-size and $\mathcal{P}_{\mathcal{W}}$ represents the projection onto the closed convex set \mathcal{W} and we set $\mathcal{W} = \{\mathbf{w}_j \mid \|\mathbf{w}_j\|_\infty \leq \alpha_1\}$, α_1 is a positive constant. And \tilde{t}_j is the last iteration during which worker j was active. It is seen that $\mathbf{w}_j^t = \mathbf{w}_j^{\tilde{t}_j}$ and $\phi_j^t = \phi_j^{\tilde{t}_j}, \forall j \in \mathbf{Q}^{t+1}$. Then, the active workers ($j \in \mathbf{Q}^{t+1}$) transmit their local model parameters \mathbf{w}_j^{t+1} and loss $f_j(\mathbf{w}_j)$ to the master.

2) After receiving the updates from active workers, the master updates the global consensus variable \mathbf{z} , additional variable h and dual variables λ_l as follows,

$$\mathbf{z}^{t+1} = \mathcal{P}_{\mathcal{Z}}(\mathbf{z}^t - \eta_z \nabla_{\mathbf{z}} \tilde{L}_p(\{\mathbf{w}_j^{t+1}\}, \mathbf{z}^t, h^t, \{\lambda_l^t\}, \{\phi_j^t\})), \quad (12)$$

$$h^{t+1} = \mathcal{P}_{\mathcal{H}}(h^t - \eta_h \nabla_h \tilde{L}_p(\{\mathbf{w}_j^{t+1}\}, \mathbf{z}^{t+1}, h^t, \{\lambda_l^t\}, \{\phi_j^t\})), \quad (13)$$

$$\lambda_l^{t+1} = \mathcal{P}_\Lambda(\lambda_l^t + \rho_1 \nabla_{\lambda_l} \tilde{L}_p(\{\mathbf{w}_j^{t+1}\}, \mathbf{z}^{t+1}, h^{t+1}, \{\lambda_l^t\}, \{\phi_j^t\})), \quad l=1, \dots, |\mathbf{A}^t|, \quad (14)$$

where η_z , η_h and ρ_1 represent the step-size. \mathcal{P}_Z , $\mathcal{P}_\mathcal{H}$ and \mathcal{P}_Λ respectively represent the projection onto the closed convex sets Z , \mathcal{H} and Λ . We set $Z = \{z \mid \|z\|_\infty \leq \alpha_1\}$, $\mathcal{H} = \{h \mid 0 \leq h \leq \alpha_2\}$ and $\Lambda = \{\lambda_l \mid 0 \leq \lambda_l \leq \alpha_3\}$, where α_2 and α_3 are positive constants. $|\mathbf{A}^t|$ denotes the number of cutting planes. Then, master broadcasts \mathbf{z}^{t+1} , h^{t+1} , $\{\lambda_l^{t+1}\}$ to the active workers.

3) *Active workers* update the local dual variables ϕ_j as follows,

$$\phi_j^{t+1} = \begin{cases} \mathcal{P}_\Phi(\phi_j^t + \rho_2 \nabla_{\phi_j} \tilde{L}_p(\{\mathbf{w}_j^{t+1}\}, \mathbf{z}^{t+1}, h^{t+1}, \{\lambda_l^{t+1}\}, \{\phi_j^t\})), \forall j \in \mathbf{Q}^{t+1}, \\ \phi_j^t, \forall j \notin \mathbf{Q}^{t+1}, \end{cases} \quad (15)$$

where ρ_2 represents the step-size and \mathcal{P}_Φ represents the projection onto the closed convex set Φ and we set $\Phi = \{\phi_j \mid \|\phi_j\|_\infty \leq \alpha_4\}$, α_4 is a positive constant. And master can also obtain $\{\phi_j^{t+1}\}$ according to Eq. (15). It is seen that the projection operation in each step is computationally simple since the closed convex sets have simple structures [4].

4 Iterative Active Set Method

Cutting plane methods may give rise to numerous linear constraints and lots of extra message passing [55]. Moreover, more iterations are required to obtain the ε -stationary point when the size of a set containing cutting planes increases (which corresponds to a larger M), which can be seen in Theorem 1. To improve the computational efficiency and speed up the convergence, we consider removing the inactive cutting planes. The proposed **Iterative Active SEt** method (EASE) can be divided into the two steps: during T_1 iterations, 1) solving the cutting plane generation subproblem to generate cutting plane, and 2) removing the inactive cutting plane every k iterations, where $k > 0$ is a pre-set constant and can be controlled flexibly.

The cutting planes are generated according to the uncertainty set. For example, if we employ ellipsoid uncertainty set, the cutting plane is generated via solving a SOCP. In this paper, we propose CD -norm uncertainty set, which can be expressed as follows,

$$\mathcal{P} = \{\mathbf{p} : -\tilde{p}_j \leq p_j - q_j \leq \tilde{p}_j, \sum_j \left| \frac{p_j - q_j}{\tilde{p}_j} \right| \leq \Gamma, \mathbf{1}^\top \mathbf{p} = 1\}, \quad (16)$$

