

---

# On-Demand Sampling: Learning Optimally from Multiple Distributions \*

---

Nika Haghtalab, Michael I. Jordan, and Eric Zhao

University of California, Berkeley

## Abstract

Societal and real-world considerations such as robustness, fairness, social welfare and multi-agent tradeoffs have given rise to multi-distribution learning paradigms, such as *collaborative* [5], *group distributionally robust* [36], and *fair federated learning* [27]. In each of these settings, a learner seeks to minimize its worst-case loss over a set of  $n$  predefined distributions, while using as few samples as possible. In this paper, we establish the optimal sample complexity of these learning paradigms and give algorithms that meet this sample complexity. Importantly, our sample complexity bounds exceed that of the sample complexity of learning a single distribution only by an additive factor of  $\frac{n \log(n)}{\epsilon^2}$ . These improve upon the best known sample complexity of agnostic federated learning by Mohri et al. [27] by a multiplicative factor of  $n$ , the sample complexity of collaborative learning by Nguyen and Zakynthinou [29] by a multiplicative factor  $\frac{\log n}{\epsilon^3}$ , and give the first sample complexity bounds for the *group DRO* objective of Sagawa et al. [36]. To achieve optimal sample complexity, our algorithms learn to sample and learn from distributions on demand. Our algorithm design and analysis extends stochastic optimization techniques to solve zero-sum games in a new stochastic setting.

## 1 Introduction

Pervasive needs for robustness, fairness, and multi-agent collaboration in learning have given rise to multi-distribution learning paradigms (e.g., [5, 36, 27, 12]). In these settings, we seek to learn a model that performs well on *any distribution* in a pre-defined set of interest. For fairness considerations, these distributions may represent heterogeneous populations of different protected or socio-economic attributes; in robustness applications, they may capture a learner’s uncertainty regarding the true underlying task; and in multi-agent collaborative or federated applications, they may represent agent-specific learning tasks. In these applications, the performance and optimality of a model is measured by its worst test-time performance on a distribution in the set. We are concerned with this fundamental problem of designing sample-efficient multi-distribution learning algorithms.

The sample complexity of multi-distribution learning differs from that of learning a single distribution in several ways. On one hand, learning tasks of varying difficulty require different numbers of samples. On the other hand, similarity or overlap among learning tasks may obviate the need to sample from some distributions. This makes the use of a fixed per-distribution sample budget highly inefficient and suggests that optimal multi-distribution learning algorithms should *sample on demand*. That is, algorithms should take additional samples *whenever they need them* and *from whichever distribution* they want them. On-demand sampling is especially appropriate when some population data is scarce (as in fairness mechanisms in which samples are amended [32]); when the designer can actively

---

\*Authors are ordered alphabetically. Addresses: nika@berkeley.edu, jordan@cs.berkeley.edu, eric.zh@berkeley.edu.

Problem	Sample Complexity	Thm	Best Previous Result
Collab. Learning UB	$\varepsilon^{-2} (\log  \mathcal{H}  + n \log(\frac{n}{\delta}))$	[4.1]	$\varepsilon^{-5} \log(\frac{1}{\varepsilon}) \log(\frac{n}{\delta}) (\log  \mathcal{H}  + n)$ [29]
Collab. Learning LB	$\varepsilon^{-2} (\log  \mathcal{H}  + n \log(\frac{k}{\delta}))$	[4.2]	$\varepsilon^{-1} n \log(k/\delta)$ [5]
GDRO/AFL UB	$\varepsilon^{-2} (\log  \mathcal{H}  + n \log(\frac{n}{\delta}))$	[4.1]	$\varepsilon^{-2} (n \log  \mathcal{H}  + n \log(\frac{n}{\delta}))$ [27]
GDRO/AFL UB	$\varepsilon^{-2} (D_{\mathcal{H}} + n \log(\frac{n}{\delta}))$	[5.1]	N/A
(Training error convg.)	$\varepsilon^{-2} (D_{\mathcal{H}} + n \log(\frac{n}{\delta}))$	[5.2]	$\varepsilon^{-2} D_{\mathcal{H}}$ (expected convergence only) [36]

Table 1: This table gives upper (*UB*) and lower bounds (*LB*) on the sample complexity of learning model class  $H$  on  $n$  distributions. For the collaborative learning and AFL settings, the sample complexity upper bounds refer to the problem of learning a randomized model of worst-case error  $\text{OPT} + \varepsilon$  or a deterministic classifier of worst-case error  $2\text{OPT} + \varepsilon$ . For the GDRO setting, sample complexity refers to learning a deterministic model with worst-case error of  $\text{R-OPT} + \varepsilon$ , where  $\text{R-OPT}$  is the best worst-case error attainable in a convex compact model space  $H$ .  $D_{\mathcal{H}}$  denotes the Bregman radius of  $H$ , and  $k = \min \{n, \log |\mathcal{H}|\}$ . Sample complexity bounds of Collaborative and Agnostic federated learning in existing works, extend to VC dimension and Rademacher complexity. Our results also extend to VC dimension under some assumptions.

perturb datasets towards rare or atypical instances (such as in robustness applications [21, 44]); or when sample sets represent agents' contributions to an interactive multi-agent system [27, 6].

Blum et al. [5] demonstrated the benefit of on-demand sampling in the *collaborative learning* setting, where all data distributions are realizable with respect to the same target classifier. This line of work established that learning  $n$  distributions on-demand takes  $\tilde{O}(\log(n))$  times the sample complexity of learning a single realizable distribution [5, 8, 29], whereas relying on batched uniform convergence takes  $\tilde{\Omega}(n)$  times that of learning a single distribution [5]. However, beyond the realizable setting, the best known multi-distribution learning results fall short of this promise: existing on-demand sample complexity bounds for agnostic collaborative learning have highly suboptimal dependence on  $\varepsilon$ , requiring  $\tilde{O}(\log(n)/\varepsilon^3)$  times the sample complexity of agnostically learning a single distribution [29]. On the other hand, agnostic federated learning bounds [27] have been studied only for algorithms that sample in one large batch and thus require  $\tilde{\Omega}(n)$  times the sample complexity of a single learning task. Moreover, the test-time performance of some key multi-distribution methods, such as group distributionally robust optimization [36], have not been studied from a provable or mathematical perspective before.

In this paper, we give a general framework for obtaining *optimal and on-demand sample complexity* for three multi-distribution learning settings. Table 1 summarizes our results. All three settings consider a set  $\mathcal{D}$  of  $n$  distributions and a model class  $\mathcal{H}$ . They evaluate the performance of a model  $h$  (or a distribution over models) by its worst-case performance,  $\max_{D \in \mathcal{D}} \text{Risk}_D(h)$ . As a benchmark, they consider the worst-case loss of the best model, i.e.,  $\text{OPT} = \min_{h^* \in \mathcal{H}} \max_{D \in \mathcal{D}} \text{Risk}_D(h^*)$ . Importantly, all of our sample complexity upper bounds demonstrate only an *additive increase of  $\varepsilon^{-2} n \log(n/\delta)$  over the sample complexity of a single learning task*, compared to the multiplicative factor increase required by existing works.

- *Collaborative learning of Blum et al. [5]*: For agnostic collaborative learning, our Theorem 4.1 gives a randomized and a deterministic model that achieve performance guarantees of  $\text{OPT} + \varepsilon$  and  $2\text{OPT} + \varepsilon$ , respectively. Our algorithms have an optimal sample complexity of  $O(\frac{1}{\varepsilon^2} (\log(|\mathcal{H}|) + n \log(\frac{n}{\delta})))$ . This improves upon the work of Nguyen and Zakynthinou [29] in two ways. First, it provides error bounds of  $\text{OPT} + \varepsilon$  for randomized classifiers, where only  $2\text{OPT} + \varepsilon$  was previously established. Second, it improves the upper bound of Nguyen and Zakynthinou [29] by a multiplicative factor of  $\log(n)/\varepsilon^3$ . In Theorem 4.2, we give a matching lower bound on this sample complexity, thereby establishing the optimality of our algorithms.
- *Group distributionally robust learning (group DRO) of Sagawa et al. [36]*: For group DRO, we consider a convex and compact model space  $\mathcal{H}$ . Our Theorem 5.1 studies a model that achieves an  $\text{OPT} + \varepsilon$  guarantee on the worst-case test-time performance of the model with an on-demand sample complexity of  $O(\frac{1}{\varepsilon^2} (D_{\mathcal{H}} + n \log(\frac{n}{\delta})))$ . Our results also imply a high-probability bound

for the convergence of group DRO *training error* that improves upon the (expected) convergence guarantees of Sagawa et al. [36] by a factor of  $n$ .

- *Agnostic federated learning of [27]*: For agnostic federated learning, we consider a finite class of hypotheses. Our Theorems 4.1 and 5.1 show that on-demand sampling can accelerate the generalization of agnostic federated learning by a factor of  $n$  compared to batch results established by Mohri et al. [27]. Our results also imply matching high-probability bounds to Mohri et al. [27] on the convergence of the training error in the batched setting.

To achieve these results, we contribute new insights and techniques for solving stochastic zero-sum games with sources of randomization that differ in both cost and quality. We frame the multi-distribution learning problems as a stochastic zero-sum game with uncertain payoffs and utilize stochastic mirror descent and a variational perspective to solve the game. In this case, the maximizing player can be interpreted as a weight vector for distributions  $\mathcal{D}$ , specifying from which distributions future on-demand samples should be taken. These on-demand samples form a stochastic gradient for the players. However, the quality of these estimators, the number of samples needed for them, and whether they can be reused later on, differs between the two players. We extend the Stochastic Mirror Descent framework to optimally trade off these asymmetric needs for samples. In Section 3 we give an overview of this approach and its technical challenges and contributions.

## 1.1 Related Work

**Learning models.** There are many lines of work that study multi-distribution learning but which have evolved independently in separate communities. The field of *collaborative learning* concerns the learning of a shared machine learning model by multiple *stakeholders* that each desire a model with low error on their own data distribution. The line of work initiated by Blum et al. [5] studies on-demand sample complexity bounds for realizable collaborative learning and was later extended to several related settings (e.g., [29, 8, 7, 30]). The agnostic federated learning framework of Mohri et al. [27] poses an equivalent of the multi-distribution learning objective as a fair and intuitive target for federated learning algorithms, and studies it in the offline setting with data-dependent analysis. Multi-distribution learning also arises in distributionally robust optimization [4] as the Group DRO problem [17], which is motivated by deep learning applications with multiple deployment domains or protected demographics. These works focus on an empirical perspective, but have discussed training error convergence in offline settings [17, 36, 37]. Multi-distribution learning is also related to a line of work on multi-source domain adaptation (e.g., [3, 24]) and multi-group fairness notions (e.g., [35, 38, 13]). We describe these parallel threads in more detail in Section A.

**Stochastic game equilibria.** Our approach relates to a line of research on using online algorithms to find min-max equilibria by playing no-regret algorithms against one another [34, 15, 31, 9, 10]. Online mirror descent (OMD) is one well-studied family of methods that can find approximate minima of convex functions, and also approximate min-max equilibria of convex-concave games, with high probability using noisy first-order information [33, 28, 16, 2]. We bring these online learning tools to bear on the problem of finding saddle points in robust optimization formulations. The primary technical difference between multi-distribution learning and traditional saddle-point optimization problems is that we have sample access to distributions instead of noisy local gradients.

## 2 Preliminaries

Let  $\mathcal{X}$  be an instance space,  $\mathcal{Y}$  a label space, and  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$  a space of datapoints. A data distribution  $D$  is a joint probability distribution over  $\mathcal{Z}$ . We consider a hypothesis class  $\mathcal{H}$  of a subset of functions mapping  $\mathcal{X}$  to  $\mathcal{Y}$ . With each distribution  $D$ , define a loss function  $\ell_D : \mathcal{H} \times \mathcal{Z} \rightarrow [0, 1]$  measuring the loss of hypothesis  $h$  on data point  $z \in \mathcal{Z}$ . We write  $\ell_D$  as  $\ell$  when  $D$  is clear from context. We denote the expected loss, i.e. risk, of a hypothesis  $h \in \mathcal{H}$  under a data distribution  $D \in \mathcal{D}$  by:

$$\text{Risk}_D(h) := \mathbb{E}_{(x,y) \sim D} [\ell_D(h, (x, y))].$$

Importantly, we only assume that  $\ell_D$ 's are bounded and make no other assumptions on losses or distributions. For a distribution over the hypothesis class,  $p \in \Delta\mathcal{H}$ , and a distribution over data distributions,  $q \in \Delta\mathcal{D}$ , we refer to their expected loss by  $\text{Risk}_q(p) := \mathbb{E}_{D \sim q} [\mathbb{E}_{h \sim p} [\text{Risk}_D(h)]]$ .

**Collaborative Learning.** We will use the *collaborative PAC learning model* of Blum et al. [5] and its agnostic extensions by Nguyen and Zakynthinou [29]. In this setting, the goal is to guarantee small risk for *every* distribution in a collection. Formally, given a set of data distributions  $\mathcal{D} := \{D_1, \dots, D_n\}$ , the goal of the learner is to learn a hypothesis  $h$  such that, with probability  $1 - \delta$ ,

$$\max_{D \in \mathcal{D}} \text{Risk}_D(h) \leq \text{OPT} + \varepsilon, \text{ where } \text{OPT} := \min_{h \in \mathcal{H}} \max_{D \in \mathcal{D}} \text{Risk}_D(h). \quad (1)$$

**Group Distribution Robustness.** We will also study the closely related setting of *group distributionally robust optimization (Group DRO)* of Sagawa et al. [36]. Formally, the group DRO setting considers a model set  $\Theta$  that is a convex compact subset of the Euclidean space and a convex loss function  $\ell : \Theta \times \mathcal{Z} \rightarrow [0, 1]$  that is assumed to be differentiable over  $\Theta$ . Given a set of data distributions  $\mathcal{D} := \{D_1, \dots, D_n\}$ , the learner seeks a model  $\theta \in \Theta$ , such that, with probability  $1 - \delta$ ,

$$\max_{D \in \mathcal{D}} \mathbb{E}_{(x,y) \sim D} [\ell(\theta, (x, y))] \leq \text{R-OPT} + \varepsilon, \text{ where } \text{R-OPT} := \min_{\theta \in \Theta} \max_{D \in \mathcal{D}} \mathbb{E}_{(x,y) \sim D} [\ell(\theta, (x, y))]. \quad (2)$$

There is a close relationship between the Group DRO setting and collaborative learning. In particular, when  $\Theta = \Delta(\mathcal{H})$  and  $\mathcal{H}$  is finite, the two goals are analogous but with two exceptions: first, the Group DRO could return a distribution over functions while collaborative learning requires the solution to be a deterministic function, and second, R-OPT is potentially more competitive than OPT since it allows randomization. We note that the group DRO setting is equivalent to the agnostic federated learning framework of [27], thus our results for DRO extend to that setting as well.

**Sample complexity.** We are interested in the design of algorithms that achieve the above goals while using a small number of samples from distributions  $D_1, \dots, D_n$ . We formalize the sample complexity by the total number of calls made to *example oracles*  $\text{EX}(D_i)$ . Each call  $\text{EX}(D)$  produces an i.i.d. sample from  $D$ . We note that these example oracles also allow us to sample from any mixture distribution  $q \in \Delta\mathcal{D}$ , e.g., by first selecting a  $D_i$  according to the mixture and then calling  $\text{EX}(D_i)$ .

## 2.1 Technical Background

We will use tools and definitions from the literature on zero-sum games and no-regret learning throughout the paper. This section provides a brief overview of these concepts.

**Zero-Sum Games.** A finite two-player zero-sum game is described by the tuple  $(A, A_+, \phi)$  where  $A = \{1, \dots, n\}$  and  $A_+ = \{1, \dots, m\}$  are finite sets of actions and where  $\phi : A \times A_+ \rightarrow [0, C]$ . In this game, the players choose *mixed strategies* over actions sets. These are distributions that are denoted by a vector of probabilities  $p \in \Delta A$  and  $q \in \Delta A_+$ . The expected payoff of mixed strategies is denoted by  $\phi(p, q) = \mathbb{E}_{i \sim p, j \sim q} [\phi(i, j)]$ . The goal of the minimizing player is to minimize this expected payoff and the maximizer seeks to maximize the expected payoff; that is, to solve

$$\min_{p \in \Delta A} \max_{q \in \Delta A_+} \phi(p, q).$$

A pair  $(p, q)$  that solves this optimization problem is called a *min-max equilibrium*. Similarly, a solution is called an  $\varepsilon$ -*min-max equilibrium* if neither player can unilaterally improve their objective by more than  $\varepsilon$ . Formally,  $(p, q)$  is an  $\varepsilon$ -min-max equilibrium if both players' regrets are at most  $\varepsilon$ , i.e.,  $\text{Reg-Min}(p, q) := \phi(p, q) - \min_{i^* \in A} \phi(i^*, q) \leq \varepsilon$  and  $\text{Reg-Max}(p, q) := \max_{j^* \in A_+} \phi(p, j^*) - \phi(p, q) \leq \varepsilon$ . We will next describe methods that find  $\varepsilon$ -min-max equilibria by finding solutions  $(p, q)$  for which  $\text{Reg-Min}(p, q) + \text{Reg-Max}(p, q)$  is at most  $\varepsilon$ . We describe a more general formulation for convex-concave zero-sum games in Appendix B.1 which we will use for the Group DRO problem.

**No-Regret Learning.** We consider an online setting where an arbitrary set of *operators*,  $g^{(1)}, \dots, g^{(T)} \in \mathcal{E}^*$ , is revealed sequentially to a learner who must choose a matching sequence of actions,  $w^{(1)}, \dots, w^{(T)}$ , from a convex compact set  $Z \subseteq \mathcal{E}$ . Here,  $\mathcal{E}$  and  $\mathcal{E}^*$  respectively refer to an arbitrary Euclidean space and its dual. We focus on a setting where an online learner commits to action  $w^{(t)} \in Z$  before seeing  $g^{(t)}, g^{(t+1)}, \dots$  and aims to achieve vanishing *variational error*  $\text{Err}_V(w^{(1:T)})$  defined by

$$\text{Err}_V(w^{(1:T)}) := \max_{w^* \in Z} \frac{1}{T} \sum_{t=1}^T \langle g^{(t)}, w^{(t)} - w^* \rangle. \quad (3)$$

We will denote no-regret algorithms by their update rule  $\mathcal{Q} : \{Z \times \mathcal{E}^*\} \rightarrow Z$ , where  $\{Z \times \mathcal{E}^*\}$  denotes the space of arbitrary length sequences of action-operator pairs. Given a history sequence  $w^{(1)}, \dots, w^{(t)} \in Z$  and operator sequence  $g^{(1)}, \dots, g^{(t)} \in \mathcal{E}^*$ , the algorithm returns  $w^{(t+1)} = \mathcal{Q}(\{w^{(1)}, g^{(1)}\}, \dots, \{w^{(t)}, g^{(t)}\})$ . When the history is clear from context, we write  $w^{(t+1)} = \mathcal{Q}(w^{(t)}, g^{(t)})$  as shorthand. For the particular case where  $Z = \Delta^n$  is a probability simplex, one such algorithm is Exponential Gradient Descent (also known as Hedge):

$$\mathcal{Q}_{\text{hedge}}\left(\{w^{(1)}, g^{(1)}\}, \dots, \{w^{(t)}, g^{(t)}\}\right) := \frac{\tilde{w}}{\|\tilde{w}\|_1} \text{ where } \tilde{w}_i := w_i^{(t)} \exp\{-\eta g_i^{(t)}\}, \tilde{w} \in \mathbb{R}^n \quad (4)$$

where  $\eta$  is a user-defined step size, and  $w_1$  is a user-defined initial iterate. By default, we take  $w_1 = [\frac{1}{n}]^n$ . The following lemma is a classical result on the variational error of exponential gradient descent.

**Lemma 2.1** ([40]). *Let  $g^{(1)}, \dots, g^{(T)} \in \mathbb{R}^n$  and  $Z = \Delta^n$ . Further assume  $\|g^{(t)}\|_\infty \leq C$  for all timesteps  $t = 1, \dots, T$ . Choosing  $\eta = \sqrt{\log n / T}$ , after  $T$  iterations of exponential gradient descent, the outputs  $w^{(1)}, \dots, w^{(T)}$  satisfies,*

$$\text{Err}_V(w^{(1:T)}) \leq \frac{3C}{2} \sqrt{\frac{KL(w^{(T)} || w^{(1)})}{T}}.$$

### 3 Technical Overview of Our Approach

In this section, we provide an overview of our technical approach for addressing the sample complexity of collaborative learning and group DRO problems. In later sections, we will refer to the approach outlined in this section to sketch proofs and design algorithms. We will focus our exposition on collaborative learning and briefly indicate how the same approach applies to the group DRO setting.

At a high level, we first frame collaborative learning as a zero-sum game with uncertain payoffs and aim to use a variational perspective to learn its minmax equilibrium. We specifically choose the variational perspective (instead of an arbitrary online learning approach), since it allows us to linearize the effect of uncertain payoffs on the resulting error. We then use stochastic gradients to solve the variational problem. Our stochastic gradients will rely on i.i.d. samples from the distributions to estimate gradients both with respect to distributions over  $\mathcal{H}$  and mixtures over  $\mathcal{D}$  but with an asymmetric bound on the bias and variance of the estimates. Along the way, we develop tools and formalisms that handle the asymmetric cost of stochastic gradients and obtain optimal sample complexity results. We now address these steps in more detail.

**Collaborative Learning as Zero-Sum Games.** When the hypothesis class  $\mathcal{H}$  is finite, the collaborative learning problem with distribution set  $\mathcal{D}$  corresponds to a zero-sum game  $(A, A_+, \phi)$  with  $A = \mathcal{H}$ ,  $A_+ = \mathcal{D}$ , and  $\phi(i, j) = \text{Risk}_j(i)$ , where  $i \in A$  and  $j \in A_+$ . Observe that the value of the min-max solution is equivalent to R-OPT. It is not hard to see that any  $\varepsilon$ -min-max equilibrium  $(p, q)$  of this game corresponds to a  $2\varepsilon$  collaborative learning solution, i.e.,

$$\mathbb{E}_{h \sim p} \left[ \max_{D \in \mathcal{D}} \text{Risk}_D(h) \right] \leq \text{OPT} + 2\varepsilon. \quad (5)$$

This enables us to use tools that have been developed for solving zero-sum games in order to address collaborative learning and group DRO settings. We will use a similar construction when hypothesis class  $\mathcal{H}$  has finite VC dimension, where  $A$  will instead refer to an appropriate  $\varepsilon$ -cover of  $\mathcal{H}$ .

**Using the Variational Error to deal with Payoff Uncertainty.** A sufficient condition for minimizing regret, and thus finding  $\varepsilon$ -min-max equilibrium, is minimizing the variational error (Equation 3). In particular, for any finite zero-sum game  $(A, A_+, \phi)$ , defining  $Z = [\Delta A, \Delta A_+]$  and operators

$$g^{(t)} = \left[ \left\{ \partial_{p_i} \phi(p^{(t)}, q^{(t)}) \right\}_{i \in A}, \left\{ -\partial_{q_j} \phi(p^{(t)}, q^{(t)}) \right\}_{j \in A_+} \right], \quad (6)$$

ensures that variational error provides an upper bound on regret:  $\text{Err}_V(w^{(1:T)}) \geq \text{Reg-Min}(p, q) + \text{Reg-Max}(p, q)$ , where  $w = (p, q)$  (see Fact C.1). In collaborative learning, when  $p^{(t)}$  is the min-player's distribution over hypotheses and  $q^{(t)}$  is max-player's distribution over the mixtures, the

gradient vectors refer to the risks of each hypothesis or distribution under  $q^{(t)}$  or  $p^{(t)}$  respectively:

$$g^{(t)} = [g_-^{(t)}, g_+^{(t)}], \quad g_-^{(t)} = \{\text{Risk}_{q^{(t)}}(h)\}_{h \in \mathcal{H}}, \quad g_+^{(t)} = \{\text{Risk}_D(p^{(t)})\}_{D \in \mathcal{D}}. \quad (7)$$

In the collaborative learning setting, we can only create noisy estimates  $\hat{g}$  for these gradients from samples. No-regret algorithms are advantageous in this setting as they choose their  $t$ th iterate  $w^{(t)}$  before seeing the  $t$ th gradient  $g^{(t)}$ . This means that  $w^{(t)}$  is independent of gradient noise,  $\varepsilon^{(t)} := g^{(t)} - \hat{g}^{(t)}$ . We can thus linearize the noise and decompose variational error into the *training* and *generalization* errors as follows

$$\text{Err}_V(w^{(1:T)}) \leq \max_{w^* \in \Delta^n} \frac{1}{T} \sum_{t=1}^T \langle \hat{g}^{(t)}, w^{(t)} - w^* \rangle + \max_{w^* \in \Delta^n} \frac{1}{T} \sum_{t=1}^T \langle \varepsilon^{(t)}, w^{(t)} - w^* \rangle. \quad (8)$$

In contrast, generic no-regret algorithms that do not solve the variational inequality (e.g., when one player plays Hedge and another plays clairvoyant best-response as used in existing work in collaborative learning due to Blum et al. [5], Nguyen and Zakyntinou [29], Chen et al. [8]) nest the generalization and training errors which leads to a multiplicative increase in sample complexity.

**Leveraging Noisy Stochastic Gradients.** We will work with stochastic estimators of  $g$ . These are functions  $\hat{g} : \xi \times \Delta_{-} \times \Delta_{+}$  of some external source of randomness,  $\xi \in \xi$ , and a strategy profile of interest. For collaborative learning, the randomness source  $\xi$  is an i.i.d.-sampled data point from an appropriate mixture of distributions and the estimator  $\hat{g}$  is then the empirical loss on this sample, which is an unbiased and bounded estimator in the range of the loss function, i.e.,  $[0, 1]$ .

Interestingly, estimators of these stochastic gradients have an asymmetric need for data. As seen in Equation 7, the min-player's gradient  $g_-(p, q)$  includes the risk of every hypothesis  $h \in \mathcal{H}$  for the same data distribution  $q$ . Therefore, an unbiased estimator  $\hat{g}_-(p, q)$  can be constructed from a single call to an example oracle  $\text{EX}(q)$ . We call this source of randomness  $\xi^q$  and say that its cost is  $r_- = 1$ . While  $\xi^q$  costs 1 unit, the randomness it provides is specialized to the point of inquiry, that is, it cannot be used for estimating other  $\hat{g}_-(p, q')$ . We call this source of randomness and its associated unbiased estimation a *locally* unbiased estimator.

On the other hand, the max-player's gradient  $g_+(p, q)$  includes the risk of the same hypothesis  $p$  on every distribution  $D \in \mathcal{D}$ . Therefore, an unbiased estimator  $\hat{g}_+(p, q)$  requires  $n$  samples, i.e., a call to every example oracle  $\text{EX}(D_i)$ . We call this source of randomness producing  $n$  samples  $\xi^p$  and say that its cost is  $r_+ = n$ . Importantly, while  $\xi^p$  costs  $n$  unit, the randomness it provides can be reused for estimating other gradients, that is, it can provide unbiased and bounded estimators for all  $\hat{g}_+(p', q')$ . We call this source of randomness and its associated unbiased estimator a *globally* unbiased estimator. To emphasize the fact that this source of randomness is agnostic to  $(p, q)$  we refer to it by  $\xi^\perp$  hereafter. We refer the reader to Appendix B.2 for a more formal definition and description of these asymmetries.

**Minimizing Regret with Asymmetric Cost.** With the goal of minimizing sample complexity in mind, it is essential that we reuse randomness  $\xi^\perp$  across  $n$  time steps of variational algorithms. To do this, we introduce a stochastic variational approach in Algorithm 1 that accommodates different sampling frequencies for the minimizing and maximizing players. This will decouple the sample complexity of the minimizing agent (who requires a time horizon of at least  $\log(A_-) \approx \log(\mathcal{H})$ ) and the maximizing agent. Lemma 3.1 proves this decoupling allows us to find an  $\varepsilon$ -min-max equilibrium with an additive  $n + \log(\mathcal{H})$  sample complexity instead of a multiplicative  $n \log(\mathcal{H})$ .

Algorithm 1 uses the same randomness  $\xi^{\perp(a)}$  of cost  $r$  for estimating  $g_+(p^t, q^t)$  for all  $t \in [ar + 1, \dots, a(r + 1)]$ . On the other hand, the algorithm uses fresh randomness  $\xi^{(t)}$  of cost 1 to estimate  $g_-(p^t, q^t)$  for every time step  $t$ . The total randomness cost of this algorithm is thus  $2T$  because iteration of the outer loop incurs  $2r$  cost.

**Lemma 3.1.** *Let  $(A_-, A_+, \phi)$  be a finite zero-sum game. Assume there exists  $\xi^{q^{(t)}}$  of cost 1 providing locally unbiased estimates  $\hat{g}_-(\cdot)$  and there exists  $\xi^{\perp(a)}$  of cost  $r$  providing globally unbiased estimates  $\hat{g}_+(\cdot)$ . With probability  $1 - \delta$ , Algorithm 1 returns an  $\varepsilon$ -min-max equilibrium of the game, so long as*

$$T \geq \frac{18}{\varepsilon^2} \left( \max \left\{ \frac{9 \log |A_-|}{4}, 8 \log \left( \frac{r + 1}{\delta} \right) \right\} + \max \left\{ \frac{9 \log |A_+|}{4}, \frac{8r^2}{r + 1} \log \left( \frac{r + 1}{\delta} \right) \right\} \right). \quad (9)$$

*Moreover, the total cost of randomness incurred by the algorithm is at most  $2T$ .*

---

**Algorithm 1** Finding Equilibria in Finite Zero-Sum Games with Asymmetric Costs.

---

**Output:** Mixed strategy profile  $(p, q) \in \Delta A_- \times \Delta A_+$ ;  
**Input:** Action sets  $A_-$ ,  $A_+$ , cost  $r \in \mathbb{Z}_+$ , timesteps  $T$ , iterates  $p^{(1)}, q^{(1)}$ , gradient estimators  $\hat{g}_-, \hat{g}_+$ ;  
**for**  $a = 1, 2, \dots, \lceil T/r \rceil$  **do**  
    Realize  $\xi^{\perp(a)}$  at cost  $r$ ;                      // Sample datapoints from every distribution.  
    **for**  $t = ar + 1 - r, \dots, ar$  **do**  
        Realize  $\xi^{q^{(t)}}$  at cost 1;                      // Sample from adversary-selected distribution.  
        Estimate gradients:  $\hat{g}_+^{(t)} = \hat{g}_+ \left( \xi^{\perp(a)}, p^{(t)}, q^{(t)} \right)$ ,  $\hat{g}_-^{(t)} = \hat{g}_- \left( \xi^{q^{(t)}}, p^{(t)}, q^{(t)} \right)$ ;  
        Run Hedge updates:  $p^{(t+1)} = \mathcal{Q}_{\text{hedge}} \left( p^{(t)}, \hat{g}_+^{(t)} \right)$ ,  $q^{(t+1)} = \mathcal{Q}_{\text{hedge}} \left( q^{(t)}, \hat{g}_-^{(t)} \right)$ ;  
    **end for**  
**end for**  
Return the uniformly mixed strategies  $\bar{p} = \frac{1}{T} \sum_{t=1}^T p^{(t)}$  and  $\bar{q} = \frac{1}{T} \sum_{t=1}^T q^{(t)}$ ;

---

*Proof sketch.* Our approach uses Equation 8 to decompose the variational error into training error and generalization error. Since exponential gradient descent is known to bound the training error (as shown in Lemma C.4), it only remains to bound the generalization error (the second term in Equation 3). We note that in expectation each summand  $\langle \varepsilon^{(t)}, w^{(t)} - w^* \rangle$  is zero. This is because  $\varepsilon^{(t)} = g^{(t)} - \hat{g}^{(t)}$  and  $\hat{g}^{(t)}$  are unbiased estimators. Therefore, the sum of these terms has an intuitive martingale interpretation and could be bounded by the Azuma-Hoeffding inequality.

There is a subtlety here, however. When we reuse the maximizing player's randomness over  $r$  rounds, we create correlations between these terms in the generalization error that cannot be directly accommodated by a martingale. The trick here is to note that these correlations are entirely contained in  $r$ -length periods. So, we can partition our sequence to  $r$  martingales and bound each one. This completes the proof. See Appendix C.1 for detailed proof of this lemma.  $\square$

**Derandomization.** The  $\varepsilon$ -min-max equilibria  $(\bar{p}, \bar{q})$  returned by Exponentiated Gradient Descent gives a probability distribution  $\bar{p}$  over the hypothesis class  $\mathcal{H}$  that achieves the collaborative learning bound. To obtain a deterministic hypothesis, we can instead work with  $h_p^{Maj}$  whose predictions are  $p$ -weighted majority votes over the hypotheses in  $\mathcal{H}$ . As stated below, the error of this deterministic classifier is approximately bounded by the expected error of  $\bar{p}$ .

**Lemma 3.2.** For any  $p \in \Delta \mathcal{H}$ ,  $\max_{D \in \mathcal{D}} \text{Risk}_D(h_p^{Maj}) \leq 2 \max_{D \in \mathcal{D}} \text{Risk}_D(p)$ .

This lemma in particular implies that for any  $\varepsilon$ -min-max equilibria  $(\bar{p}, \bar{q})$ , we have

$$\max_{D \in \mathcal{D}} \text{Risk}_D(h_{\bar{p}}^{Maj}) \leq 2\text{R-OPT} + 4\varepsilon \leq 2\text{OPT} + 4\varepsilon.$$

## 4 Collaborative Learning Bounds

In this section, we characterize the sample complexity of collaborative learning by providing tight upper and lower bounds for this problem. We describe Algorithm 2, which attains near-optimal sample complexity by on-demand sampling: iteratively selecting distributions to sample from.

### 4.1 Sample Complexity Upper Bounds

We are now prepared to describe our collaborative learning algorithm and guarantees, using the tools we developed in Section 3. Algorithm 2 is a direct application of Algorithm 1 to a zero-sum game with action sets  $A_- = \mathcal{H}$ ,  $A_+ = \mathcal{D}$  and payoff  $\phi(h, D) = \text{Risk}_D(h)$ . Here,  $\xi^{q^{(t)}}$  makes one call to  $\text{EX}(q^{(t)})$  and  $\xi^{\perp(a)}$  makes one call to  $\text{EX}(D)$  for each  $D \in \mathcal{D}$ . In other words, Algorithm 2 constructs distributions  $p^{(t)} \in \Delta \mathcal{H}$  and  $q^{(t)} \in \Delta \mathcal{D}$  by running the Hedge update. The gradient estimators used by Hedge are the empirical losses on a set of independent random variables. In particular, the minimizing player uses gradients  $\ell_D(h, z^{(t)})$  for all  $h \in \mathcal{H}$  for a single sample  $z^{(t)} \sim \text{EX}(D)$  with  $D \sim q^{(t)}$  and the maximizing player uses gradients  $\ell_D(p^{(t)}, z_D^a)$  for all distributions  $D \in \mathcal{D}$

---

**Algorithm 2** On-Demand Agnostic Collaborative Learning.

---

**Input:** Hypothesis class  $\mathcal{H}$ , distribution set  $\mathcal{D}$  with  $n := |\mathcal{D}|$ ;  
**Initialize:**  $p^{(1)} = [1/|\mathcal{H}|]^{|\mathcal{H}|}$ ,  $q^{(1)} = [1/n]^n$ , and iterations  $T = \frac{36}{\varepsilon^2} (9 \log(|\mathcal{H}|) + 35n \log(n/\delta))$ ;  
**for**  $a = 1, 2, \dots, \lceil T/n \rceil$  **do**  
  For all  $D \in \mathcal{D}$ , sample datapoint  $z_D^a$  from  $\text{EX}(D)$ .  
  **for**  $t = an + 1 - n, \dots, an$  **do**  
    Sample  $z^{(t)}$  from  $\text{EX}(D)$  with  $D \sim q^{(t)}$  and estimate  $\hat{g}^{(t)} = [\ell_D(h, z^{(t)})]_{h \in \mathcal{H}}$ ,  $\hat{g}_+^{(t)} = [\ell_D(p^{(t)}, z_D^a)]_{D \in \mathcal{D}}$ ;  
    Run Hedge updates:  $p^{(t+1)} = \mathcal{Q}_{\text{hedge}}(p^{(t)}, \hat{g}^{(t)})$ ,  $q^{(t+1)} = \mathcal{Q}_{\text{hedge}}(q^{(t)}, \hat{g}_+^{(t)})$ ;  
  **end for**  
**end for**  
**Return:** probability distribution over  $\mathcal{H}$  given by the uniform mixture  $\frac{1}{T} \sum_{t=1}^T p^{(t)}$ .