where $\Gamma \in \mathbb{R}^1$ can flexibly control the level of robustness, $\mathbf{q} = [q_1, \dots, q_N] \in \mathbb{R}^N$ represents the prior distribution, $-\tilde{p}_j$ and \tilde{p}_j ($\tilde{p}_j \geq 0$) represent the lower and upper bounds for $p_j - q_j$, respectively. The setting of \mathbf{q} and $\tilde{p}_j, \forall j$ are based on the prior knowledge. D -norm is a classical uncertainty set (which is also called as budget uncertainty set) [5]. We call Eq. (16) CD -norm uncertainty set since \mathbf{p} is a probability vector so all the entries of this vector are non-negative and add up to exactly one, *i.e.*, $\mathbf{1}^\top \mathbf{p} = 1$. Due to the special structure of CD -norm, the cutting plane generation subproblem is easy to solve and the level of robustness in terms of the outage probability, *i.e.*, probabilistic bounds of the violations of constraints can be flexibly adjusted via a single parameter Γ . We claim that l_1 -norm (or twice total variation distance) uncertainty set is closely related to CD -norm uncertainty set. Nevertheless, there are two differences: 1) CD -norm uncertainty set could be regarded as a weighted l_1 -norm with additional constraints. 2) CD -norm uncertainty set can flexibly set the lower and upper bounds for every p_j (*i.e.*, $q_j - \tilde{p}_j \leq p_j \leq p_j + \tilde{p}_j$), while $0 \leq p_j \leq 1, \forall j$ in l_1 -norm uncertainty set. Based on the CD -norm uncertainty set, the cutting plane can be derived as follows,

1) Solve the following problem,

$$\begin{aligned} \mathbf{p}^{t+1} &= \arg \max_{p_1, \dots, p_N} \sum_j (p_j - \bar{p}) f_j(\mathbf{w}_j) \\ \text{s.t.} \quad &\sum_j \left| \frac{p_j - q_j}{\tilde{p}_j} \right| \leq \Gamma, \quad -\tilde{p}_j \leq p_j - q_j \leq \tilde{p}_j, \forall j, \quad \sum_j p_j = 1 \\ &\text{var.} \quad p_1, \dots, p_N \end{aligned} \quad (17)$$

where $\mathbf{p}^{t+1} = [p_1^{t+1}, \dots, p_N^{t+1}] \in \mathbb{R}^N$. Let $\tilde{\mathbf{a}}^{t+1} = \mathbf{p}^{t+1} - \bar{\mathbf{p}}$, where $\bar{\mathbf{p}} = [\bar{p}, \dots, \bar{p}] \in \mathbb{R}^N$. This first step aims to obtain the distribution $\tilde{\mathbf{a}}^{t+1}$ by solving problem in Eq. (17). This problem can be effectively solved through combining merge sort [13] (for sorting $\tilde{p}_j f_j(\mathbf{w}_j), j=1, \dots, N$) with few basic arithmetic operations (for obtaining $p_j^{t+1}, j=1, \dots, N$). Since N is relatively large in

Algorithm 1 ASPIRE-EASE

Initialization: iteration $t = 0$, variables $\{\mathbf{w}_j^0\}, \mathbf{z}^0, h^0, \{\lambda_l^0\}, \{\phi_j^0\}$ and set \mathbf{A}^0 .

repeat

for active worker do

 updates local \mathbf{w}_j^{t+1} according to Eq. (11);

end for

active workers transmit local model parameters and loss to *master*;

master receives updates from *active workers* **do**

 updates $\mathbf{z}^{t+1}, h^{t+1}, \{\lambda_l^{t+1}\}, \{\phi_j^{t+1}\}$ in master according to Eq. (12), (13), (14), (15);

master broadcasts $\mathbf{z}^{t+1}, h^{t+1}, \{\lambda_l^{t+1}\}$ to *active workers*;

for active worker do

 updates local ϕ_j^{t+1} according to Eq. (15);

end for

if $(t + 1) \bmod k == 0$ and $t < T_1$ **then**

master updates \mathbf{A}^{t+1} according to Eq. (19) and (20), and broadcast parameters to all workers;

end if

$t = t + 1$;

until convergence

distributed system, the arithmetic complexity of solving problem in Eq. (17) is dominated by merge sort, which can be regarded as $\mathcal{O}(N \log(N))$.

2) Let $\mathbf{f}(\mathbf{w}) = [f_1(\mathbf{w}_1), \dots, f_N(\mathbf{w}_N)] \in \mathbb{R}^N$, check the feasibility of the following constraints:

$$\tilde{\mathbf{a}}^{t+1 \top} \mathbf{f}(\mathbf{w}) \leq \max_{\mathbf{a}_l \in \mathbf{A}^t} \mathbf{a}_l^\top \mathbf{f}(\mathbf{w}), \quad (18)$$

3) If Eq. (18) is violated, $\tilde{\mathbf{a}}^{t+1}$ will be added into \mathbf{A}^t :

$$\mathbf{A}^{t+1} = \begin{cases} \mathbf{A}^t \cup \{\tilde{\mathbf{a}}^{t+1}\}, & \text{if Eq.(18) is violated,} \\ \mathbf{A}^t, & \text{otherwise,} \end{cases} \quad (19)$$

when a new cutting plane is added, its corresponding dual variable $\lambda_{|\mathbf{A}^t|+1}^{t+1} = 0$ will be generated. After the cutting plane subproblem is solved, the inactive cutting plane will be removed, that is:

$$\mathbf{A}^{t+1} = \begin{cases} \mathbb{C}_{\mathbf{A}^{t+1}}\{\mathbf{a}_l\}, & \text{if } \lambda_l^{t+1} = 0 \text{ and } \lambda_l^t = 0, 1 \leq l \leq |\mathbf{A}^t|, \\ \mathbf{A}^{t+1}, & \text{otherwise,} \end{cases} \quad (20)$$

where $\mathbb{C}_{\mathbf{A}^{t+1}}\{\mathbf{a}_l\}$ is the complement of $\{\mathbf{a}_l\}$ in \mathbf{A}^{t+1} , and the dual variable will be removed. Then master broadcasts $\mathbf{A}^{t+1}, \{\lambda_l^{t+1}\}$ to all workers. Details of algorithm are summarized in Algorithm 1.

5 Convergence Analysis

Definition 1 (Stationarity gap) *Following [52, 32, 53], the stationarity gap of our problem at t^{th} iteration is defined as:*

$$\nabla G^t = \begin{bmatrix} \left\{ \frac{1}{\eta_w} (\mathbf{w}_j^t - \mathcal{P}_{\mathcal{W}}(\mathbf{w}_j^t - \eta_w \nabla_{\mathbf{w}_j} L_p(\{\mathbf{w}_j^t\}, \mathbf{z}^t, h^t, \{\lambda_l^t\}, \{\phi_j^t\}))) \right\} \\ \frac{1}{\eta_z} (\mathbf{z}^t - \mathcal{P}_{\mathcal{Z}}(\mathbf{z}^t - \eta_z \nabla_{\mathbf{z}} L_p(\{\mathbf{w}_j^t\}, \mathbf{z}^t, h^t, \{\lambda_l^t\}, \{\phi_j^t\}))) \\ \frac{1}{\eta_h} (h^t - \mathcal{P}_{\mathcal{H}}(h^t - \eta_h \nabla_h L_p(\{\mathbf{w}_j^t\}, \mathbf{z}^t, h^t, \{\lambda_l^t\}, \{\phi_j^t\}))) \\ \left\{ \frac{1}{\rho_1} (\lambda_l^t - \mathcal{P}_{\Lambda}(\lambda_l^t + \rho_1 \nabla_{\lambda_l} L_p(\{\mathbf{w}_j^t\}, \mathbf{z}^t, h^t, \{\lambda_l^t\}, \{\phi_j^t\}))) \right\} \\ \left\{ \frac{1}{\rho_2} (\phi_j^t - \mathcal{P}_{\Phi}(\phi_j^t + \rho_2 \nabla_{\phi_j} L_p(\{\mathbf{w}_j^t\}, \mathbf{z}^t, h^t, \{\lambda_l^t\}, \{\phi_j^t\}))) \right\} \end{bmatrix}, \quad (21)$$

where ∇G^t is the simplified form of $\nabla G(\{\mathbf{w}_j^t\}, \mathbf{z}^t, h^t, \{\lambda_l^t\}, \{\phi_j^t\})$.

Definition 2 (ε -stationary point) $(\{\mathbf{w}_j^t\}, \mathbf{z}^t, h^t, \{\lambda_l^t\}, \{\phi_j^t\})$ is an ε -stationary point ($\varepsilon \geq 0$) of a differentiable function L_p , if $\|\nabla G^t\| \leq \varepsilon$. $T(\varepsilon)$ is the first iteration index such that $\|\nabla G^t\| \leq \varepsilon$, i.e., $T(\varepsilon) = \min\{t \mid \|\nabla G^t\| \leq \varepsilon\}$.

Assumption 1 (Smoothness/Gradient Lipschitz) L_p has Lipschitz continuous gradients. We assume that there exists $L > 0$ satisfying

$$\begin{aligned} & \|\nabla_{\theta} L_p(\{\mathbf{w}_j\}, \mathbf{z}, h, \{\lambda_l\}, \{\phi_j\}) - \nabla_{\theta} L_p(\{\hat{\mathbf{w}}_j\}, \hat{\mathbf{z}}, \hat{h}, \{\hat{\lambda}_l\}, \{\hat{\phi}_j\})\| \\ & \leq L\|\mathbf{w}_{\text{cat}} - \hat{\mathbf{w}}_{\text{cat}}; \mathbf{z} - \hat{\mathbf{z}}; h - \hat{h}; \boldsymbol{\lambda}_{\text{cat}} - \hat{\boldsymbol{\lambda}}_{\text{cat}}; \boldsymbol{\phi}_{\text{cat}} - \hat{\boldsymbol{\phi}}_{\text{cat}}\|, \end{aligned}$$

where $\theta \in \{\{\mathbf{w}_j\}, \mathbf{z}, h, \{\lambda_l\}, \{\phi_j\}\}$ and $[\cdot]$ represents the concatenation. $\mathbf{w}_{\text{cat}} - \hat{\mathbf{w}}_{\text{cat}} = [\mathbf{w}_1 - \hat{\mathbf{w}}_1; \dots; \mathbf{w}_N - \hat{\mathbf{w}}_N] \in \mathbb{R}^{pN}$, $\boldsymbol{\lambda}_{\text{cat}} - \hat{\boldsymbol{\lambda}}_{\text{cat}} = [\lambda_1 - \hat{\lambda}_1; \dots; \lambda_{|\mathbf{A}^t|} - \hat{\lambda}_{|\mathbf{A}^t|}] \in \mathbb{R}^{|\mathbf{A}^t|}$, $\boldsymbol{\phi}_{\text{cat}} - \hat{\boldsymbol{\phi}}_{\text{cat}} = [\phi_1 - \hat{\phi}_1; \dots; \phi_N - \hat{\phi}_N] \in \mathbb{R}^{pN}$.

Assumption 2 (Boundedness) Before obtaining the ε -stationary point (i.e., $t \leq T(\varepsilon) - 1$), we assume variables in master satisfy that $\|\mathbf{z}^{t+1} - \mathbf{z}^t\|^2 + \|h^{t+1} - h^t\|^2 + \sum_l \|\lambda_l^{t+1} - \lambda_l^t\|^2 \geq \vartheta$, where $\vartheta > 0$ is a relative small constant. The change of the variables in master is upper bounded within τ iterations:

$$\|\mathbf{z}^t - \mathbf{z}^{t-k}\|^2 \leq \tau k_1 \vartheta, \quad \|h^t - h^{t-k}\|^2 \leq \tau k_1 \vartheta, \quad \sum_l \|\lambda_l^t - \lambda_l^{t-k}\|^2 \leq \tau k_1 \vartheta, \quad \forall 1 \leq k \leq \tau,$$

where $k_1 > 0$ is a constant.