---

where a single sample  $z_D^a \sim \text{EX}(D)$  is drawn per distribution and is reused for all time steps  $t \in [(a-1)n+1, \dots, an]$ .

Our main result in this section bounds the sample complexity of Algorithm 2.

**Theorem 4.1.** *For any finite hypothesis class  $\mathcal{H}$  and unknown set of distributions  $\mathcal{D}$ , with probability  $1 - \delta$ , Algorithm 2 returns a distribution  $\bar{p} \in \Delta\mathcal{H}$  such that*

$$\mathbb{E}_{h \sim \bar{p}} \left[ \max_{D \in \mathcal{D}} \text{Risk}_D(h) \right] \leq \text{OPT} + \varepsilon \quad \text{and} \quad \max_{D \in \mathcal{D}} \text{Risk}_D(h_{\bar{p}}^{\text{Maj}}) \leq 2\text{OPT} + \varepsilon,$$

using a number of samples that is  $\mathcal{O}\left(\frac{\log|\mathcal{H}| + n \log(n/\delta)}{\varepsilon^2}\right)$ .

*Proof sketch.* By construction, Lemma 3.1 guarantees that with probability at least  $1 - \delta$ , the pair  $(\bar{p}, \bar{q})$  is an  $\varepsilon/2$ -min-max equilibrium for the corresponding zero-sum game. As shown by Equation 5,  $\bar{p}$  is a randomized classifier that meets the collaborative learning objective, i.e., its expected worst-case error is  $\text{OPT} + \varepsilon$ . By Lemma 3.2, the corresponding deterministic classifier  $h_{\bar{p}}^{\text{Maj}}$  has worst-case error of  $2\text{OPT} + \varepsilon$ . This bounds the error of the resulting classifier.

To bound the sample complexity, Lemma 3.1 shows that the randomness cost of Algorithm 1 is at most  $2t$ . Since the cost of randomness is exactly the total number of samples we take from our example oracles, the total sample complexity of Algorithm 2 is  $2t \in \mathcal{O}(\varepsilon^{-2} (\log|\mathcal{H}| + n \log(n/\delta)))$ .  $\square$

An analogue of Theorem 4.1 (Theorem C.3) holds for the case of infinite hypothesis classes of bounded Littlestone dimension with a sample complexity of  $\mathcal{O}(\varepsilon^{-2} (\text{Little}(\mathcal{H}) + n \log(n/\delta)))$ . A similar result also holds with dependence on the VC dimension of  $\mathcal{H}$  only (which is smaller than its Littlestone dimension) when additional assumptions hold. For example, if a hypothesis class  $\mathcal{H}'$  is known in advance that is an  $\varepsilon$ -net of  $\mathcal{H}$  with respect to every distribution in  $\mathcal{D}$ , one can instead run Algorithm 2 with a hypothesis class  $\mathcal{H}'$ . Such an  $\varepsilon$ -net of size  $n\varepsilon^{-\mathcal{O}(\text{VCD}(\mathcal{H}))}$  necessarily exists; for example, the union of  $\varepsilon$ -nets with respect to each distribution  $D \in \mathcal{D}$ . It is also not strictly necessary to know an  $\varepsilon$ -net in advance. Instead, one can compute a net from samples or from other information about distributions in  $\mathcal{D}$ . In Appendix C.5, we explore a range of assumptions that allow us to compute such an  $\varepsilon$ -net from samples, without incurring a significant increase in the sample complexity of Theorem 4.1.

We end this subsection with a few remarks about our sample complexity upper bound.

**Remark 4.1.** *One question left open by these results is, for agnostic multi-distribution learning, whether it is possible to achieve sample complexity rates of  $\mathcal{O}(\varepsilon^{-2} (\log(n) \text{VCD}(\mathcal{H}) + n \log(n/\delta)))$  without any additional assumptions or a priori knowledge of an  $\varepsilon$ -net. It also remains open whether the  $\log(n)$  factor in the  $\log(n) \text{VCD}(\mathcal{H})/\varepsilon^2$  term is necessary for some VC classes, as Theorem 4.1 proves it is not necessary for some (e.g., finite) VC classes.*

**Remark 4.2.** *Theorem 4.1 improves over the best-known sample complexity for agnostic collaborative learning by Nguyen and Zakynthinou [29] in two ways, giving an  $\text{OPT} + \varepsilon$  bound for randomized classifiers instead of their  $2\text{OPT} + \varepsilon$  bound, and improving their sample complexity of  $\mathcal{O}(\frac{1}{\varepsilon^5} (\log(n) \log(|\mathcal{H}|) \log(\frac{1}{\varepsilon}) + n \log(\frac{n}{\delta})))$  by a multiplicative factor of  $\frac{1}{\varepsilon^3} \log(n) \log(\frac{1}{\varepsilon})$ .*



**Remark 4.3.** For constants  $\varepsilon$  and  $\delta$ , our sample complexity of  $\mathcal{O}(\log(|\mathcal{H}|) + n \log n)$  appears to violate the lower bound of  $\Omega(\log(|\mathcal{H}|) \log n + n \log \log |\mathcal{H}|)$  due to Chen, Zhang, and Zhou [8]. This discrepancy is due to a small error in the proof of that lower bound, which we have verified in private communications with the authors. In the next subsection, we give lower bounds on the sample complexity of collaborative learning that match our upper bounds.

## 4.2 Sample Complexity Lower Bound

We now provide matching lower bounds for agnostic collaborative learning. Our lower bounds hold for collaborative learning algorithms obtaining error of  $R\text{-OPT} + \varepsilon$ , using a randomized or deterministic hypothesis. We call an algorithm an  $(\varepsilon, \delta)$ -collaborative learning algorithm if for any collaborative instances it attains an error of  $R\text{-OPT} + \varepsilon$  with probability at least  $1 - \delta$ .

**Theorem 4.2.** Take any  $n, d \in \mathbb{Z}_+$ ,  $\varepsilon, \delta \in (0, 1/8)$ , and  $(\varepsilon, \delta)$ -collaborative learning algorithm  $A$ . There exists a collaborative learning problem  $(\mathcal{H}, \mathcal{D})$  with  $|\mathcal{D}| = n$  and  $|\mathcal{H}| = 2^d$ , on which  $A$  takes at least  $\Omega\left(\frac{1}{\varepsilon^2} (\log |\mathcal{H}| + |\mathcal{D}| \log(\min\{|\mathcal{D}|, \log |\mathcal{H}|\} / \delta))\right)$  samples.

*Proof sketch.* We defer the formal proof of this theorem to Appendix C.3 and sketch the main ideas here. Let  $\mathcal{X} = \{1, \dots, d\}$ ,  $\mathcal{Y} = \{+, -\}$ , and  $\mathcal{H}$  be the set of all functions  $\mathcal{X} \rightarrow \mathcal{Y}$ . Our construction combines two sets of hard distributions. Consider the case when  $n = d \cdot \eta$  for some  $\eta \in \mathbb{Z}$ . First, we can reduce to a  $d$ -armed multi-arm bandit exploration problem giving us an  $\Omega(d \log(1/\delta)/\varepsilon^2)$ . Second, we construct  $\eta$  hard instances on  $\eta$  corresponding points. Since the learning algorithms has to solve each problem it has to incur a loss of  $\eta \cdot d \log(d/\delta)/\varepsilon^2$ .  $\square$

## 5 Group DRO and Agnostic Federated Learning

The results we describe in the collaborative learning setting can be generalized to the group DRO setting, and equivalently, agnostic federated learning.

**Theorem 5.1.** Consider a group distributionally robust problem  $(\Theta, \mathcal{D})$  with convex compact unit-diameter parameter space  $\Theta$  of Bregman radius  $D_\Theta$  (Definition B.11), and convex loss  $\ell : \Theta \times \mathcal{Z} \rightarrow [0, C]$ . A variant of Algorithm 2 (in particular Algorithm 4 in Appendix 4.1), returns  $\hat{\theta} \in \Theta$  such that  $\max_{D \in \mathcal{D}} \mathbb{E}_{z \sim D} [\ell(\hat{\theta}, z)] \leq R\text{-OPT} + \varepsilon$ , using a number of samples that is  $\mathcal{O}\left(\frac{D_\Theta C^2 + n C^2 \log(n/\delta)}{\varepsilon^2}\right)$ .

The proof of this lemma is deferred to Appendix 4.1 and is similar to the proof of Theorem 4.1 except that it uses a generalization of Lemma 3.1 for general convex-concave games. This theorem establishes a generalization bound for the problem of group distributionally robust optimization [36] and improves, by a factor of  $n$ , existing sample complexity bounds for agnostic federated learning [27]. This improvement is attained by sampling data on-demand, whereas [27] only chooses a fixed distribution over groups/clients to sample from; this highlights the importance of adapting one's sampling strategy on-the-fly when learning robust models.

Another important question is how fast the training error of stochastic gradient descent converges for the group DRO/AFL settings and was considered by Sagawa et al. [36]. We can transfer our generalization guarantees for on-demand settings into batch settings and achieve the following corollary, which improves on the convergence guarantees of Sagawa et al. [36] by a factor of  $n$ .

**Corollary 5.2.** Under the same assumptions of Theorem 5.1, we give a procedure (see Appendix 4.1) that minimizes GDRO/AFL training error within  $\varepsilon$  of  $R\text{-OPT}$  with probability at least  $1 - \delta$  in fewer samples than  $\mathcal{O}\left(\frac{D_\Theta C^2 + n C^2 \log(n/\delta)}{\varepsilon^2}\right)$ .

## 6 Empirical Analysis of On-Demand Sampling for Group DRO

This section describes experiments where we adapt our on-demand sampling-based multi-distribution learning algorithm for deep learning applications. In particular, we compare our algorithm against the de-facto standard multi-distribution learning algorithm for deep learning, Group DRO (GDRO) [36]. As GDRO is designed for use with offline-collected datasets, to provide an accurate comparison, we modify our algorithm to work on offline datasets (i.e., with no on-demand sample access).

**Resampling Multi-Distribution Learning (R-MDL).** To adapt our multi-distribution learning algorithm, Algorithm 2, for deep learning applications, we replace its hypothesis-selecting no-regret learning algorithm with a minibatch gradient descent algorithm. We can further adapt our algorithm to offline datasets by simulating on-demand sampling on the empirical distributions of datasets. This modified algorithm, R-MDL, is described in full in Algorithm 5.

In contrast, the GDRO algorithm is also minibatch gradient descent but samples minibatches uniformly from all distributions. Datapoints in each minibatch are importance weighted according to their distribution of origin, where a no-regret algorithm adversarially weights each distribution. Though effective, this GDRO method is brittle and requires tricks like unconventionally strong regularization [36]. Our theory of on-demand sampling suggests that R-MDL should mollify this brittleness.

**Experiment Setting** In Table 2, we replicate the Group DRO experiments of Sagawa et al. [36] and compare the standard GDRO algorithm with our R-MDL algorithm (Algorithm 5). We fine-tune Resnet-50 models (convolutional neural networks) [18] and BERT models (transformer-based network) [11] on the image classification datasets Waterbirds [36, 41] and CelebA [23] and the natural language dataset MultiNLI [42] respectively. We train these models in 3 settings: with standard hyperparameters, under strong weight decay ( $\ell_2$ ) regularization, or under early stopping.

		Worst-Group Accuracy			Gap in Avg. vs. Worst-Group Acc.		
		ERM	GDRO	R-MDL	ERM	GDRO	R-MDL
Standard Reg.	Waterbirds	60.0 (1.9)	76.9 (1.7)	<b>86.4 (1.4)</b>	37.3 (1.9)	20.5 (1.7)	<b>8.1 (1.4)</b>
	CelebA	41.1 (3.7)	41.7 (3.7)	<b>88.9 (2.3)</b>	53.7 (3.7)	53 (3.7)	<b>3.4 (2.3)</b>
	MultiNLI	66.3 (1.6)	66.6 (1.6)	<b>70.3 (1.5)</b>	16.2 (1.6)	15.6 (1.6)	<b>4.5 (1.5)</b>
Strong Reg.	Waterbirds	21.3 (1.6)	84.6 (1.4)	<b>89.4 (1.2)</b>	74.4 (1.6)	12 (1.4)	<b>0.4 (1.3)</b>
	CelebA	37.8 (3.6)	86.7 (2.5)	<b>88.8 (2.3)</b>	58 (3.6)	6.8 (2.5)	<b>1.2 (2.3)</b>
Early Stop	Waterbirds	6.7 (1.0)	85.8 (1.4)	<b>87.1 (1.3)</b>	87.1 (1.0)	7.4 (1.4)	<b>5.6 (1.3)</b>
	CelebA	25.0 (3.2)	88.3 (2.4)	<b>90.6 (2.2)</b>	69.6 (3.2)	3.5 (2.4)	<b>0.7 (2.2)</b>
	MultiNLI	66.0 (1.6)	<b>77.7 (1.4)</b>	43.1 (1.7)	16.8 (1.6)	<b>3.7 (1.4)</b>	18.3 (1.7)

Table 2: Worst-group accuracy (our primary performance metric) and the gap between worst-group accuracy and average accuracy, of empirical risk minimization (ERM), Group DRO (GDRO), and our R-MDL algorithm in three experiment settings—standard hyperparameters (Standard Reg.), inflated weight decay regularization (Strong Reg.), and early stopping (Early Stop)—and on three datasets—Waterbirds, CelebA, and MultiNLI. Figures are percentages evaluated on the test split of each dataset, with standard deviation in parentheses. R-MDL consistently outperforms GDRO and performs reliably with or without strong regularization.

**R-MDL consistently outperforms GDRO and ERM.** In every dataset and in almost every setting, R-MDL significantly outperforms GDRO and ERM in worst-group accuracy. In addition, whereas GDRO and ERM have large gaps between worst-group accuracy and average accuracy, R-MDL has almost matching worst-group and average accuracies. This indicates that R-MDL is more effective at prioritizing learning on difficult groups.

**R-MDL is robust to regularization strength.** R-MDL retains high worst-group accuracy even without strong regularization. These results challenge the observation of Sagawa et al. [36] that strong regularization is critical for the performance of Group DRO methods. This suggests that the brittleness of GDRO is due to reweighting rendering the adversary too weak. In contrast, R-MDL provides a robust multi-distribution learning method with significantly less hyperparameter sensitivity.

## 7 Acknowledgments

This work was supported in part by the National Science Foundation under grant CCF-2145898, a C3.AI Digital Transformation Institute grant, and the Mathematical Data Science program of the Office of Naval Research. This work was partially done while Haghtalab and Zhao were visitors at the Simons Institute for the Theory of Computing.

## References

- [1] N. Alon, O. Ben-Eliezer, Y. Dagan, S. Moran, M. Naor, and E. Yogev. Adversarial laws of large numbers and optimal regret in online classification. In *Proceedings of the Annual ACM Symposium on Theory of Computing (STOC)*, pages 447–455. ACM, 2021.
- [2] A. Beck and M. Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Oper. Res. Lett.*, 31(3):167–175, 2003.
- [3] S. Ben-David and R. Schuller. Exploiting task relatedness for multiple task learning. In *Learning theory and kernel machines*, pages 567–580. Springer, 2003.
- [4] A. Ben-Tal, L. El Ghaoui, and A. Nemirovski. *Robust optimization*, volume 28. Princeton university press, 2009.
- [5] A. Blum, N. Haghtalab, A. D. Procaccia, and M. Qiao. Collaborative PAC Learning. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, pages 2392–2401, 2017.
- [6] A. Blum, N. Haghtalab, R. L. Phillips, and H. Shao. One for One, or All for All: Equilibria and Optimality of Collaboration in Federated Learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 139 of *Proceedings of Machine Learning Research*, pages 1005–1014. PMLR, 2021.
- [7] A. Blum, S. Heinecke, and L. Reyzin. Communication-Aware Collaborative Learning. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 35, pages 6786–6793, 2021.
- [8] J. Chen, Q. Zhang, and Y. Zhou. Tight Bounds for Collaborative PAC Learning via Multiplicative Weights. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, pages 3602–3611, 2018.
- [9] C. Daskalakis, A. Deckelbaum, and A. Kim. Near-Optimal No-Regret Algorithms for Zero-Sum Games. In *Proceedings of the Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 235–254. SIAM, 2011.
- [10] C. Daskalakis, M. Fishelson, and N. Golowich. Near-Optimal No-Regret Learning in General Games. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, volume 34, pages 27604–27616. Curran Associates, Inc., 2021.
- [11] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics*, pages 4171–4186, Minneapolis, Minnesota, 2019. Association for Computational Linguistics.
- [12] J. C. Duchi and H. Namkoong. Learning Models with Uniform Performance via Distributionally Robust Optimization. *CoRR*, abs/1810.08750, 2018. arXiv: 1810.08750.
- [13] C. Dwork, M. P. Kim, O. Reingold, G. N. Rothblum, and G. Yona. Outcome indistinguishability. In *Proceedings of the Annual ACM Symposium on Theory of Computing (STOC)*, pages 1095–1108. ACM, 2021.
- [14] A. Ehrenfeucht, D. Haussler, M. J. Kearns, and L. G. Valiant. A General Lower Bound on the Number of Examples Needed for Learning. *Inf. Comput.*, 82(3):247–261, 1989.
- [15] Y. Freund and R. E. Schapire. A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting. *J. Comput. Syst. Sci.*, 55(1):119–139, 1997.
- [16] S. Hart and A. Mas-Colell. A simple adaptive procedure leading to correlated equilibrium. *Econometrica*, 68(5):1127–1150, 2000. Publisher: Wiley Online Library.
- [17] T. B. Hashimoto, M. Srivastava, H. Namkoong, and P. Liang. Fairness Without Demographics in Repeated Loss Minimization. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 80 of *Proceedings of Machine Learning Research*, pages 1934–1943. PMLR, 2018.