Setting 1 (Bounded $|\mathbf{A}^t|$) $|\mathbf{A}^t| \leq M, \forall t$, i.e., an upper bound is set for the number of cutting planes.

Setting 2 (Setting of c_1^t, c_2^t) $c_1^t = \frac{1}{\rho_1(t+1)^{\frac{1}{6}}} \geq \underline{c}_1$ and $c_2^t = \frac{1}{\rho_2(t+1)^{\frac{1}{6}}} \geq \underline{c}_2$ are nonnegative non-increasing sequences, where \underline{c}_1 and \underline{c}_2 are positive constants and meet $M\underline{c}_1^2 + N\underline{c}_2^2 \leq \frac{\varepsilon^2}{4}$.

Theorem 1 (Iteration complexity) Suppose Assumption 1 and 2 hold, if we set the step-sizes of primal variables as $\eta_{\mathbf{w}} = \eta_{\mathbf{z}} = \eta_h = \frac{2}{L + \rho_1 |\mathbf{A}^t| L^2 + \rho_2 N L^2 + 8(\frac{|\mathbf{A}^t| \gamma L^2}{\rho_1 (c_1^t)^2} + \frac{N \gamma L^2}{\rho_2 (c_2^t)^2})}$. And we set the step-sizes of

dual variables as $\rho_1 < \min\{\frac{2}{L+2c_1^t}, \frac{1}{15\tau k_1 N L^2}\}$ and $\rho_2 \leq \frac{2}{L+2c_2^t}$, respectively. For a given ε , we have:

$$T(\varepsilon) \sim \mathcal{O}\left(\max\left\{\left(\frac{4M\sigma_1^2}{\rho_1^2} + \frac{4N\sigma_2^2}{\rho_2^2}\right)^{\frac{1}{3}} \frac{1}{\varepsilon^6}, \left(\frac{4(d_6 + \frac{\rho_2(N-S)L^2}{2})^{\frac{2}{3}}}{\varepsilon^2} (d + k_d(\tau-1)) d_5 + (T_1 + 2)^{\frac{1}{3}}\right)^3\right\}\right), \quad (22)$$

where $\sigma_1, \sigma_2, \gamma, k_d, \bar{d}, d_5, d_6$ and T_1 are constants. The detailed proof is given in Appendix A.

There exists a wide array of works regarding the convergence analysis of various algorithms for nonconvex/convex optimization problems involved in machine learning [25, 53]. Our analysis, however, differs from existing works in two aspects. First, we solve the non-convex PD-DRO in an *asynchronous distributed manner*. To our best knowledge, there are few works focusing on solving the DRO in a distributed manner. Compared to solving the non-convex PD-DRO in a centralized manner, solving it in an *asynchronous distributed manner* poses significant challenges in algorithm design and convergence analysis. Secondly, we do not assume the inner problem can be solved nearly optimally for each outer iteration, which is numerically difficult to achieve in practice [4]. Instead, ASPIRE-EASE is *single loop* and involves simple gradient projection operation at each step.

6 Experiment

In this section, we conduct experiments on four real-world datasets to assess the performance of the proposed method. Specifically, we evaluate the robustness against data heterogeneity, robustness against malicious attacks and efficiency of the proposed method. Ablation study is also carried out to demonstrate the excellent performance of ASPIRE-EASE.

6.1 Datasets and Baseline Methods

We compare the proposed ASPIRE-EASE with baseline methods based on SHL [20], Person Activity [26], Single Chest-Mounted Accelerometer (SM-AC) [9] and Fashion MNIST [51] datasets. The baseline methods include Ind_j (learning the model from an individual worker j), Mix_{Even} (learning the model from all workers with even weights using ASPIRE), FedAvg [33], AFL [35] and DRFA-Prox [16]. The detailed descriptions of datasets and baselines are given in Appendix C.

In our empirical studies, since the downstream tasks are multi-class classification, the cross entropy loss is used on each worker (i.e., $\mathcal{L}_j(\cdot), \forall j$). For SHL, Person Activity, and SM-AC datasets, we adopt the deep multilayer perceptron [49] as the base model. And we use the same logistic regression model as in [35, 16] for Fashion MNIST dataset. The base models are trained with SGD. More details are given in Appendix C. Following related works in this direction [41, 35, 16], worst case performance are reported for the comparison of robustness. Specifically, we use Acc_w and Loss_w to represent the worst case test accuracy and training loss (i.e., the test accuracy and training loss on the worker with worst performance), respectively. We also report the standard deviation Std of $[\text{Acc}_1, \dots, \text{Acc}_N]$ (the test accuracy on every worker). In the experiment, S is set as 1, that means

Table 1: Performance comparisons based on Acc_w (%) \uparrow , Loss_w \downarrow and Std \downarrow (\uparrow and \downarrow respectively denote higher scores represent better performance and lower scores represent better performance). The boldfaced digits represent the best results, “—” represents not available.

Model	SHL			Person Activity			SC-MA			Fashion MNIST		
	$\text{Acc}_w \uparrow$	$\text{Loss}_w \downarrow$	$\text{Std} \downarrow$	$\text{Acc}_w \uparrow$	$\text{Loss}_w \downarrow$	$\text{Std} \downarrow$	$\text{Acc}_w \uparrow$	$\text{Loss}_w \downarrow$	$\text{Std} \downarrow$	$\text{Acc}_w \uparrow$	$\text{Loss}_w \downarrow$	$\text{Std} \downarrow$
$\max\{\text{Ind}_j\}$	19.06±0.65	—	29.1	49.38±0.08	—	8.32	22.56±0.78	—	17.5	—	—	—
Mix _{Even}	69.87±3.10	0.806±0.018	4.81	56.31±0.69	1.165±0.017	3.00	49.81±0.21	1.424±0.024	6.99	66.80±0.18	0.784±0.003	10.1
FedAvg [33]	69.96±3.07	0.802±0.023	5.21	56.28±0.63	1.154±0.019	3.13	49.53±0.96	1.441±0.015	7.17	66.58±0.39	0.781±0.002	10.2
AFL [35]	78.11±1.99	0.582±0.021	1.87	58.39±0.37	1.081±0.014	0.99	54.56±0.79	1.172±0.018	3.50	77.32±0.15	0.703±0.001	1.86
DRFA-Prox [16]	78.34±1.46	0.532±0.034	1.85	58.62±0.16	1.096±0.037	1.26	54.61±0.76	1.151±0.039	4.69	77.95±0.51	0.702±0.007	1.34
ASPIRE-EASE	79.16±1.13	0.515±0.019	1.02	59.43±0.44	1.053±0.010	0.82	56.31±0.29	1.127±0.021	3.16	78.82±0.07	0.696±0.004	1.01
ASPIRE-EASE _{per}	78.94±1.27	0.521±0.023	1.36	59.54±0.21	1.051±0.016	0.79	56.71±0.16	1.119±0.028	3.48	78.73±0.06	0.698±0.006	1.09

the master will make an update once it receives a message. Each experiment is repeated 10 times, both mean and standard deviations are reported. We implement our model with PyTorch and conduct all the experiments on a server with two TITAN V GPUs.

6.2 Results

Robustness against Data Heterogeneity. We first assess the robustness of the proposed ASPIRE-EASE by comparing it with baseline methods when data are heterogeneously distributed across different workers. Specifically, we compare the Acc_w , Loss_w and Std of different methods on all datasets. The performance comparison results are shown in Table 1. In this table, we can observe that $\max\{\text{Ind}_j\}$, which represents the best performance of individual training over all workers, exhibits the worst robustness on SHL, Person Activity, and SC-MA. This is because individual training ($\max\{\text{Ind}_j\}$) only learns from the data in its local worker and cannot generalize to other workers due to different data distributions. Note that $\max\{\text{Ind}_j\}$ is unavailable for Fashion MNIST since each worker only contains one class of data and cross entropy loss cannot be used in this case. $\max\{\text{Ind}_j\}$ also does not have Loss_w , since Ind_j is trained only on individual worker j . The FedAvg and Mix_{Even} exhibit better performance than $\max\{\text{Ind}_j\}$ since they consider the data from all workers. Nevertheless, FedAvg and Mix_{Even} only assign the fixed weight for each worker. AFL is more robust than FedAvg and Mix_{Even} since it not only utilizes the data from all workers but also considers optimizing the weight of each worker. DRFA-Prox outperforms AFL since it also considers the prior distribution and regards it as a regularizer in the objective function. Finally, we can observe that the proposed ASPIRE-EASE shows excellent robustness, which can be attributed to two factors: 1) ASPIRE-EASE considers data from all workers and can optimize the weight of each worker; 2) compared with DRFA-Prox which uses prior distribution as a regularizer, the prior distribution is incorporated within the constraint in our formulation (Eq. 4), which can be leveraged more effectively. And it is seen that ASPIRE-EASE can perform periodic communication since ASPIRE-EASE_{per}, which represents ASPIRE-EASE with periodic communication, also has excellent performance.

Within ASPIRE-EASE, the level of robustness can be controlled by adjusting Γ . Specially, when $\Gamma = 0$, we obtain a nominal optimization problem in which no adversarial distribution is considered. The size of the uncertainty set will increase with Γ (when $\Gamma \leq N$), which enhances the adversarial robustness of the model. As shown in Figure 1, the robustness of ASPIRE-EASE can be gradually enhanced when Γ increases. More results are available in Figure C2 of Appendix C.

Robustness against Malicious Attacks. To assess the model robustness against malicious attacks, malicious workers with backdoor attacks [1, 48], which attempt to mislead the model training process, are added to the distributed system. Following [14], we report the success attack rate of backdoor attacks for comparison. It can be calculated by checking how many instances in the backdoor dataset can be misled and categorized into the target labels. Lower success attack rates indicate more robustness against backdoor attacks. The comparison results are summarized in Table 2 and more detailed settings of backdoor attacks are available in Appendix C. In Table 2, we observe that AFL can be attacked easily since it could assign higher weights to malicious workers. Compared to AFL, FedAvg and Mix_{Even} achieve relatively lower success attack rates since they assign equal weights to the malicious workers and other workers. DRFA-Prox can achieve even lower success attack rates since it can leverage the prior distribution to assign lower weights for malicious workers. The proposed ASPIRE-EASE achieves the lowest success attack rates since it can leverage the prior distribution more effectively. Specifically, it will assign lower weights to malicious workers with tight theoretical guarantees.

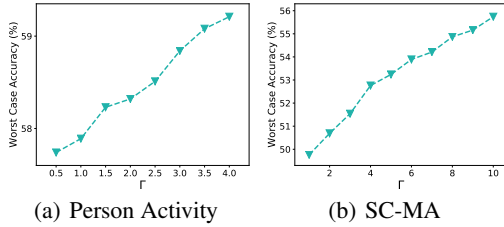


Figure 1: Γ control the degree of robustness (worst case performance in the problem) on (a) Person Activity, (b) SC-MA datasets.