- [18] K. He, X. Zhang, S. Ren, and J. Sun. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778. IEEE Computer Society, 2016.
- [19] L. Hu, C. Peale, and O. Reingold. Metric Entropy Duality and the Sample Complexity of Outcome Indistinguishability. In *Proceedings of the Algorithmic Learning Theory*, volume 167 of *Proceedings of Machine Learning Research*, pages 515–552. PMLR, 2022.
- [20] A. Juditsky, A. Nemirovski, and C. Tauvel. Solving variational inequalities with stochastic mirror-prox algorithm. *Stochastic Systems*, 1(1):17–58, 2011. Publisher: INFORMS.
- [21] A. Kar, A. Prakash, M.-Y. Liu, E. Cameracci, J. Yuan, M. Rusiniak, D. Acuna, A. Torralba, and S. Fidler. Meta-Sim: Learning to Generate Synthetic Datasets. In *Proceedings of the International Conference on Computer Vision*, pages 4550–4559. IEEE, 2019.
- [22] R. M. Karp and R. Kleinberg. Noisy binary search and its applications. In *Proceedings of the Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 881–890. SIAM, 2007.
- [23] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep Learning Face Attributes in the Wild. In *Proceedings of the International Conference on Computer Vision*, pages 3730–3738. IEEE Computer Society, 2015.
- [24] Y. Mansour, M. Mohri, and A. Rostamizadeh. Domain Adaptation with Multiple Sources. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, pages 1041–1048. Curran Associates, Inc., 2008.
- [25] S. Marcel and Y. Rodriguez. Torchvision the machine-vision package of torch. In *Proceedings of the International Conference on Multimedia*, pages 1485–1488. ACM, 2010.
- [26] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y. Arcas. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 54 of *Proceedings of Machine Learning Research*, pages 1273–1282. PMLR, 2017.
- [27] M. Mohri, G. Sivek, and A. T. Suresh. Agnostic Federated Learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 97 of *Proceedings of Machine Learning Research*, pages 4615–4625. PMLR, 2019.
- [28] A. S. Nemirovskij and D. B. Yudin. *Problem complexity and method efficiency in optimization*. Wiley-Interscience, 1983.
- [29] H. L. Nguyen and L. Zakynthinou. Improved Algorithms for Collaborative PAC Learning. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, pages 7642–7650, 2018.
- [30] M. Qiao. Do Outliers Ruin Collaboration? In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 80 of *Proceedings of Machine Learning Research*, pages 4177–4184. PMLR, 2018.
- [31] A. Rakhlin and K. Sridharan. Optimization, Learning, and Games with Predictable Sequences. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, pages 3066–3074, 2013.
- [32] V. V. Ramaswamy, S. S. Kim, and O. Russakovsky. Fair attribute classification through latent space de-biasing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9301–9310, 2021.
- [33] H. Robbins and S. Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951. Publisher: JSTOR.
- [34] J. Robinson. An iterative method of solving a game. *Annals of mathematics*, pages 296–301, 1951. Publisher: JSTOR.

- [35] G. N. Rothblum and G. Yona. Multi-group Agnostic PAC Learnability. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 139 of *Proceedings of Machine Learning Research*, pages 9107–9115. PMLR, 2021.
- [36] S. Sagawa, P. W. Koh, T. B. Hashimoto, and P. Liang. Distributionally Robust Neural Networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*. OpenReview.net, 2020.
- [37] S. Sagawa, A. Raghunathan, P. W. Koh, and P. Liang. An Investigation of Why Overparameterization Exacerbates Spurious Correlations. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 119 of *Proceedings of Machine Learning Research*, pages 8346–8356. PMLR, 2020.
- [38] C. J. Tosh and D. Hsu. Simple and near-optimal algorithms for hidden stratification and multi-group learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 162 of *Proceedings of Machine Learning Research*, pages 21633–21657. PMLR, 2022.
- [39] L. G. Valiant. A Theory of the Learnable. In *Proceedings of the Annual ACM Symposium on Theory of Computing (STOC)*, pages 436–445. ACM, 1984.
- [40] N. K. Vishnoi. *Algorithms for Convex Optimization*. Cambridge University Press, 2021.
- [41] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011 dataset. Technical report, California Institute of Technology, 2011. Publisher: California Institute of Technology.
- [42] A. Williams, N. Nangia, and S. Bowman. A Broad-Coverage Challenge Corpus for Sentence Understanding through Inference. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics*, pages 1112–1122, New Orleans, Louisiana, 2018. Association for Computational Linguistics.
- [43] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. Le Scao, S. Gugger, M. Drame, Q. Lhoest, and A. Rush. Transformers: State-of-the-Art Natural Language Processing. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 38–45, Online, 2020. Association for Computational Linguistics.
- [44] S. Zakharov, W. Kehl, and S. Ilic. DeceptionNet: Network-Driven Domain Randomization. In *Proceedings of the International Conference on Computer Vision*, pages 532–541. IEEE, 2019.
- [45] C. Zhang. Information-theoretic lower bounds of PAC sample complexity, 2019.

## 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]
  - (b) Did you describe the limitations of your work? [Yes]
  - (c) Did you discuss any potential negative societal impacts of your work? [N/A]
  - (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]
  - (b) Did you include complete proofs of all theoretical results? [Yes]
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]
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes]
  - (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]
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? [Yes]
  - (b) Did you mention the license of the assets? [Yes]
  - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [Yes]
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [Yes]
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]

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Related Work . . . . .	3
<b>2</b>	<b>Preliminaries</b>	<b>3</b>
2.1	Technical Background . . . . .	4
<b>3</b>	<b>Technical Overview of Our Approach</b>	<b>5</b>
<b>4</b>	<b>Collaborative Learning Bounds</b>	<b>7</b>
4.1	Sample Complexity Upper Bounds . . . . .	7
4.2	Sample Complexity Lower Bound . . . . .	9
<b>5</b>	<b>Group DRO and Agnostic Federated Learning</b>	<b>9</b>
<b>6</b>	<b>Empirical Analysis of On-Demand Sampling for Group DRO</b>	<b>9</b>
<b>7</b>	<b>Acknowledgments</b>	<b>10</b>
<b>A</b>	<b>Extended Related Work</b>	<b>16</b>
<b>B</b>	<b>Full Formulation</b>	<b>18</b>
B.1	Convex-Concave Zero-Sum Game . . . . .	18
B.2	Stochastic Settings . . . . .	18
B.3	Multi-Distribution Learning . . . . .	19
<b>C</b>	<b>Omitted Proofs</b>	<b>21</b>
C.1	Proof of Lemma C.1 (Generalization of Lemma 3.1) . . . . .	21
C.2	Proof of Theorem C.2 (Generalization of Theorems 4.1 and 5.1) . . . . .	26
C.3	Proof of Theorem 4.2 . . . . .	30
C.4	Proof of Lemmas C.4 and C.10 . . . . .	33
C.5	Proof of Theorem C.12 . . . . .	35
<b>D</b>	<b>Experiment Details</b>	<b>38</b>

## A Extended Related Work

There are many lines of work that study multi-distribution learning but which have evolved independently in separate communities.

**Collaborative and Federated Learning** Blum, Haghtalab, Procaccia, and Qiao [5] posed the first fully general description of multi-distribution learning, motivated by the application of *collaborative PAC learning*. The field of *collaborative learning* concerns the learning of a shared machine learning model by multiple *stakeholders* that each desire a model with low error on their own data distribution. The line of work studies on-demand sample complexity bounds for the setting where stakeholders collect data so as to minimize the error of the worst-off stakeholder [5, 29, 8, 7]. This setting, stated in its full generality, yields the multi-distribution learning problem as presented in this paper. Blum et al. [5] established a  $\log(n)$  factor blowup for the realizable case and the best known sample complexity guarantees for the general agnostic setting experiences a factor  $\log(k)/\varepsilon^3$  blowup and is due to Nguyen and Zakynthinou [29]. In comparison, our work establishes a tight additive increase in the sample complexity (which is comparable to  $\log(k)$  multiplicative factor blowup with no dependence on  $\varepsilon$ ). A related line of work concerns the strategic considerations of collaborative learning and seeks incentive-aware mechanisms for collecting data used in collaborative learning [6].

The field of *federated learning* concerns a closely related motivating application where the goal is to learn a model from data dispersed across multiple devices but querying data from each device is expensive [26]. The agnostic federated learning framework of Mohri, Sivek, and Suresh [27] poses (an equivalent of) the multi-distribution learning objective as a fair and intuitive target for federated learning algorithms, and studies it in the offline setting with data-dependent analysis.

**Group Distributionally Robust Optimization (Group DRO)** Multi-distribution learning also arises in distributionally robust optimization [4] under the name of Group DRO, a class of DRO problems where the distributional uncertainty set is finite [17]. Group DRO literature is motivated by applications where these distributions correspond to deployment domains or protected demographics that a machine learning model should avoid spuriously linking to labels [17, 36, 37]. Although Group DRO—like collaborative learning—is mathematically an instance of multi-distribution learning, prior work on group DRO focus on training error convergence in offline settings as they focus on deep learning applications. As we later discuss, theoretical aspects of online multi-distribution learning can translate into actionable insights for Group DRO applications.

**Multi-Group Fairness and Learning Notions** Multi-distribution learning is also related to the fields of multi-group learning [35, 38] and multi-group fairness [13, 19]. These works study offline learning settings with a single distribution  $D$  and implicitly consider distribution  $D_i$  to be the conditional distribution on a subset of the support representing group  $i$ . In these settings, the learner often does not have explicit access to example oracles for distributions  $D_1, \dots, D_n$  and instead uses rejection sampling to collect data from  $D_1, \dots, D_n$ . As a result, they experience a sub-optimal sample complexity blowup with a factor of  $n$ . This blowup may not be obvious upon first glance, as these works provide theoretical guarantees for each group in terms of the number of datapoints from said group. Rothblum and Yona [35], Tosh and Hsu [38] consider a similar problem to multi-distribution learning; by assuming that there exists a hypothesis that is simultaneously  $\varepsilon$ -optimal on every distribution (an assumption not made in our setting), they compare their learned hypothesis against the best hypothesis for each individual distribution.

**Multi-Source Domain Adaptation and Learning** Multi-source domain adaptation, or multi-task, learning is another related line of work that concerns using data from multiple different training distributions to learn some target distribution, under the assumption that the training and target distributions share some task relatedness [3, 24]. Multi-distribution learning can be framed similarly as using a finite set of training distributions to simultaneously learn the convex hull of the training distributions. Interestingly, multi-distribution learning’s requirement of learning the entire convex hull implicitly obviates the task relatedness assumptions of multi-source learning.

**Stochastic game equilibria.** Our technical approach to multi-distribution learning relates to a line of research on using online algorithms to find min-max equilibria using stochastic feedback [34, 15, 31, 9, 10]. Online mirror descent (OMD) is one well-studied family of methods that can find



approximate minima of convex functions, and also approximate min-max equilibria of convex-concave games, using noisy first-order information [33, 28, 16, 2]. However, applying OMD analysis to the saddle-points in multi-distribution learning is a non-trivial affair. While optimization theory typically thinks about noisy zeroth/first-order oracles (e.g., Juditsky et al. [20]), our learning-theoretic analysis uses oracles that provide noisy semi-local information (i.e., datapoints).

## B Full Formulation

In this section, we formally describe our formulations of stochastic convex-concave games and multi-distribution learning problems.

### B.1 Convex-Concave Zero-Sum Game

In this subsection, we give a formal definition of a convex-concave zero-sum game and its min-max equilibria. We also introduce assumptions on these games for efficiently finding saddle-points.

**Definition B.1.** A convex-concave two-player zero-sum game is described by the tuple  $(A_-, A_+, \phi)$ , where  $A_- \subset \mathcal{E}_-$  is a subset of Euclidian space  $\mathcal{E}_-$ ,  $A_+ \subset \mathcal{E}_+$  is a subset of Euclidian space  $\mathcal{E}_+$ , and  $\phi : A_- \times A_+ \rightarrow \mathbb{R}$  is a Lipschitz continuous convex-concave function.

On a convex-concave two-player zero-sum game  $(A_-, A_+, \phi)$ , we can define both exact and approximate notions of min-max equilibria in terms of player regrets.

**Definition B.2.** The minimizing and maximizing player's regrets at a strategy profile  $(p, q) \in A_- \times A_+$  are denoted  $\text{Reg-Min}$ ,  $\text{Reg-Max}$  respectively, and defined as,

$$\text{Reg-Min}(p, q) := \phi(p, q) - \min_{p^* \in A_-} \phi(p^*, q), \quad \text{Reg-Max}(p, q) := \max_{q^* \in A_+} \phi(p, q^*) - \phi(p, q).$$

**Definition B.3.** A strategy profile  $(p, q) \in A_- \times A_+$  is a min-max equilibrium if both players have zero regret:  $\text{Reg-Min}(p, q) = 0$  and  $\text{Reg-Max}(p, q) = 0$ . More weakly,  $(p, q) \in A_- \times A_+$  is an  $\varepsilon$ -min-max equilibrium if both players have at most  $\varepsilon$  regret:  $\text{Reg-Min}(p, q) \leq \varepsilon$  and  $\text{Reg-Max}(p, q) \leq \varepsilon$ .

In this paper, we may also impose the following assumptions on a convex-concave zero-sum game.

**Assumption 1.** The action sets  $A_-$ ,  $A_+$  are compact, convex, and have diameters  $R_-$ ,  $R_+$  respectively:

$$\forall p, p' \in A_- : \|p - p'\| \leq R_-, \quad \forall q, q' \in A_+ : \|q - q'\| \leq R_+.$$

**Assumption 2.** At any  $p, q \in A_- \times A_+$ , the partial subdifferential of the payoff function  $\phi$  is non-empty. Furthermore, every partial subgradient vector has a bounded norm:

$$\|\partial_p \phi(p, q)\|_{\mathcal{E}_-^*} \leq C_-, \quad \|\partial_q \phi(p, q)\|_{\mathcal{E}_+^*} \leq C_+.$$

### B.2 Stochastic Settings

In this subsection, we give a formal definition of an asymmetric stochastic setting for a zero-sum game. Our formulation of stochastic first-order oracles observes the convention of representing all randomness in stochastic oracles—and by extension, in any stochastic optimization process—in terms of an i.i.d. sequence of random variables. One nuance our formulation addresses is how randomness can be re-used by stochastic first-order oracles. We do this by formalizing our stochastic setting in terms of multiple i.i.d. sequences of random variables, where the sequence to which a random variable belongs specifies how randomness corresponding to the random variable can be used.

We begin by introducing the notion of a coupled random variable. In the context of a two-player game, a random variable may be coupled to a minimizing player's strategy profile, a maximizing player's strategy profile, neither or both. Our definition formalizes the notion that a random variable can only be interpreted in the context of the mixed strategy to which it is coupled.

**Definition B.4.** For any  $p \in A_-$ , we define a random variable  $\eta$  to be  $p$ -coupled if its range is a measurable space  $E_p$  defined by  $p$ . Similarly, for any  $q \in A_+$ , we define a random variable  $\eta$  to be  $q$ -coupled if its range is a measurable space defined by  $q$ . A random variable  $\eta$  is  $(p, q)$ -coupled if its range is a measurable space defined by  $(p, q)$ .

For convenience, we will denote  $p$ -coupled random variables with superscript  $\eta^p$  and, similarly,  $q$ -coupled random variables with superscript  $\eta^q$ . Random variables that are not coupled will be denoted by  $\eta^\perp$  when such clarification is necessary.

We will now define stochastic first-order oracles that express their randomness in terms of sequences of i.i.d. coupled random variables.

**Definition B.5.** In a zero-sum two-player game, the minimizing player's randomness source is defined as a set  $\xi_- \subseteq \{\xi_-^q \mid q \in A_+ \cup \{\perp\}\}$ , where  $\xi_-^q := \{\xi_{-,i}^q\}_{i \in \mathbb{Z}}$  is a sequence of i.i.d. random variables all coupled with  $q \in A_+$ . In addition, all random variables in all sequences in  $\xi_-$  are independent.

**Definition B.6.** In a zero-sum two-player game, the maximizing player's randomness source is defined as a set  $\xi_+ \subseteq \{\xi_+^p \mid p \in A_- \cup \{\perp\}\}$ , where  $\xi_+^p := \{\xi_{+,i}^p\}_{i \in \mathbb{Z}}$  is a sequence of i.i.d. random variables all coupled with  $p \in A_-$ . In addition, all random variables in all sequences in  $\xi_+$  are independent.

**Definition B.7.** For any  $q \in A_+$ , consider the function  $\hat{g}_-^q : E^q \times A_- \times A_+ \rightarrow \mathcal{E}_-^*$ . The minimizing player has a locally unbiased first-order oracle if there exists, for all  $q \in A_+$ , a  $\hat{g}_-^q$  such that for all  $p \in A_-$  and  $i \in \mathbb{Z}$ :

$$\mathbb{E}_{\xi_{-,i}^q} [\hat{g}_-^q(\xi_{-,i}^q, p, q)] = \partial_p \phi(p, q).$$

We analogously define locally unbiased oracles for the maximizing player.

When  $q$  is clear from context, we write  $\hat{g}_-^q$  as  $\hat{g}_-$ . We can also define a globally unbiased oracle.

**Definition B.8.** For any  $q \in A_+$ , consider the function  $\hat{g}_-^\perp : A_- \times A_+ \rightarrow \mathcal{E}_-^*$ . The minimizing player has a globally unbiased first-order oracle if there exists  $\hat{g}_-^\perp$  where for all  $q \in A_+$  and  $p \in A_-$  and  $i \in \mathbb{Z}$ :

$$\mathbb{E}_{\xi_{-,i}^\perp} [\hat{g}_-^\perp(\xi_{-,i}^\perp, p, q)] = \partial_p \phi(p, q).$$

We analogously define globally unbiased first-order oracles for the maximizing player.