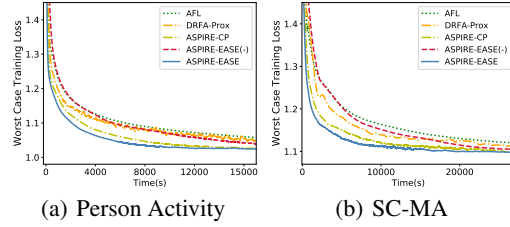


Figure 2: Comparison of the convergence time on worst case worker on (a) Person Activity, (b) SC-MA datasets.

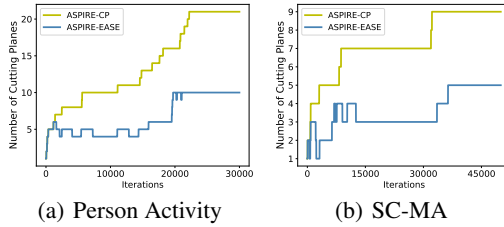


Figure 3: Comparison of ASPIRE-CP and ASPIRE-EASE regarding the number of cutting planes on (a) Person Activity, (b) SC-MA datasets.

Table 2: Performance comparisons about the success attack rate (%) \downarrow . The boldfaced digits represent the best results.

Model	SHL	Person Activity	SC-MA	Fashion MNIST
MixEven	36.21 \pm 2.23	34.32 \pm 2.18	52.14 \pm 2.89	83.18 \pm 2.07
FedAvg [33]	38.15 \pm 3.02	33.25 \pm 2.49	55.39 \pm 3.13	82.04 \pm 1.84
AFL [35]	68.63 \pm 4.24	43.66 \pm 3.87	75.81 \pm 4.03	90.04 \pm 2.52
DRFA-Prox [16]	21.23 \pm 3.63	27.27 \pm 3.31	30.79 \pm 3.65	63.24 \pm 2.47
ASPIRE-EASE	9.17\pm1.65	22.36\pm2.33	14.51\pm3.21	45.10\pm1.64

Efficiency. In Figure 2, we compare the convergence speed of the proposed ASPIRE-EASE with AFL and DRFA-Prox by considering different communication and computation delays for each worker. The proposed ASPIRE-EASE has two variants, ASPIRE-CP (ASPIRE with cutting plane method), ASPIRE-EASE(-) (ASPIRE-EASE without asynchronous setting). More results are available in Figure C3 of Appendix C. Based on the comparison, we can observe that the proposed ASPIRE-EASE generally converges faster than baseline methods and its two variants. This is because 1) compared with AFL, DRFA-Prox, and ASPIRE-EASE(-), ASPIRE-EASE is an asynchronous algorithm in which the master updates its parameters only after receiving the updates from active workers instead of all workers; 2) unlike DRFA-Prox, the master in ASPIRE-EASE only needs to communicate with active workers once per iteration; 3) compared with ASPIRE-CP, ASPIRE-EASE utilizes active set method instead of cutting plane method, which is more efficient. It is seen from Figure 2 that, the convergence speed of ASPIRE-EASE mainly benefits from the asynchronous setting.

Ablation Study. For ASPIRE, compared with cutting plane method, EASE is more efficient since it considers removing the inactive cutting planes. To demonstrate the efficiency of EASE, we firstly compare ASPIRE-EASE with ASPIRE-CP concerning the number of cutting planes used during the training. In Figure 3, we can observe that ASPIRE-EASE uses fewer cutting planes than ASPIRE-CP, thus is more efficient. The convergence speed of ASPIRE-EASE and ASPIRE-CP in Figure 2 also suggests that ASPIRE-EASE converges much faster than ASPIRE-CP. More results are available in Figure C3 and C4, Appendix C.

7 Conclusion

In this paper, we present ASPIRE-EASE method to effectively solve the distributed distributionally robust optimization problem with non-convex objectives. In addition, CD -norm uncertainty set has been proposed to effectively incorporate the prior distribution into the problem formulation, which allows for flexible adjustment of the degree of robustness of DRO. Theoretical analysis has also been conducted to analyze the convergence properties and the iteration complexity of ASPIRE-EASE. ASPIRE-EASE exhibits strong empirical performance on multiple real-world datasets and is effective in tackling DRO problems in a fully distributed and asynchronous manner. In the future work, more uncertainty sets could be designed for our framework and more update rule for variables in ASPIRE could be considered.

Acknowledgments and Disclosure of Funding

The work of Yang Jiao and Kai Yang was supported in part by the Fundamental Research Funds for the Central Universities of China, in part by the Shenzhen Institute of Artificial Intelligence and Robotics for Society (AIRS), in part by the National Natural Science Foundation of China under Grant 61771013, and in part by the Fundamental Research Funds of Shanghai Jiading District.

References

- [1] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov. How to backdoor federated learning. In *International Conference on Artificial Intelligence and Statistics*, pages 2938–2948. PMLR, 2020.
- [2] A. Ben-Tal and A. Nemirovski. Robust solutions of uncertain linear programs. *Operations research letters*, 25(1):1–13, 1999.
- [3] A. Ben-Tal, L. El Ghaoui, and A. Nemirovski. *Robust optimization*. Princeton university press, 2009.
- [4] D. P. Bertsekas. Nonlinear programming. *Journal of the Operational Research Society*, 48(3): 334–334, 1997.
- [5] D. Bertsimas and M. Sim. The price of robustness. *Operations research*, 52(1):35–53, 2004.
- [6] D. Bertsimas, I. Dunning, and M. Lubin. Reformulation versus cutting-planes for robust optimization. *Computational Management Science*, 13(2):195–217, 2016.
- [7] J. Blanchet and K. Murthy. Quantifying distributional model risk via optimal transport. *Mathematics of Operations Research*, 44(2):565–600, 2019.
- [8] L. Bottou, F. E. Curtis, and J. Nocedal. Optimization methods for large-scale machine learning. *Siam Review*, 60(2):223–311, 2018.
- [9] P. Casale, O. Pujol, and P. Radeva. Personalization and user verification in wearable systems using biometric walking patterns. *Personal and Ubiquitous Computing*, 16(5):563–580, 2012.
- [10] T.-H. Chang, M. Hong, W.-C. Liao, and X. Wang. Asynchronous distributed ADMM for large-scale optimization—Part I: Algorithm and convergence analysis. *IEEE Transactions on Signal Processing*, 64(12):3118–3130, 2016.
- [11] Y. Chen, Y. Ning, M. Slawski, and H. Rangwala. Asynchronous online federated learning for edge devices with Non-IID data. In *2020 IEEE International Conference on Big Data (Big Data)*, pages 15–24. IEEE, 2020.
- [12] A. Cohen, A. Daniely, Y. Drori, T. Koren, and M. Schain. Asynchronous stochastic optimization robust to arbitrary delays. *Advances in Neural Information Processing Systems*, 34:9024–9035, 2021.
- [13] R. Cole. Parallel merge sort. *SIAM Journal on Computing*, 17(4):770–785, 1988.
- [14] J. Dai, C. Chen, and Y. Li. A backdoor attack against LSTM-based text classification systems. *IEEE Access*, 7:138872–138878, 2019.
- [15] E. Delage and Y. Ye. Distributionally robust optimization under moment uncertainty with application to data-driven problems. *Operations research*, 58(3):595–612, 2010.
- [16] Y. Deng, M. M. Kamani, and M. Mahdavi. Distributionally robust federated averaging. *arXiv preprint arXiv:2102.12660*, 2021.
- [17] J. C. Duchi and H. Namkoong. Learning models with uniform performance via distributionally robust optimization. *The Annals of Statistics*, 49(3):1378–1406, 2021.
- [18] R. Gao and A. J. Kleywegt. Distributionally robust stochastic optimization with wasserstein distance. *arXiv preprint arXiv:1604.02199*, 2016.

- [19] G. Geraci, M. Wildemeersch, and T. Q. Quek. Energy efficiency of distributed signal processing in wireless networks: A cross-layer analysis. *IEEE Transactions on Signal Processing*, 64(4): 1034–1047, 2015.
- [20] H. Gjoreski, M. Ciliberto, L. Wang, F. J. O. Morales, S. Mekki, S. Valentin, and D. Roggen. The university of sussex-huawei locomotion and transportation dataset for multimodal analytics with mobile devices. *IEEE Access*, 6:42592–42604, 2018.
- [21] B. L. Gorissen, İ. Yanıkoğlu, and D. den Hertog. A practical guide to robust optimization. *Omega*, 53:124–137, 2015.
- [22] F. Haddadpour, M. M. Kamani, M. Mahdavi, and A. Karbasi. Learning distributionally robust models at scale via composite optimization. *arXiv preprint arXiv:2203.09607*, 2022.
- [23] Y. Hu, X. Chen, and N. He. On the bias-variance-cost tradeoff of stochastic optimization. *Advances in Neural Information Processing Systems*, 34, 2021.
- [24] J. Jiang, W. Zhang, J. Gu, and W. Zhu. Asynchronous decentralized online learning. *Advances in Neural Information Processing Systems*, 34:20185–20196, 2021.
- [25] C. Jin, P. Netrapalli, and M. Jordan. What is local optimality in nonconvex-nonconcave minimax optimization? In *International Conference on Machine Learning*, pages 4880–4889. PMLR, 2020.
- [26] B. Kaluža, V. Mirchevska, E. Dovgan, M. Luštrek, and M. Gams. An agent-based approach to care in independent living. In *International joint conference on ambient intelligence*, pages 177–186. Springer, 2010.
- [27] S. P. Karimireddy, S. Kale, M. Mohri, S. J. Reddi, S. U. Stich, and A. T. Suresh. SCAFFOLD: Stochastic Controlled Averaging for On-Device Federated Learning. 2019.
- [28] D. Kuhn, P. M. Esfahani, V. A. Nguyen, and S. Shafieezadeh-Abadeh. Wasserstein distributionally robust optimization: Theory and applications in machine learning. In *Operations Research & Management Science in the Age of Analytics*, pages 130–166. INFORMS, 2019.
- [29] D. Levy, Y. Carmon, J. C. Duchi, and A. Sidford. Large-scale methods for distributionally robust optimization. *Advances in Neural Information Processing Systems*, 33:8847–8860, 2020.
- [30] W.-H. Liao and Y.-T. Huang. Investigation of DNN model robustness using heterogeneous datasets. In *2020 25th International Conference on Pattern Recognition (ICPR)*, pages 4393–4397. IEEE, 2021.
- [31] T. Lin, C. Jin, and M. Jordan. On gradient descent ascent for nonconvex-concave minimax problems. In *International Conference on Machine Learning*, pages 6083–6093. PMLR, 2020.
- [32] S. Lu, I. Tsaknakis, M. Hong, and Y. Chen. Hybrid block successive approximation for one-sided non-convex min-max problems: algorithms and applications. *IEEE Transactions on Signal Processing*, 68:3676–3691, 2020.
- [33] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017.
- [34] S. Mehrotra and D. Papp. A cutting surface algorithm for semi-infinite convex programming with an application to moment robust optimization. *SIAM Journal on Optimization*, 24(4): 1670–1697, 2014.
- [35] M. Mohri, G. Sivek, and A. T. Suresh. Agnostic federated learning. In *International Conference on Machine Learning*, pages 4615–4625. PMLR, 2019.
- [36] A. Nedic and A. Ozdaglar. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control*, 54(1):48–61, 2009.
- [37] Y. Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2003.