Finally, we may impose the following norm-bound assumption on the first-order oracles we discuss.

**Assumption 3.** Every globally unbiased first-order oracle has a range with bounded norm:  $\|\hat{g}_-^\perp(\cdot)\|_{\mathcal{E}_-^*} \leq C_-$ ,  $\|\hat{g}_+^\perp(\cdot)\|_{\mathcal{E}_+^*} \leq C_+$ . Furthermore, every locally unbiased first-order oracle also has a range with bounded norm: for all  $p, q \in A_-$ ,  $A_+$ ,  $\|\hat{g}_-^q(\cdot)\|_{\mathcal{E}_-^*} \leq C_-$ ,  $\|\hat{g}_+^p(\cdot)\|_{\mathcal{E}_+^*} \leq C_+$ .

### B.3 Multi-Distribution Learning

In this subsection, we give a formal definition of multi-distribution learning that unifies the problem formulations of collaborative learning [5], agnostic federated learning [27], and group DRO [36]. We further introduce assumptions that characterize two special cases of multi-distribution learning: *convex multi-distribution learning* and *binary classifier multi-distribution learning*.

We begin by reviewing some common definitions from convex optimization.

**Definition B.9.** Let  $Z$  be a convex compact subset of a Euclidian space  $\mathcal{E}$  with norm  $\|\cdot\|$ . A distance generating function on  $Z$  is a function  $\omega : Z \rightarrow \mathbb{R}$ , where:

1.  $\omega$  is continuous and strongly convex, modulus 1, w.r.t to  $\|\cdot\|$  on  $Z$ .
2. There exists a non-empty subset  $Z^\circ \subset Z$  where the subdifferential  $\partial\omega$  is non-empty and  $\partial\omega$  admits a continuous selection on  $Z^\circ$ .

Furthermore, the center of  $Z$  w.r.t.  $\omega$  is defined as  $z^c := \arg \min_{z \in Z^\circ} \omega(z)$ .

**Definition B.10.** The prox function  $V : Z^\circ \times Z \rightarrow \mathbb{R}^+$  associated with a distance generating function  $\omega : Z \rightarrow \mathbb{R}$  is defined as:

$$V(z, u) := \omega(u) - \omega(z) - \langle \omega'(z), u - z \rangle.$$

The prox function is also known as the Bregman divergence.

**Definition B.11.** Given a convex set  $Z$  with a distance generating function  $\omega$  satisfying Definition B.9, the Bregman radius is defined as  $\max_{u \in Z} V(z^c, u) \leq D_Z$ , where  $z^c$  is the center of  $Z$  as defined in Definition B.9.

We state our most general formulation of multi-distribution learning as follows.

**Definition B.12.** Let  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$  be a space of datapoints and  $\mathcal{D} = \{D_i\}_{i=1}^n$  be a finite multi-set of  $n$  joint probability distributions over  $\mathcal{Z}$ . Let  $\Theta$  denote a set of parameters and  $\ell_D : \Theta \times \mathcal{Z} \rightarrow [0, L]$  be a loss function. Then the tuple  $(\Theta, \mathcal{D}, \ell)$  describes a multi-distribution learning problem, w.r.t.  $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$ .

We will use  $\text{Unique}(\mathcal{D})$  to denote a subset of  $\mathcal{D}$  excluding duplicates of distributions which we know a-priori to have a multiplicity of more than one in  $\mathcal{D}$ .

One case of multi-distribution learning we study in this paper is *convex multi-distribution learning*, which includes as special cases the problem formulations of Sagawa et al. [36] and Mohri et al. [27]. *Convex multi-distribution learning* also encompasses the problem formulation of Blum et al. [5] for finite hypothesis spaces, i.e., when  $|\mathcal{H}| < \infty$ .

**Definition B.13.** The tuple  $(\Theta, \mathcal{D}, \ell)$  describes a convex multi-distribution learning problem when  $\Theta$  is convex compact, each  $\ell_D$  is convex in  $\Theta$ , and there exists a distance generating function  $\omega : \Theta \rightarrow \mathbb{R}$  on our parameter space  $\Theta$ .

**Definition B.14.** The diameter of the parameter space  $\Theta$  is an  $R_\Theta > 0$  satisfying:

$$\forall \theta, \theta' \in \Theta : \|\theta - \theta'\| \leq R_\Theta.$$

**Assumption 4.** Given a convex multi-distribution learning problem  $(\Theta, \mathcal{D}, \ell)$ , we assume that, for any datapoint  $z$  in the supports of the distributions in  $\mathcal{D}$  and any  $\theta \in \Theta$ , the partial subgradient of  $\ell(\theta, z)$  w.r.t.  $\theta$  has bounded norm:

$$\|\partial_\theta \ell(\theta, z)\|_{\mathcal{E}^*} \leq C.$$

**Assumption 5.** Given a convex multi-distribution learning problem  $(\Theta, \mathcal{D}, \ell)$ , we assume there exists a distance generating function  $\omega$  where  $\Theta$  has bounded Bregman radius  $D_\Theta$ .

**Remark B.1.** As  $\omega$  is strongly convex modulus 1 by definition, any  $\Theta$  satisfying Assumption 5 has a finite diameter  $R_\Theta \leq 2\sqrt{2D_\Theta}$ .

Another important case of multi-distribution learning is *binary classifier multi-distribution learning*, which includes the problem formulations of Blum et al. [5] both for finite hypothesis spaces ( $|\mathcal{H}| < \infty$ ) and finite VC dimension hypothesis spaces ( $\text{VCD}(\mathcal{H}) < \infty$ ).

**Definition B.15.** The tuple  $(\Theta, \mathcal{D})$  describes a binary classifier multi-distribution learning problem when  $\mathcal{Z} = \mathcal{X} \times \{0, 1\}$ ,  $\Theta$  is the set of probability distributions over a set of binary classification rules  $\mathcal{H} : \mathcal{X} \rightarrow \{0, 1\}$  and  $\ell(\theta, (x, y)) := \mathbb{E}_{h \sim \theta} [1[h(x) \neq y]]$ .

**Remark B.2.** A binary classifier multi-distribution learning problem  $(\Theta, \mathcal{D}, \ell)$  is equivalent to a convex multi-distribution learning problem  $(\Theta, \mathcal{D}, \ell)$  when the support of  $\Theta$  is finite, i.e.,  $\Theta$  is a probability distribution over a finite number of binary classification rules.

Finally, we can define a multi-distribution analogue to probably-approximately-correct learning [39].

**Definition B.16.** An example oracle  $EX(D)$  is an infinite set of i.i.d. samples from a probability distribution  $D$  over datapoints. Colloquially, a “new call” to example oracle  $EX(D)$  refers to realizing a previously unrealized sample in  $EX(D)$ .

**Definition B.17.** A learning algorithm  $A$  is an  $(\varepsilon, \delta)$  multi-distribution learning algorithm for a set of multi-distribution learning problems  $\mathbb{V} := \{(\Theta_i, \mathcal{D}_i, \ell_i)\}_i$  if, given any problem  $(\Theta_i, \mathcal{D}_i, \ell_i) \in \mathbb{V}$ , accessing only the tuple  $(\Theta_i, \ell_i, \{EX(D) \mid D \in \Delta\mathcal{D}\})$ ,  $A$  outputs a parameter  $\theta \in \Theta_i$  that satisfies, with probability at least  $1 - \delta$ :

$$\max_{D \in \mathcal{D}} \text{Risk}_D(\theta) \leq \inf_{\theta^* \in \Theta} \max_{D \in \mathcal{D}} \text{Risk}_D(\theta^*) + \varepsilon.$$

We use  $(\varepsilon, \delta)$ -algorithm as a shorthand for  $(\varepsilon, \delta)$  multi-distribution learning algorithm.

**Definition B.18.** A multi-distribution learning algorithm  $A$  has a sample complexity of  $N$  (or “takes  $N$  samples”) on a set of multi-distribution learning problems  $\mathbb{V} := \{(\Theta_i, \mathcal{D}_i, \ell_i)\}_i$  if  $N$  is the smallest integer such that, given any problem  $V \in \mathbb{V}$ , the event that  $A$  takes more than  $N$  samples is measure-zero. If no such  $N$  exists, we say  $A$  has infinite sample complexity.

## C Omitted Proofs

### C.1 Proof of Lemma C.1 (Generalization of Lemma 3.1)

---

**Algorithm 3** Finding Equilibria in Convex-Concave Zero-Sum Games with Asymmetric Costs.

---

**Output:** Mixed strategy profile  $(p, q) \in A_- \times A_+$ .

**Input:** Action sets  $A_-, A_+$ , cost  $r \in \mathbb{Z}_+$ , timesteps  $T$ , initial actions  $p^{(1)}, q^{(1)}$ , and no-regret learning algorithms  $\mathcal{Q}_- : \{A_- \times \mathcal{E}_+^*\} \rightarrow A_-$ ,  $\mathcal{Q}_+ : \{A_+ \times \mathcal{E}_-^*\} \rightarrow A_+$ .

**for**  $t = 1, 2, \dots, T$  **do**

Realize randomness  $\xi_{-,t}^{q^{(t)}}$  and  $\xi_{+,t}^\perp$ .

Take estimates  $\hat{g}_-^{(t)} := \hat{g}_- \left( \xi_{-,t}^{q^{(t)}}, p^{(t)}, q^{(t)} \right)$  and  $\hat{g}_+^{(t)} := \hat{g}_+ \left( \xi_{+,t}^\perp, p^{(t)}, q^{(t)} \right)$ .

Run the no-regret updates:

$$p^{(t+1)} = \mathcal{Q}_- \left( \left\{ p^{(1)}, \hat{g}_-^{(1)} \right\}, \dots, \left\{ p^{(t)}, \hat{g}_-^{(t)} \right\} \right)$$

$$q^{(t+1)} = \mathcal{Q}_+ \left( \left\{ q^{(1)}, \hat{g}_+^{(1)} \right\}, \dots, \left\{ q^{(t)}, \hat{g}_+^{(t)} \right\} \right)$$

**end for**

Return the uniformly mixed strategies  $\bar{p} = \frac{1}{T} \sum_{t=1}^T p^{(t)}$  and  $\bar{q} = \frac{1}{T} \sum_{t=1}^T q^{(t)}$ .

---

We first recall the following lemma from Section 3.

**Lemma 3.1.** *Let  $(A_-, A_+, \phi)$  be a finite zero-sum game. Assume there exists  $\xi^{q^{(t)}}$  of cost 1 providing locally unbiased estimates  $\hat{g}_-(\cdot)$  and there exists  $\xi^{\perp(a)}$  of cost  $r$  providing globally unbiased estimates  $\hat{g}_+(\cdot)$ . With probability  $1 - \delta$ , Algorithm 1 returns an  $\varepsilon$ -min-max equilibrium of the game, so long as*

$$T \geq \frac{18}{\varepsilon^2} \left( \max \left\{ \frac{9 \log |A_-|}{4}, 8 \log \left( \frac{r+1}{\delta} \right) \right\} + \max \left\{ \frac{9 \log |A_+|}{4}, \frac{8r^2}{r+1} \log \left( \frac{r+1}{\delta} \right) \right\} \right). \quad (9)$$

Moreover, the total cost of randomness incurred by the algorithm is at most  $2T$ .

We will prove a more general result, Lemma C.1, that implies Lemma 3.1 as a special case. Lemma C.1 provides sample complexity upper bounds for Algorithm 3, an algorithm for approximating the saddle-point of a convex-concave game with high-probability. Algorithm 3 is also a generalization of Algorithm 1.

**Lemma C.1** (Generalization of Lemma 3.1). *Let  $(A_-, A_+, \phi)$  be a convex-concave game satisfying Definition B.1 and Assumptions 1 and 2. Suppose the minimizing player has a locally unbiased first-order oracle  $\hat{g}_-$  and the maximizing player has a globally unbiased first-order oracle  $\hat{g}_+$ , with both oracles satisfying Assumptions 3. Take  $\mathcal{Q}_-$  to be any no-regret algorithms with the guarantee that for, any sequence  $g^{(1)}, \dots, g^{(T)} \in \mathcal{E}_+^*$ , if  $\|g^{(i)}\|_{\mathcal{E}_+^*} \leq C$  for all  $i \in [T]$ , the  $\mathcal{Q}_-$ -learned sequence  $w^{(1)}, \dots, w^{(T)}$  satisfies:*

$$\text{Err}_V(p^{(1:T)}) \leq \sqrt{\frac{\gamma_-(T, A_-, C)}{T}}.$$

Take  $\mathcal{Q}_+$  to be any no-regret algorithms with the guarantee that for, any sequence  $g^{(1)}, \dots, g^{(T)} \in \mathcal{E}_+^*$ , if  $\|g^{(i)}\|_{\mathcal{E}_+^*} \leq C$  for all  $i \in [T]$ , the  $\mathcal{Q}_+$ -learned sequence  $w^{(1)}, \dots, w^{(T)}$  satisfies:

$$\text{Err}_V(w^{(1:T)}) \leq \sqrt{\frac{\gamma_+(T, A_+, C)}{T}}.$$

Then, the mixed strategy profile  $(\bar{p}, \bar{q})$  outputted by Algorithm 3 is an  $\varepsilon$ -min-max equilibrium with probability at least  $1 - \delta$  so long as:

$$T \geq \frac{9}{\varepsilon^2} \left( \gamma_-(T, A_-, 2C_-) + 8R_-^2 C_-^2 \log \left( \frac{r+1}{\delta} \right) + \gamma_+(T, A_+, 2C_+) + \frac{8R_+^2 C_+^2 r^2}{r+1} \log \left( \frac{r+1}{\delta} \right) \right). \quad (10)$$

Moreover, exactly  $T$  elements of  $\xi_-$  and  $\lceil T/r \rceil$  elements of  $\xi_+$  (defined in Definitions B.7 and B.8) will be realized. This means that if sampling from  $\xi_-$  incurs a unit cost and sampling from  $\xi_+$  incurs at most  $r$  unit cost, total cost will be at most  $2r\lceil T/r \rceil$ .

Before proving Lemma C.1, we review the following technical results.

First, we note an immediate consequence of working with a convex payoff function.

**Fact C.1.** Let  $\phi : Z \rightarrow \mathbb{R}$  be a convex function on a convex compact domain  $Z$  and  $g^{(t)} = \partial\phi(w^{(t)})$  be a partial subgradient of  $\phi$  at  $w^{(t)}$ . Then, for any  $[w^{(t)}]_{t=1}^T \in Z$ :

$$\phi\left(\sum_{t=1}^T w^{(t)}\right) - \min_{w^* \in Z} \phi(w^*) \leq \text{Err}_V(w^{(1:T)}) := \max_{w^* \in Z} \sum_{t=1}^T \langle g^{(t)}, w^{(t)} - w^* \rangle.$$

*Proof.* Fix any  $w^* \in Z$ . By our choice of  $g$ , we know that

$$\sum_{t=1}^T \langle g^{(t)}, w^{(t)} - w^* \rangle = \sum_{t=1}^T \langle \partial_{w^{(t)}} \phi(w^{(t)}), w^{(t)} - w^* \rangle.$$

By the convexity of  $\phi$ , it follows that

$$\sum_{t=1}^T \langle \partial_{w^{(t)}} \phi(w^{(t)}), w^{(t)} - w^* \rangle \geq \sum_{t=1}^T \phi(w^{(t)}) - \phi(w^*) \geq \phi\left(\sum_{t=1}^T w^{(t)}\right) - \phi(w^*).$$

with equality when  $\phi$  is bilinear.  $\square$

**Fact C.2.** Let  $\phi : A \times A_+ \rightarrow \mathbb{R}$  be a convex-concave function on convex compact domains  $A, A_+$  and define the operators  $g_-^{(t)} := \partial_{p^{(t)}} \phi(p^{(t)}, q^{(t)})$  and  $g_+^{(t)} := \partial_{q^{(t)}} \phi(p^{(t)}, q^{(t)})$ . Given sequences  $p^{(1)}, \dots, p^{(T)} \in A$  and  $q^{(1)}, \dots, q^{(T)} \in A_+$ , their ergodic averages  $p^{(1:T)} := \frac{1}{T} \sum_{t=1}^T p^{(t)}$  and  $q^{(1:T)} := \frac{1}{T} \sum_{t=1}^T q^{(t)}$  constitute an  $\varepsilon$ -equilibrium if  $\text{Err}_V(p^{(1:T)}) \leq \varepsilon$  and  $\text{Err}_V(q^{(1:T)}) \leq \varepsilon$ .

*Proof.* By Fact C.1, when variational errors are bounded as  $\text{Err}_V(p^{(1:T)}) \leq \varepsilon$  and  $\text{Err}_V(q^{(1:T)}) \leq \varepsilon$ , we know player regrets are bounded:  $\text{Reg-Min}(p^{(1:T)}, q^{(1:T)}) \leq \varepsilon$  and  $\text{Reg-Max}(p^{(1:T)}, q^{(1:T)}) \leq \varepsilon$ . This satisfies our Definition B.3 for an  $\varepsilon$ -min-max equilibria.  $\square$

We now claim concentration results for locally unbiased and globally unbiased first-order oracles.

**Fact C.3.** Let  $(A, A_+, \phi)$  be a convex-concave game satisfying Definition B.1 and Assumptions 1 and 2. Without loss of generality, let our player of interest be the minimizing player. Consider a play sequence  $\{p^{(t)}, q^{(t)}\}_{t=1}^T$  with some complementary sequence  $\{y^{(t)}\}_{t=1}^T \in A_-$ . Suppose, at each timestep, the minimizing player uses a random variable  $\hat{g}_-^{(t)}$  to estimate  $g_-^{(t)} := \partial_{p^{(t)}} \phi(p^{(t)}, q^{(t)})$ . If the following assumptions hold:

1. For every  $t \in [T]$ , the subsequences  $\{p^{(\tau)}, q^{(\tau)}, y^{(\tau)}\}_{\tau=1}^t$  is independent of  $\hat{g}_-^{(t)}, \dots, \hat{g}_-^{(T)}$ .
2. All estimates  $\hat{g}_-^{(1)}, \dots, \hat{g}_-^{(T)}$  are independent.
3.  $\hat{g}_-^{(t)}$  is an unbiased estimate of  $g_-^{(t)}$  and additionally satisfies Assumption 3.

We can then bound the error of the stochastic oracle,  $\varepsilon_-^{(t)} := g_-^{(t)} - \hat{g}_-^{(t)}$ , with respect to our play sequence as follows. With probability at least  $1 - \delta$ ,

$$\max_{p^* \in A_-} \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_-^{(t)}, p^{(t)} - y^{(t)} \rangle \leq \sqrt{\frac{8R_-^2 C_-^2 \log\left(\frac{1}{\delta}\right)}{T}}. \quad (11)$$

*Proof.* Define the filtration  $\{\mathcal{F}^{(t)}\}_{t=0}^T$  as the sigma algebra generated by  $\{\hat{g}_-^{(t)}\}_{t=1}^T$ , with  $\mathcal{F}^{(0)}$  being a singleton containing only the superset of our sigma algebra. Observe that  $\varepsilon_-^{(t)}$  is independent of  $\{p^{(\tau)}\}_{\tau=1}^t$  and  $\{y^{(\tau)}\}_{\tau=1}^t$  by assumption. As  $\hat{g}_-(\cdot)$  is unbiased, for any  $t' = 0, \dots, t-1$ :

$$\mathbb{E} \left[ \langle \varepsilon_-^{(t)}, p^{(t)} - y^{(t)} \rangle \mid \mathcal{F}^{(t')} \right] = 0.$$

We can thus construct the Doob martingale:

$$U = \left\{ \mathbb{E} \left[ \sum_{t=1}^T \langle \varepsilon_-^{(t)}, p^{(t)} - y^{(t)} \rangle \mid \mathcal{F}^{(t')} \right], \mathcal{F}^{(t')} \right\}_{t'=0}^T,$$

and bound the difference sequence accordingly. For any  $t' \in [T]$ , we have the deterministic bound:

$$\begin{aligned} & \mathbb{E} \left[ \sum_{t=1}^T \langle \varepsilon_-^{(t)}, p^{(t)} - y^{(t)} \rangle \mid \mathcal{F}^{(t')} \right] - \mathbb{E} \left[ \sum_{t=1}^T \langle \varepsilon_-^{(t)}, p^{(t)} - y^{(t)} \rangle \mid \mathcal{F}^{(t'-1)} \right] \\ &= \left| \langle \varepsilon_-^{(t')}, p^{(t')} - y^{(t')} \rangle \right| \\ &\leq \left\| \varepsilon_-^{(t')} \right\|_* \left\| p^{(t')} - y^{(t')} \right\|, \end{aligned}$$

where the final inequality is Holder's. Since, by Assumption 1, the diameter of our action sets are bounded by  $R_-$ , we know  $\left\| p^{(t')} - y^{(t')} \right\| \leq R_-$ . Invoking Assumptions 2 and 3, we know  $\left\| \varepsilon_-^{(t')} \right\|_* \leq 2C_-$ . By the Azuma-Hoeffding inequality, we can thus bound, for any  $\varepsilon > 0$ ,

$$\Pr \left( \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_-^{(t)}, p^{(t)} - y^{(t)} \rangle \geq \varepsilon \right) \leq \exp \left( -\frac{\varepsilon^2}{8TC_-^2R_-^2} \right).$$

□

**Fact C.4.** Let  $(A_-, A_+, \phi)$  be a convex-concave game satisfying Definition B.1 and Assumptions 1 and 2. Without loss of generality, let our player of interest be the minimizing player. Consider a play sequence  $\{p^{(t)}, q^{(t)}\}_{t=1}^T$ . If the following assumptions hold:

1.  $\hat{g}_-^{(t)}$  is an estimate of  $g_-^{(t)} := \partial_{p^{(t)}} \phi(p^{(t)}, q^{(t)})$  that satisfies Assumption 3.
2. There exists the no-regret algorithm  $\mathcal{Q}_-$  satisfying the assumptions of Lemma C.1.

We can then bound the error of the stochastic oracle,  $\varepsilon_-^{(t)} := g_-^{(t)} - \hat{g}_-^{(t)}$ , with respect to our play sequence as follows. Define  $y^{(t+1)} := \mathcal{Q}_- \left( \{y^{(\tau)}, \varepsilon_-^{(\tau)}\}_{\tau=1}^t \right)$ . With probability at least  $1 - \delta$ ,

$$\max_{p^* \in A_-} \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_-^{(t)}, p^{(t)} - p^* \rangle \leq \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_-^{(t)}, p^{(t)} - y^{(t)} \rangle + \sqrt{\frac{\gamma_-(T, A_-, 2C_-)}{T}}.$$

*Proof.* We can rewrite Equation 11 with respect to a sequence

$$\begin{aligned} \max_{p^* \in A_-} \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_-^{(t)}, p^{(t)} - p^* \rangle &= \max_{p^* \in A_-} \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_-^{(t)}, p^{(t)} - y^{(t)} + y^{(t)} - p^* \rangle \\ &= \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_-^{(t)}, p^{(t)} - y^{(t)} \rangle + \max_{p^* \in A_-} \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_-^{(t)}, y^{(t)} - p^* \rangle \end{aligned}$$

We will first bound the summand:

$$\max_{p^* \in A_-} \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_-^{(t)}, y^{(t)} - p^* \rangle.$$

By definition, our sequence  $y^{(1)}, \dots, y^{(T)}$  is a  $\mathcal{Q}_-$ -learned sequence for the operator errors  $\varepsilon_-^{(1)}, \dots, \varepsilon_-^{(T)}$ . By Assumptions 2 and 3, we enforce that all operators and operator estimates have bounded norm, i.e.,  $\|g_-(\cdot)\|_{\mathcal{E}^*} \leq C_-$  and  $\|\hat{g}_-(\cdot)\|_{\mathcal{E}^*} \leq C_-$ . By triangle inequality, we can bound  $\left\| \varepsilon_-^{(t')} \right\|_* \leq 2C_-$ . Hence, the guarantees of  $\mathcal{Q}_-$  imply:

$$\max_{p^* \in A_-} \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_-^{(t)}, y^{(t)} - p^* \rangle \leq \sqrt{\frac{\gamma_-(T, A_-, 2C_-)}{T}}.$$

□

We now prove our general claim.

*Proof of Lemma C.1.* By Fact C.2, it suffices to prove the variational error bounds for each player:

$$\text{Err}_V(p^{(1:T)}) \leq \varepsilon, \quad \text{Err}_V(q^{(1:T)}) \leq \varepsilon,$$

with respect to the operators,

$$g_-^{(t)} := \partial_{p^{(t)}} \phi(p^{(t)}, q^{(t)}), \quad g_+^{(t)} := \partial_{q^{(t)}} \phi(p^{(t)}, q^{(t)}).$$

In Algorithm 3, we estimate the true operators  $\{g_-^{(t)}\}_{t=1}^T, \{g_+^{(t)}\}_{t=1}^T$  with the stochastic estimates:

$$\hat{g}_-^{(t)} := \hat{g}_- \left( \xi_{-,t}^{q^{(t)}}, p^{(t)}, q^{(t)} \right), \quad \hat{g}_+^{(t)} := \hat{g}_+ \left( \xi_{+,\lceil \frac{t}{r} \rceil}^\perp, p^{(t)}, q^{(t)} \right).$$

Let  $\varepsilon_-^{(t)} := g_-^{(t)} - \hat{g}_-^{(t)}$  and let  $\varepsilon_+^{(t)} := g_+^{(t)} - \hat{g}_+^{(t)}$  denote the difference between our true and estimated operators at each timestep. We can thus divide each variational error into a *training error* and *generalization error* component:

$$\begin{aligned} \text{Err}_V(p^{(1:T)}) &\leq \max_{p^* \in A_-} \frac{1}{T} \sum_{t=1}^T \langle \hat{g}_-^{(t)}, p^{(t)} - p^* \rangle + \max_{p^* \in A_-} \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_-^{(t)}, p^{(t)} - p^* \rangle \\ \text{Err}_V(q^{(1:T)}) &\leq \max_{q^* \in A_+} \frac{1}{T} \sum_{t=1}^T \langle \hat{g}_+^{(t)}, q^{(t)} - q^* \rangle + \max_{q^* \in A_+} \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_+^{(t)}, q^{(t)} - q^* \rangle \end{aligned}$$

We handle the training error first. Recall that  $p^{(1)}, \dots, p^{(T)}$  is a  $\mathcal{Q}_-$ -learned sequence for the operator sequence  $\hat{g}_-^{(1)}, \dots, \hat{g}_-^{(T)}$ . By Assumption 3, we enforce that all operator estimates have bounded norm, i.e.,  $\|\hat{g}_-(\cdot)\|_{\mathcal{E}^*} \leq C_-$ . Hence, the guarantees of  $\mathcal{Q}_-$  imply:

$$\max_{p^* \in A_-} \frac{1}{T} \sum_{t=1}^T \langle \hat{g}_-^{(t)}, p^{(t)} - p^* \rangle \leq \sqrt{\frac{\gamma_-(T, A_-, C_-)}{T}}. \quad (12)$$

Similarly,  $q^{(1)}, \dots, q^{(T)}$  is a  $\mathcal{Q}_+$ -learned algorithms enjoying  $\mathcal{Q}_+$ 's guarantee:

$$\max_{q^* \in A_+} \frac{1}{T} \sum_{t=1}^T \langle \hat{g}_+^{(t)}, q^{(t)} - q^* \rangle \leq \sqrt{\frac{\gamma_+(T, A_+, C_+)}{T}}. \quad (13)$$

We now handle the generalization error. We first consider the minimization player. Observe that, for every  $t \in [T]$ , the play sequence  $\{p^{(\tau)}, q^{(\tau)}\}_{\tau=1}^t$  is measurable by  $\left\{ \xi_{-, \tau}^{q^{(\tau)}}, \xi_{+,\lceil \tau/r \rceil}^\perp \right\}_{\tau=1}^{t-1}$ , which  $\xi_{-,t}^{q^{(t)}}$  is independent of by construction. We can thus invoke Facts C.3 and C.4 to bound, with probability at least  $1 - \delta$ :

$$\max_{p^* \in A_-} \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_-^{(t)}, p^{(t)} - p^* \rangle \leq \sqrt{\frac{8R_-^2 C_-^2 \log(\frac{1}{\delta})}{T}} + \sqrt{\frac{\gamma_-(T, A_-, 2C_-)}{T}}. \quad (14)$$

We now consider the maximization player. First, we invoke Fact C.4 to separate:

$$\max_{q^* \in A_+} \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_+^{(t)}, q^{(t)} - q^* \rangle \leq \frac{1}{T} \sum_{t=1}^T \langle \varepsilon_+^{(t)}, q^{(t)} - y^{(t)} \rangle + \sqrt{\frac{\gamma_+(T, A_+, 2C_+)}{T}}. \quad (15)$$

For notional convenience, let  $i(j) = (j-1)r + i$  denote the  $i$ th timestep of the  $j$ th period. Also let  $m_i := \left| \{i(j)\}_{j=1}^\infty \cup [T] \right|$  denote the number of valid timesteps that can be written as  $i(j)$ . Observe that  $m_i \leq \lceil T/r \rceil$ . Fix a choice of  $i \in [r]$ . Observe that, for every  $k \in [m_i]$ , the play sequence  $\{p^{(i(j))}, q^{(i(j))}\}_{j=1}^k$  is measurable by  $\left\{ \xi_{-, i(j)}^{q^{(i(j))}}, \xi_{+,\lceil j \rceil}^\perp \right\}_{j=1}^{k-1}$ , which  $\xi_{+,\lceil j \rceil}^\perp$  is independent of by construction. We can thus again invoke Fact C.3 to bound, with probability at least  $1 - \delta$ :

$$\max_{q^* \in A_+} \frac{1}{m_i} \sum_{j=1}^{m_i} \langle \varepsilon_+^{(i(j))}, q^{(i(j))} - y^{(i(j))} \rangle \leq \sqrt{\frac{8R_+^2 C_+^2 \log(\frac{r}{\delta})}{m_i}}.$$



Taking a union bound over said Azuma inequality for all  $i \in [r]$ , we have that with probability at least  $1 - \delta$ ,

$$\begin{aligned}
\max_{q^* \in A_+} \sum_{t=1}^T \left\langle \varepsilon_+^{(t)}, q^{(t)} - y^{(t)} \right\rangle &= \sum_{i=1}^r \sum_{j=1}^{m_i} \left\langle \varepsilon_+^{(i(j))}, q^{(i(j))} - y^{(i(j))} \right\rangle \\
&\leq \sum_{i=1}^r \sqrt{8m_i R_+^2 C_+^2 \log(r/\delta)} \\
&\leq \sqrt{8 \frac{r^2}{r-1} T R_+^2 C_+^2 \log(r/\delta)} \tag{16} \\
\text{(optional: assuming } r \geq 1) \quad &\leq \sqrt{8(r+2) T R_+^2 C_+^2 \log(r/\delta)}.
\end{aligned}$$

Gluing together our bounds on training error (Equation 12, Equation 13) and generalization error (Equation 14, Equation 15, Equation 16) with triangle inequalities and union bounds, we have with probability at least  $1 - \delta$ ,

$$\begin{aligned}
\text{Err}_V(p^{(1:T)}) &\leq \sqrt{\frac{\gamma_-(T, A_-, C_-)}{T}} + \sqrt{\frac{\gamma_-(T, A_-, 2C_-)}{T}} + \sqrt{\frac{8R_-^2 C_-^2 \log\left(\frac{r+1}{\delta}\right)}{T}} \\
\text{Err}_V(q^{(1:T)}) &\leq \sqrt{\frac{\gamma_+(T, A_+, C_+)}{T}} + \sqrt{\frac{\gamma_+(T, A_+, 2C_+)}{T}} + \sqrt{\frac{8r^2 R_+^2 C_+^2 \log\left(\frac{r+1}{\delta}\right)}{(r-1)T}}.
\end{aligned}$$

□

## C.2 Proof of Theorem C.2 (Generalization of Theorems 4.1 and 5.1)

**Fact C.5.** Let  $(\Theta, \mathcal{D}, \ell)$  be a multi-distribution learning problem satisfying Definition B.12. Define a corresponding convex-concave game  $(A, A_+, \phi)$  where:

$$A = \Theta, \quad A_+ = \Delta\mathcal{D}, \quad \phi(p, q) = \text{Risk}_q(p).$$

The minimizing player's mixed strategy  $\bar{p}$  in any  $\varepsilon$ -min-max equilibria (Definition B.3) constitutes an  $2\varepsilon$ -error solution to  $(\Theta, \mathcal{D}, \ell)$ .

*Proof.* If  $(\bar{p}, \bar{q})$  is an  $\varepsilon$ -min-max equilibria, the following holds by definition

$$\text{Risk}_{\bar{q}}(\bar{p}) \leq \min_{p \in \Theta} \text{Risk}_{\bar{q}}(p) + \varepsilon \text{ and } \text{Risk}_{\bar{q}}(\bar{p}) \geq \max_{q \in \Delta\mathcal{D}} \text{Risk}_q(\bar{p}) - \varepsilon.$$

Equivalently, by the min-max theorem,

$$\begin{aligned} \max_{q \in \Delta\mathcal{D}} \text{Risk}_q(\bar{p}) - \varepsilon &\leq \min_{p \in \Theta} \text{Risk}_{\bar{q}}(p) + \varepsilon \\ &\leq \min_{p \in \Theta} \max_{q \in \Delta\mathcal{D}} \text{Risk}_q(p) + \varepsilon. \end{aligned}$$

□

---

### Algorithm 4 On-Demand Multi-Distribution Learning.

---

**Input:** Parameter space  $\Theta$  with distance generating function  $\omega$ , distribution set  $\mathcal{D}$  with  $n := |\mathcal{D}|$  and  $n' := |\text{Unique}(\mathcal{D})|$ , and loss function  $\ell : \Theta \times \mathcal{Z} \rightarrow [0, L]$ , all satisfying Definition B.13 and Assumptions 5 and 4;

**Initialize:** minimizing iterate  $p^{(1)} = \Theta^o$  where  $\theta^o$  is as defined in Definition B.9, maximizing iterate  $q^{(1)} = [1/n]^n$ , and iteration cap:

$$T = \frac{36}{\varepsilon^2} \left( 9C^2 D_\Theta + 8R_\Theta^2 C^2 \log \left( \frac{n+1}{\delta} \right) + 32L^2(n' + 2.1) \log \left( \frac{n' + 1}{\delta} \right) \right);$$

**for**  $a = 1, 2, \dots, \lceil T/n' \rceil$  **do**

For all  $D \in \text{Unique}(\mathcal{D})$ , sample datapoint  $z_D^a$  from  $\text{EX}(D)$  ;

**for**  $t = an' + 1 - n', \dots, an'$  **do**

Sample datapoint  $z^{(t)}$  from  $\text{EX}(D)$  with  $D \sim q^{(t)}$ ;

Define the estimates  $\hat{g}_-^{(t)} = \partial_\theta \ell_D(p^{(t)}, z^{(t)})$  and  $\hat{g}_+^{(t)} = [\ell_D(p^{(t)}, z_D^a)]_{D \in \mathcal{D}}$ ;

Update iterates:  $p^{(t+1)} = \mathcal{Q}_{\text{omd}, \omega}(p^{(t)}, \hat{g}_-^{(t)})$ ,  $q^{(t+1)} = \mathcal{Q}_{\text{hedge}}(q^{(t)}, \hat{g}_+^{(t)})$ ;

**end for**

**end for**

**Return:** parameter  $\bar{\theta} := \frac{1}{T} \sum_{t=1}^T p^{(t)} \in \Theta$ .

---

**Theorem C.2** (Generalization of Theorems 4.1 and 5.1). *Algorithm 4 is an  $(\varepsilon, \delta)$  multi-distribution learning algorithm for any convex multi-distribution learning problem  $(\Theta, \mathcal{D}, \ell)$  satisfying Definitions B.12 and B.13 and Assumptions 5 and 4. In other words, Algorithm 4 returns an  $\bar{\theta} \in \Theta$  such that:*

$$\max_{D \in \mathcal{D}} \text{Risk}_D(\bar{\theta}) \leq \inf_{\theta^* \in \Theta} \max_{D \in \mathcal{D}} \text{Risk}_D(\theta^*) + \varepsilon.$$

Furthermore, the sample complexity of Algorithm 4 is in  $\mathcal{O} \left( \frac{D_\Theta C^2 + (R_\Theta^2 C^2 + n' L^2) \log(n'/\delta)}{\varepsilon^2} \right)$  where  $n' = |\text{Unique}(\mathcal{D})|$ .