- [38] Q. Qi, Z. Guo, Y. Xu, R. Jin, and T. Yang. An online method for a class of distributionally robust optimization with non-convex objectives. *Advances in Neural Information Processing Systems*, 34:10067–10080, 2021.
- [39] J. Qian, X. Fafoutis, and L. K. Hansen. Towards federated learning: Robustness analytics to data heterogeneity. *arXiv preprint arXiv:2002.05038*, 2020.
- [40] J. Qian, L. K. Hansen, X. Fafoutis, P. Tiwari, and H. M. Pandey. Robustness analytics to data heterogeneity in edge computing. *Computer Communications*, 164:229–239, 2020.
- [41] Q. Qian, S. Zhu, J. Tang, R. Jin, B. Sun, and H. Li. Robust optimization over multiple domains. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4739–4746, 2019.
- [42] H. Rahimian and S. Mehrotra. Distributionally robust optimization: A review. *arXiv preprint arXiv:1908.05659*, 2019.
- [43] S. Sicari, A. Rizzardi, L. A. Grieco, and A. Coen-Portisini. Security, privacy and trust in Internet of Things: The road ahead. *Computer networks*, 76:146–164, 2015.
- [44] K. Singhal, H. Sidahmed, Z. Garrett, S. Wu, J. Rush, and S. Prakash. Federated reconstruction: Partially local federated learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- [45] T. Subramanya and R. Riggio. Centralized and federated learning for predictive VNF autoscaling in multi-domain 5G networks and beyond. *IEEE Transactions on Network and Service Management*, 18(1):63–78, 2021.
- [46] J. Sun, T. Chen, G. B. Giannakis, and Z. Yang. Communication-efficient distributed learning via lazily aggregated quantized gradients. *arXiv preprint arXiv:1909.07588*, 2019.
- [47] K. K. Thekumparampil, P. Jain, P. Netrapalli, and S. Oh. Efficient algorithms for smooth minimax optimization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [48] B. Wang, Y. Yao, S. Shan, H. Li, B. Viswanath, H. Zheng, and B. Y. Zhao. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 707–723. IEEE, 2019.
- [49] Z. Wang, W. Yan, and T. Oates. Time series classification from scratch with deep neural networks: A strong baseline. In *2017 International joint conference on neural networks (IJCNN)*, pages 1578–1585. IEEE, 2017.
- [50] W. Wiesemann, D. Kuhn, and B. Rustem. Robust Markov decision processes. *Mathematics of Operations Research*, 38(1):153–183, 2013.
- [51] H. Xiao, K. Rasul, and R. Vollgraf. Fashion-MNIST: A novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- [52] Z. Xu, H. Zhang, Y. Xu, and G. Lan. A unified single-loop alternating gradient projection algorithm for nonconvex-concave and convex-nonconcave minimax problems. *arXiv preprint arXiv:2006.02032*, 2020.
- [53] Z. Xu, J. Shen, Z. Wang, and Y. Dai. Zeroth-order alternating randomized gradient projection algorithms for general nonconvex-concave minimax problems. *arXiv preprint arXiv:2108.00473*, 2021.
- [54] K. Yang, Y. Wu, J. Huang, X. Wang, and S. Verdú. Distributed robust optimization for communication networks. In *IEEE INFOCOM 2008-The 27th Conference on Computer Communications*, pages 1157–1165. IEEE, 2008.
- [55] K. Yang, J. Huang, Y. Wu, X. Wang, and M. Chiang. Distributed robust optimization (DRO), part I: Framework and example. *Optimization and Engineering*, 15(1):35–67, 2014.
- [56] İ. Yanıkoğlu, B. L. Gorissen, and D. den Hertog. A survey of adjustable robust optimization. *European Journal of Operational Research*, 277(3):799–813, 2019.

- [57] S. Zawad, A. Ali, P.-Y. Chen, A. Anwar, Y. Zhou, N. Baracaldo, Y. Tian, and F. Yan. Curse or redemption? how data heterogeneity affects the robustness of federated learning. *arXiv preprint arXiv:2102.00655*, 2021.
- [58] R. Zhang and J. Kwok. Asynchronous distributed ADMM for consensus optimization. In *International conference on machine learning*, pages 1701–1709. PMLR, 2014.
- [59] X. Zhou. On the fenchel duality between strong convexity and lipschitz continuous gradient. *arXiv preprint arXiv:1803.06573*, 2018.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#) See Section 1.
 - (b) Did you describe the limitations of your work? [\[Yes\]](#) See Section 7.
 - (c) Did you discuss any potential negative societal impacts of your work? [\[N/A\]](#) There is no potential negative societal impact of our work.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#) See Section 5.
 - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) See Appendix A and B.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#) The references of the data used in this paper are added in Section 6.1.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) Section C.2.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[Yes\]](#) See Section 6.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) See Section 6.1.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [\[N/A\]](#)
 - (b) Did you mention the license of the assets? [\[N/A\]](#)
 - (c) Did you include any new assets either in the supplemental material or as a URL? [\[N/A\]](#)
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [\[N/A\]](#)
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[N/A\]](#)
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#)
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#)
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#)