*Proof.* The sample complexity of Algorithm 4 is immediate from its construction. Every period  $a$ , Algorithm 4 samples  $n'$  datapoints. Every iteration  $t$ , Algorithm 4 samples 1 datapoint. Thus, Algorithm 4 samples  $2n' \lceil T/n' \rceil$  datapoints exactly.

We now prove that Algorithm 4 is an  $(\varepsilon, \delta)$ -learning algorithm for any convex multi-distribution learning problem. We begin by constructing the following convex-concave game  $(A, A_+, \phi)$  where:

$$A = \Theta, \quad A_+ = \Delta\mathcal{D}, \quad \phi(p, q) = \text{Risk}_q(p).$$

We observe that this game satisfies Definition B.1 and Assumptions 1 and 2:

1. Definition B.13 defines  $\text{Risk}_q(\bar{p})$ —and by extension  $\phi(p, q)$ —to be convex in  $p$ .
2. As  $\text{Risk}_q(\bar{p}) := \sum_{D \in \mathcal{D}} q_D \text{Risk}_D(p)$  by definition,  $\phi(p, q)$  is linear and thus also concave in  $q$ .
3. In the 1-1 norm,  $\Delta\mathcal{D}$  satisfies Assumption 1 with diameter  $R_+ = 2$ .
4. Since  $\Theta$  has finite Bregman radius of  $D_\Theta$  by Assumption 5 and  $\omega$  is strongly convex modulus 1 by definition,  $\Theta$  satisfies Assumption 1 with a finite  $R_\Theta \leq 2\sqrt{2D_\Theta}$ .
5. Since  $\partial_q \phi(p, q) = [\text{Risk}_D(p)]_{D \in \mathcal{D}}$  and the range of  $\ell_D$  is  $[0, L]$ , Assumption 2 is satisfied for  $\partial_q \phi(p, q)$  by  $C_+ \leq L$  in the 1-infinity norm.
6.  $\partial_p \phi(p, q)$  satisfies Assumption 2 for some finite  $C_- = C$  directly by Assumption 4.

We now define a stochastic setting for our game. Let the minimizing player's randomness source be given by the sequences  $\xi^q = \{\text{EX}(q)_i\}_{i=1}^\infty$ ; recall that  $\text{EX}(D)_k$  refers to the  $k$ th call to an example oracle for a  $D \in \Delta\mathcal{D}$ . Let the maximizing player's randomness source be given by the sequence  $\xi_+^\perp = \{[\text{EX}(D)_i]_{D \in \mathcal{D}}\}_{i=1}^\infty$ . Next, define the first-order oracle estimators:

$$\hat{g}_-(\xi_{+,i}^q, p, q) = \partial_p \ell(p, \xi_{+,i}^q), \quad \hat{g}_+(\xi_{+,i}^\perp, p, q) = [\ell(p, (\xi_{+,i}^\perp)_D)]_{D \in \mathcal{D}}.$$

We can observe that  $\hat{g}_-$  is a locally unbiased first-order oracle (satisfying Definition B.7) and  $\hat{g}_+$  is a globally unbiased first-order oracle (satisfying Definition B.8), with both  $\hat{g}_-$  and  $\hat{g}_+$  satisfying Assumptions 3.

1. By the unbiasedness of empirical risk estimates,  $\hat{g}_+$  is globally unbiased as returns an empirical risk sample for each  $D \in \mathcal{D}$ . Similarly, by the unbiasedness of empirical risk estimates and linearity of derivatives,  $\hat{g}_-$  is locally unbiased.
2. As the range of loss function  $\ell$  is in  $[0, L]$ , empirical loss is also bounded in  $[0, L]$ ,  $\hat{g}_+$  satisfies Assumption 3 with  $C_+ \leq L$  in the 1-infinity norm.
3. By Assumption 4, empirical partial subgradients are norm-bounded by some finite  $C$ , so  $\hat{g}_-$  satisfies Assumption 3 with some finite  $C_- = C$ .

Finally, we observe that Algorithm 4 is equivalent to instantiating Algorithm 3 on our constructed game  $(A_-, A_+, \phi)$  for our constructed stochastic setting.

We will now rewrite the iteration complexity requirement of Lemma C.1 given by Equation 10 (copied below):

$$T \geq \frac{9}{\varepsilon'^2} \left( \gamma_-(T, A_-, 2C_-) + 8R_-^2 C_-^2 \log \left( \frac{r+1}{\delta} \right) + \gamma_+(T, A_+, 2C_+) + \frac{8R_+^2 C_+^2 r^2}{r+1} \log \left( \frac{r+1}{\delta} \right) \right).$$

In particular, we aim to show that the default iteration setting of Algorithm 4 satisfies it for  $\varepsilon' = \varepsilon/2$ .

By Lemmas C.4 and C.10, we can bound the efficacy of our no-regret algorithms  $\mathcal{Q}_{\text{omd}}, \mathcal{Q}_{\text{hedge}}$  by:

$$\gamma_-(T, A_-, C_-) \leq \frac{9C_-^2 D_\Theta}{4}, \quad \gamma_+(T, A_+, C_+) \leq \frac{9C_+^2 \log n}{4},$$

where  $\gamma_-(T, A_-, C_-)$  and  $\gamma_+(T, A_+, C_+)$  are as defined in Lemma C.1.

Accounting for our previous derivations of  $C_-, C_+, \varepsilon', R_+$ , to satisfy Equation 10, it suffices to set:

$$T \geq \frac{9 \cdot 4}{\varepsilon^2} \left( \frac{36C_-^2 D_\Theta}{4} + 8R_\Theta^2 C_-^2 \log \left( \frac{n+1}{\delta} \right) + \frac{9L^2 \log n}{4} + \frac{8 \cdot L^2 4(n')^2}{n'+1} \log \left( \frac{n'+1}{\delta} \right) \right),$$

or simplified further:

$$T \geq \frac{36}{\varepsilon^2} \left( 9C_-^2 D_\Theta + 8R_\Theta^2 C_-^2 \log \left( \frac{n+1}{\delta} \right) + 32L^2(n'+2.1) \log \left( \frac{n'+1}{\delta} \right) \right).$$

Thus, by Lemma C.1,  $\bar{q} := \frac{1}{T} \sum_{t=1}^T q^{(t)}$  and  $\bar{\theta}$ , the output of Algorithm 4, form an  $\frac{\varepsilon}{2}$ -min-max equilibria of our game  $(A_-, A_+, \phi)$  with probability at least  $1 - \delta$ . The Theorem then follows by Fact C.5.  $\square$

The following theorems, which are restated from the main text, are immediate corollaries of Theorem C.2.

**Theorem 4.1.** *For any finite hypothesis class  $\mathcal{H}$  and unknown set of distributions  $\mathcal{D}$ , with probability  $1 - \delta$ , Algorithm 2 returns a distribution  $\bar{p} \in \Delta\mathcal{H}$  such that*

$$\mathbb{E}_{h \sim \bar{p}} \left[ \max_{D \in \mathcal{D}} \text{Risk}_D(h) \right] \leq \text{OPT} + \varepsilon \quad \text{and} \quad \max_{D \in \mathcal{D}} \text{Risk}_D(h_{\bar{p}}^{\text{Maj}}) \leq 2\text{OPT} + \varepsilon,$$

using a number of samples that is  $\mathcal{O}\left(\frac{\log|\mathcal{H}| + n \log(n/\delta)}{\varepsilon^2}\right)$ .

*Proof.* Observe that this finite multi-distribution learning problem can be re-written as the convex multi-distribution learning problem  $(\Delta\mathcal{H}, \mathcal{D}, \ell)$ . Since  $\Delta\mathcal{H}$  is a probability simplex of dimension  $|\mathcal{H}|$ , we know it is compact, convex, with  $C = 1$  and  $L = 1$  (as the range of  $\ell$  is in  $[0, 1]$ ), and with  $R_\Theta \leq 2$ . We can then directly apply Theorem C.2, observing that Algorithm 2 is equivalent to Algorithm 4 in this setting.  $\square$

**Theorem 5.1.** *Consider a group distributionally robust problem  $(\Theta, \mathcal{D})$  with convex compact unit-diameter parameter space  $\Theta$  of Bregman radius  $D_\Theta$  (Definition B.11), and convex loss  $\ell : \Theta \times \mathcal{Z} \rightarrow [0, C]$ . A variant of Algorithm 2 (in particular Algorithm 4 in Appendix 4.1), returns  $\bar{\theta} \in \Theta$  such that  $\max_{D \in \mathcal{D}} \mathbb{E}_{z \sim D} [\ell(\bar{\theta}, z)] \leq R\text{-OPT} + \varepsilon$ , using a number of samples that is  $\mathcal{O}\left(\frac{D_\Theta C^2 + n C^2 \log(n/\delta)}{\varepsilon^2}\right)$ .*

*Proof.* Similarly to Theorem 4.1, this claim follows immediately from Theorem C.2 for unit diameter  $R_\Theta = 1$  and loss bound  $L = C$ .  $\square$

Corollary 5.2 follows in a similar fashion, running Algorithm 4 on empirical data distributions. The following proposition re-states this formally.

**Proposition C.1** (Generalization of Corollary 5.2). *Let  $(\Theta, \mathcal{D}, \ell)$  be a convex multi-distribution learning problem satisfying Definitions B.12 and B.13 and Assumptions 5 and 4. For every  $D \in \mathcal{D}$ , let  $\mathbf{B}_D \sim D$  be a non-empty batch of i.i.d. datapoint samples. Define  $\mathcal{D}' = \{D'\}_{D \in \mathcal{D}}$ , where  $D'$  is the empirical distribution of  $\mathbf{B}_D$ . It follows that  $(\Theta, \mathcal{D}', \ell)$  also satisfies Definitions B.12 and B.13 and Assumptions 5 and 4 with identical parameters. Thus, Algorithm 4, when applied to  $(\Theta, \mathcal{D}', \ell)$ , with probability at least  $1 - \delta$  returns an  $\bar{\theta}$  with a multi-distribution training error of at most  $\varepsilon$ . Furthermore, the number of iterations—and accordingly partial derivative operations—is in  $\mathcal{O}\left(\frac{D_\Theta C^2 + (R_-^2 C^2 + n L^2) \log(n/\delta)}{\varepsilon^2}\right)$ .*

A similar proof as Theorem C.2 gives an analogous result for binary classification problems with a hypothesis class of finite Littlestone dimension.

**Theorem C.3** (Littlestone Variant of Theorem 4.1). *Let  $\mathcal{Q}_{\text{Little}}$  denote the no-regret algorithm that achieves a regret of  $\Theta(\sqrt{\text{Little}(\mathcal{H})T})$  in any online learning setting with Littlestone dimension  $\text{Little}(\mathcal{H})$ ; such an algorithm exists by Theorem 2.4 in Alon et al. [1]. Replacing  $\mathcal{Q}_{\text{omd}, \omega}$  with  $\mathcal{Q}_{\text{Little}}$  in Algorithm 4 and updating the iteration cap to,*

$$T \geq \frac{36}{\varepsilon^2} \left( C' \text{Little}(\mathcal{H}) + 8R_\Theta^2 C^2 \log\left(\frac{n+1}{\delta}\right) + 32L^2(n+2.1) \log\left(\frac{n+1}{\delta}\right) \right),$$

*yields an  $(\varepsilon, \delta)$  multi-distribution learning algorithm for any binary classification multi-distribution learning problem  $(\mathcal{H}, \mathcal{D}, \ell)$  satisfying Definitions B.12 and B.15. Further assume that the hypothesis set has finite Littlestone dimension  $\text{Little}(\mathcal{H}) < \infty$ . The sample complexity of Algorithm 4 is in  $\mathcal{O}\left(\frac{\text{Little}(\mathcal{H}) + n \log(n/\delta)}{\varepsilon^2}\right)$ .*

*Proof.* The sample complexity of our modified Algorithm 4 remains immediate from its construction. Algorithm 4 samples  $2n\lceil T/n \rceil$  datapoints exactly. Theorem C.2's proof that Algorithm 4 is an  $(\varepsilon, \delta)$ -learning algorithm for multi-distribution learning also continues to hold. However, we must update the iteration complexity requirement of Lemma C.1 given by Equation 10 (copied below):

$$T \geq \frac{9}{\varepsilon^2} \left( \gamma_-(T, A_-, 2C_-) + 8 \log\left(\frac{r+1}{\delta}\right) + \gamma_+(T, A_+, 2C_+) + \frac{8R_+^2 C_+^2 r^2}{r+1} \log\left(\frac{r+1}{\delta}\right) \right).$$

In particular, we aim to show that the default iteration setting of Algorithm 4 satisfies it for  $\varepsilon' = \varepsilon/2$ . By Lemmas C.4 and by assumption on  $\mathcal{Q}_{\text{Little}}$ , we can bound the efficacy of our no-regret algorithms  $\mathcal{Q}_{\text{Little}}$ ,  $\mathcal{Q}_{\text{hedge}}$  by:

$$\gamma_{-}(T, A_{-}, C_{-}) \leq \frac{C' \text{Little}(\mathcal{H})}{4}, \quad \gamma_{+}(T, A_{+}, C_{+}) \leq \frac{9C_{+}^2 \log n}{4},$$

where  $\gamma_{-}(T, A_{-}, C_{-})$  and  $\gamma_{+}(T, A_{+}, C_{+})$  are as defined in Lemma C.1, and  $C'$  is some universal constant.

Accounting for our previous derivations of  $C_{-}, C_{+}, \varepsilon', R_{+}$  for Theorem 4.1, it suffices to set:

$$T \geq \frac{36}{\varepsilon^2} \left( C' \text{Little}(\mathcal{H}) + 8R_{\Theta}^2 C^2 \log \left( \frac{n+1}{\delta} \right) + 32L^2(n+2.1) \log \left( \frac{n+1}{\delta} \right) \right).$$

Thus, by Lemma C.1,  $\bar{q} := \frac{1}{T} \sum_{t=1}^T q^{(t)}$  and  $\bar{\theta}$ , the output of Algorithm 4, form an  $\frac{\varepsilon}{2}$ -min-max equilibria of our game  $(A_{-}, A_{+}, \phi)$  with probability at least  $1 - \delta$ . The Theorem then follows by Fact C.5.  $\square$

### C.3 Proof of Theorem 4.2

We now provide matching lower bounds for collaborative PAC learning.

We first define a notion of expected sample complexity. Take any multi-distribution learning problem  $V = (\Theta, \mathcal{D}, \ell)$ . Recall that, on this problem, the input to any multi-distribution learning algorithm is a random variable of form  $\hat{V} = (\Theta_i, \ell_i, \{\text{EX}(D)\}_{D \in \Delta \mathcal{D}})$ . Also recall that each example oracle  $\text{EX}(D)$  is an infinite sequence of i.i.d. samples from  $D$ . We will let  $X_V$  denote the probability distribution of the random variable tuple  $(\Theta_i, \ell_i, \{\text{EX}(D)\}_{D \in \Delta \mathcal{D}})$ . Further let  $N_A(\hat{V})$  denote the expected sample complexity of  $A$  given inputs  $\hat{V}$ , where expectation is taken over any randomness from the algorithm  $A$  itself. We can now define a general notion of expected sample complexity.

**Definition C.1.** Let  $A$  be a multi-distribution learning algorithm and  $\mathbb{P}$  a probability distribution over a set of multi-distribution learning problems  $\mathbb{V} := \{(\Theta_i, \mathcal{D}_i, \ell_i)\}_i$ . We define the expected sample complexity  $N_A(\mathbb{P})$  as:

$$N_A(\mathbb{P}) = \mathbb{E}_{V \sim \mathbb{P}} \left[ \mathbb{E}_{\hat{V} \sim X_V} [N_A(\hat{V})] \right].$$

The outer expectation is taken over the randomness of the problem selection, the inner expectation is taken over the randomness of datapoints, and  $N_A(\hat{V})$  takes an expectation over the internal randomness of the algorithm  $A$ .

Unless otherwise specified, we will use the shorthand:  $\mathbb{E}_{V \sim \mathbb{P}} [\mathbb{E}_{\hat{V} \sim X_V} [\cdot]] = \mathbb{E}_{\hat{V}} [\cdot]$ . We recall the following theorem from Section 4.2.

**Theorem 4.2.** Take any  $n, d \in \mathbb{Z}_+$ ,  $\varepsilon, \delta \in (0, 1/8)$ , and  $(\varepsilon, \delta)$ -collaborative learning algorithm  $A$ . There exists a collaborative learning problem  $(\mathcal{H}, \mathcal{D})$  with  $|\mathcal{D}| = n$  and  $|\mathcal{H}| = 2^d$ , on which  $A$  takes at least  $\Omega\left(\frac{1}{\varepsilon^2} (\log |\mathcal{H}| + |\mathcal{D}| \log(\min\{|\mathcal{D}|, \log |\mathcal{H}|\} / \delta))\right)$  samples.

We now prove two lemmas, Lemma C.4 and Lemma C.5, that directly imply Theorem 4.2. These lower-bound constructions are fairly generous and allow all distributions to share the exact same feature distribution and all but one distribution to share the exact same label distribution.

**Lemma C.4.** Take any  $n, d \in \mathbb{Z}_+$ ,  $\varepsilon, \delta \in (0, 1/8)$ , and  $(\varepsilon, \delta)$ -collaborative learning algorithm  $A$ . There exists a set of collaborative learning problems  $\mathbb{V}$  on which  $A$  takes at least  $\Omega\left(\frac{\log |\mathcal{H}|}{\varepsilon^2}\right)$  samples and where, for every  $(\mathcal{H}, \mathcal{D}) \in \mathbb{V}$ ,  $|\mathcal{D}| = n$  and  $|\mathcal{H}| = 2^d$ .

*Proof.* This claim follows directly from the lower bound on sample complexity of agnostic probably-approximately-correct (PAC) learning [39]. Let  $(\mathcal{H}, D)$  be an agnostic PAC learning problem. Accordingly define the collaborative learning problem  $(\mathcal{H}, \mathcal{D})$ , where  $\mathcal{D} = \{D\}_{i=1}^n$ . Observe that  $\min_{D' \in \mathcal{D}} \text{Risk}_{D'}(h) = \text{Risk}_D(h)$  for all well-defined choices of  $h$ . Thus, given an algorithm  $A$  that  $(\varepsilon, \delta)$  solves  $(\mathcal{H}, \mathcal{D}')$  with at most  $m$  samples, we can design an algorithm  $B$  that  $(\varepsilon, \delta)$  solves  $(\mathcal{H}, D)$ : run algorithm  $A$  by simulating samples from any  $D' \in \mathcal{D}$  with samples from  $D$  and return the output of  $A$ . We can thus invoke the well-known lower bound of agnostic PAC learning to observe that there exists an agnostic PAC learning problem  $\mathbb{V}'$  such that any  $(\varepsilon, \delta)$ -learning algorithm has a sample complexity of  $\Omega\left(\frac{\log |\mathcal{H}|}{\varepsilon^2}\right)$  [14]; we defer interested readers to Zhang [45] for a constructive proof of the existence of  $\mathbb{V}'$ . Thus, there exists a  $\mathbb{V}$  satisfying the assumptions of Lemma C.4 where any  $(\varepsilon, \delta)$  collaborative learning algorithm has a sample complexity of  $\Omega\left(\frac{\log |\mathcal{H}|}{\varepsilon^2}\right)$ .  $\square$

**Lemma C.5.** Take any  $n, d \in \mathbb{Z}_+$ ,  $\varepsilon, \delta \in (0, 1/8)$ , and  $(\varepsilon, \delta)$ -collaborative learning algorithm  $A$ . There exists a set of collaborative learning problems  $\mathbb{V}$  on which  $A$  takes at least  $\Omega\left(\frac{1}{\varepsilon^2} (|\mathcal{D}| \log(k/\delta))\right)$  samples and where, for every  $(\mathcal{H}, \mathcal{D}) \in \mathbb{V}$ ,  $|\mathcal{D}| = n$  and  $|\mathcal{H}| = 2^d$  with  $k := \min\{n, d\}$ .

*Proof.* We prove this constructively. We begin by defining collaborative learning problem sets  $\mathbb{V}_{u,w}$  for all  $w, \eta \in \mathbb{N}$  with  $u = w \cdot \eta$ . Problems in  $\mathbb{V}_{u,w}$  share a feature space  $\mathcal{X} = \{1, \dots, w\}$ , label space

$\mathcal{Y} = \{+, -\}$ , and hypothesis class  $\mathcal{H} = \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$  (the set of all deterministic binary labeling functions). For every  $i \in [w]$ , define distributions  $D_i^-$  and  $D_i^+$  as:

$$\Pr_{D_i^-}(i, -) = \frac{1}{2} + 2\varepsilon; \quad \Pr_{D_i^-}(i, +) = \frac{1}{2} - 2\varepsilon \text{ and } \Pr_{D_i^+}(i, -) = \frac{1}{2} - 4\varepsilon; \quad \Pr_{D_i^+}(i, +) = \frac{1}{2} + 4\varepsilon.$$

Let  $\mathbb{P}'_i$  be a probability distribution over possible choices of  $\eta$  collaborative learning problems. With  $1/2$  probability,  $\mathbb{P}'_i$  returns a set of  $\eta$  copies of  $D_i^-$ . With  $1/2$  probability,  $\mathbb{P}'_i$  returns a uniformly sampled shuffling of the set of  $\eta - 1$  copies of  $D_i^-$  and one copy of  $D_i^+$ . We then define  $\mathbb{P}_{u,w} = \mathbb{P}'_1 \times \dots \times \mathbb{P}'_w$  with  $\mathbb{V}_{u,w}$  being the support of  $\mathbb{P}_{u,w}$ . Observe that  $\mathbb{P}_{u,w}$  is a distribution over collaborative learning problems each with  $|\mathcal{H}| = 2^w$  and  $u$  distributions. The following claims characterize sample complexity lower bounds on  $\mathbb{P}_{u,w}$ .

**Claim C.1.** Choose  $\varepsilon \in (0, 1/2)$  and  $\delta \in (0, 1)$ . Take any  $(\varepsilon, \delta)$ -learning algorithm  $A$  on  $\mathbb{V}_{\eta,1}$  under  $\mathbb{P}_{1,1}$ , or equivalently, on  $\mathbb{P}'_1$ . Then, the expected sample complexity (see Definition C.1) of  $A$  is at least  $\frac{\eta}{256\varepsilon^2} \log(1/2\delta)$ .

**Claim C.2.** Choose  $\varepsilon \in (0, 1)$  and  $\delta \in (0, 1/8)$ . Suppose there exists an  $(\varepsilon, \delta)$  learning algorithm  $A$  for  $\mathbb{V}_{u,w}$  with an expected sample complexity of  $m$  under  $\mathbb{P}_{u,w}$ . Then there exists an  $(\varepsilon, \frac{8\delta}{7w})$  learning algorithm  $A'$  for  $\mathbb{V}_{1,1}$  under  $\mathbb{P}_{1,1}$  with an expected sample complexity on  $\mathbb{P}_{1,1}$  of  $\frac{8}{7w}m$ .

Since our desired lower bound are trivially weakly monotonic in  $n, d$ , fix a choice of  $\eta, d \in \mathbb{Z}_+$  and  $\varepsilon, \delta \in (0, 1/8)$  with  $n = \eta \cdot d$ . Combining claims C.1 and C.2, we see that any  $(\varepsilon, \delta)$  collaborative learning algorithm  $A$  for  $\mathbb{V}_{n,d}$  has an expected sample complexity on  $\mathbb{P}_{n,d}$  of at least  $m \geq \frac{7n}{2048\varepsilon^2} \log\left(\frac{7d}{16\delta}\right)$ . By the probabilistic method, for at least some collaborative learning problem  $V \in \mathbb{V}_{n,d}$ ,  $A$  must have a sample complexity in  $\Omega\left(\frac{7n}{2048\varepsilon^2} \log\left(\frac{7d}{16\delta}\right)\right)$ .  $\square$

*Proof of Claim C.1.* Consider the following problem: take  $\eta$  identical-looking coins. Under the  $H_0$  hypothesis, all coins are biased tails with probability  $1/2 + 2\varepsilon$ . Under the  $H_i$  hypothesis, the  $i$ th coin is instead biased heads with probability  $1/2 + 4\varepsilon$ . Let  $\Pr$  be a probability distribution on  $H \in \{H_i\}_{i=0}^\eta$  with  $\Pr(H_0) = 1/2$  and  $\Pr(H_1) = \dots = \Pr(H_\eta) = \frac{1}{2\eta}$ . Suppose that  $A$  is an  $(\varepsilon, \delta)$ -algorithm for  $\mathbb{V}_{\eta,1}$  with an expected sample complexity of  $m$  under  $\mathbb{P}'_1$ . Then, we can construct a coin algorithm  $A'$  that, conditioned on any  $H \in \{H_i\}_{i=0}^\eta$ , with probability at least  $1 - \delta$ , can identify whether  $H_0$  is false. Furthermore, under hypothesis distribution  $\Pr$ ,  $A'$  will also have an expected sample complexity of exactly  $m$ .

To see this, have  $A'$  run  $A$  by simulating draws from the  $i$ th distribution by flipping the  $i$ th coin. If all coins are biased tails with probability  $1/2 + 2\varepsilon$ , any  $\varepsilon$ -error hypothesis  $h$  must satisfy  $\Pr(h(1) = +) > 1/2$ . Conversely, if one coin is biased heads, any  $\varepsilon$ -error hypothesis  $h$  must satisfy  $\Pr(h(1) = +) < 1/2$ . It thus suffices to lower bound the sample complexity of  $A'$ , which is reminiscent of the pure-exploration multi-armed bandit problem.

Suppose  $A'$ , conditioned on  $H_0$ , correctly predicts  $H_0$  with probability at least  $1 - \delta$ . Then, suppose  $A'$ , under  $H_0$ , takes no more than  $T_i$  flips from the  $i$ th coin. Let  $p_{i,j_1:j_2}$  be a probability distribution over  $\{0, 1\}$  corresponding to the outcomes of the  $j_1$ st to  $j_2$ nd coin toss by  $A'$  under  $H_i$ . Let  $p_j^*$  be a uniform distribution over  $\{0, 1\}^j$ . We observe that  $p_{i,j:j}$  and  $p_j^*$  are Bernoulli distributions with a parameter within  $4\varepsilon$  of  $1/2$ . A standard information-theoretic result is that, for  $\varepsilon < 1/2$ ,  $\text{KL}(p_{i,j:j}, p_j^*) < 128\varepsilon^2$ ; e.g. see Zhang [45] and Karp and Kleinberg [22]. By the tensorization of KL Divergence,  $\text{KL}(p_{i,1:j}, p_j^*) < 128j\varepsilon^2$ . By Pinsker's inequality, we can bound total variation distance by  $\text{TV}(p_{i,1:j}, p_j^*) \leq 8\varepsilon\sqrt{j}$ . Let  $E$  be the set of outcomes of  $T_i$  flips under which  $A'$  predicts  $H_0$ ; by correctness under  $H_0$ , we have that  $\Pr_{H_0}(E) \geq 1 - \delta$ . Thus, total variation distance implies  $1 - \delta - 8\varepsilon\sqrt{j} < \Pr_{H_i}(E)$ . Since  $\Pr_{H_i}(E) < \delta$ , we have that  $\frac{1}{128\varepsilon^2} (1 - 2\delta)^2 < T_i$ . Thus, if  $A'$  is  $\delta$  accurate under all hypotheses, under  $H_0$ ,  $A'$  must take at least  $\frac{\eta}{128\varepsilon^2} (1 - 2\delta)^2 < \frac{\eta}{128\varepsilon^2} \log(1/2\delta)$  samples from each distribution. Thus, the average sample complexity of  $A'$  under  $\Pr$ —and similarly the sample complexity of  $A$  under  $\mathbb{P}'_1$ , must be at least  $\frac{\eta}{256\varepsilon^2} \log(1/2\delta)$ .  $\square$

*Proof of Claim C.2.* This claim is similar to the lower bounds of Blum et al. [5] and Karp and Kleinberg [22]. We construct  $A'$  as follows. Let  $I_j = [(j - 1)\eta + 1, j\eta]$ . Suppose a problem  $V' \in \mathbb{V}_{\eta,1}$  is drawn.

1.  $A'$  draws a problem  $(\mathcal{H}, \mathcal{D}) \in \mathbb{V}_{u,w}$  and chooses an index  $i \in [w]$  uniformly at random.
2.  $A'$  simulates the algorithm  $A$  on  $(\mathcal{H}, \mathcal{D})$ ; when  $A$  tries to sample a datapoint from distribution  $D_r$  with  $r \in I_i$ , return a sampled datapoint from the  $(r - (i - 1)\eta)$ st data distribution of  $V'$ .
3. When  $A$  terminates and returns a classifier  $h$ ,  $A'$  checks whether, for every  $j \neq i$ :  $\max_{r \in I_j} \text{Risk}_{D_r}(h) < \frac{1}{2}$ . If this condition is satisfied,  $A'$  returns  $h'(1) = h(\ell)$ . If not, we repeat from Step 1. We denote the number of total iterations with  $T$ .

Consider the probability  $p_i$  that, in the third step, for every  $j \neq i$  we have  $\max_{r \in I_j} \text{Risk}_{D_r}(h) < \frac{1}{2}$  but  $\max_{r \in I_i} \text{Risk}_{D_r}(h) \geq \frac{1}{2}$ . Let  $E_t$  denote the event that  $A'$  returns an at least  $\varepsilon$ -error hypothesis after  $t$  iterations of our procedure. Noting that  $E_t$  can only occur if  $A$  failed all  $t - 1$  iterations before and at the  $t$ th iteration, Step 3 fails to catch the bad hypothesis for  $D_i$ . By assumption,  $\delta \geq \sum_{i=1}^w p_i$ . By symmetry of our construction  $\mathbb{V}$  and recalling  $\delta < 1/8$ :

$$\sum_{t=1}^{\infty} \Pr(E_t) \leq \sum_{t=1}^{\infty} \delta^{t-1} \frac{1}{w} \sum_{i=1}^w p_i \leq \sum_{t=1}^{\infty} \delta^t / w \leq \frac{8\delta}{7w}$$

Thus,  $A'$  is an  $(\varepsilon, \frac{8\delta}{7w})$ -algorithm for  $\mathbb{P}_{\eta,1}$ .

We now bound the sample complexity of  $A'$ . Let  $N_{A'}(t)$  denote the number of samples that  $A'$  takes from  $V'$  on the  $t$ th iteration. Note that  $N_{A'}(1), N_{A'}(2), \dots$  are i.i.d. In addition, by the symmetry of  $\mathbb{V}$  and linearity of expectation,  $\mathbb{E}_{V' \in \mathbb{P}_{\eta,1}} [N_{A'}(t)] = m/w$ . Thus, we can write:

$$\mathbb{E}_{V'} \left[ \sum_{t=1}^T N_{A'}(t) \right] = \mathbb{E}_{V'} [T] \mathbb{E}_{V'} [N_{A'}(1)] = \mathbb{E}_{V'} [T] m/w.$$

We can upper bound  $T$  by observing that our procedure only repeats if  $A$  fails. Thus,

$$\mathbb{E}_{V'} [T] = \sum_{t=1}^{\infty} \Pr(T \geq t) \leq \sum_{t=0}^{\infty} \delta^t \leq \frac{1}{1-\delta} \leq \frac{8}{7}.$$

Thus,  $A'$  has an expected sample complexity of at most  $\frac{8m}{7w}$ . □

The following is a more general restatement of Theorem 4.2 in terms of the terminology of Section B. It follows by observing the difficult cases described in Theorem 4.2 constitute challenging cases for both convex multi-distribution learning (Definition B.13) and binary classifier multi-distribution learning (Definition B.15).

**Corollary C.6.** *Take any  $n, m \in \mathbb{N}$  and  $\varepsilon, \delta \in (0, 1/8)$ . There exists a finite set  $\mathbb{V}$  of multi-distribution learning problems where:*

1. Every  $(\Theta, \mathcal{D}, \ell) \in \mathbb{V}$  satisfies Definitions B.13 and B.15, with  $|\mathcal{D}| = n$  and  $D_{\Theta} = \log(m)$ .
2. Every  $(\varepsilon, \delta)$ -algorithm  $A$  has a sample complexity in  $\left( \frac{D_{\Theta} + n \log(\min\{n, D_{\Theta}\}/\delta)}{\varepsilon^2} \right)$ .



#### C.4 Proof of Lemmas C.4 and C.10

For completeness, this section includes standard results on exponentiated gradient descent and mirror descent more generally, proofs of which can be found in [40].

**Lemma C.7. Pinsker's inequality** For any two vectors from the same probability simplex  $w, w' \in \Delta^n$ , we can bound their generalized Kullback-Leibler divergence as,

$$KL(w||w') \geq \frac{1}{2} \|w - w'\|_1^2.$$

**Lemma C.8. Law of Cosines** Define  $x, y, z \in Z$  where  $Z$  is a convex set, and let  $H : Z \rightarrow \mathbb{R}$  be a strictly convex differentiable distance generating function and  $D_H$  the Bregman divergence of  $H$ . Then,

$$\langle \nabla H(y) - \nabla H(z), y - x \rangle = D_H(x, y) + D_H(y, z) - D_H(x, z).$$

**Lemma C.9. Pythagorean theorem for Bregman Divergence** Define  $x, y \in Z^0$  where  $Z^0$  is a closed subset of a convex set  $Z$ ,  $z \in Z$ . Let  $H : Z \rightarrow \mathbb{R}$  be a strictly convex differentiable distance generating function, and let  $D_H$  be the Bregman divergence of  $H$ . If  $y = \arg \min_{u \in Z^0} D_H(u, z)$ , then  $D_H(x, y) + D_H(y, z) \leq D_H(x, z)$ .

We now turn to proving Lemma (restated below), which concerns exponentiated gradient descent with bounded gradients.

**Lemma 2.1** ([40]). Let  $g^{(1)}, \dots, g^{(T)} \in \mathbb{R}^n$  and  $Z = \Delta^n$ . Further assume  $\|g^{(t)}\|_\infty \leq C$  for all timesteps  $t = 1, \dots, T$ . Choosing  $\eta = \sqrt{\log n / T}$ , after  $T$  iterations of exponential gradient descent, the outputs  $w^{(1)}, \dots, w^{(T)}$  satisfies,

$$\text{Err}_V(w^{(1:T)}) \leq \frac{3C}{2} \sqrt{\frac{KL(w^{(T)}||w^{(1)})}{T}}.$$

*Proof.* This proof closely follows that of Theorem 7.5 in [40]. Fix  $t \in 1, \dots, T$ . First, we provide an expression for  $g^{(t)}$  in terms of  $\tilde{w}^{(t+1)}$  and  $w^{(t)}$ , where  $\tilde{w}$  is as defined in Equation 4. For all  $i \in 1, \dots, n$ , we have by Equation 4:

$$\tilde{w}_i^{(t+1)} = w_i^{(t)} \exp(-\eta g_i^{(t)}).$$

Equivalently,

$$g_i^{(t)} = \frac{1}{\eta} \left( \log w_i^{(t)} - \log \tilde{w}_i^{(t+1)} \right).$$

Letting  $H(x) = \sum_{i=1}^n x_i \log x_i - x_i$  denote our distance generating function, generalized negative entropy, we can also write,

$$g_i^{(t)} = \frac{1}{\eta} \left( \nabla H(w^{(t)}) - \nabla H(\tilde{w}^{(t+1)}) \right)_i,$$

where logs are applied coordinate-wise. Defining  $KL(\cdot||\cdot)$  as generalized Kullback–Leibler divergence: the Bregman divergence with respect to our choice of  $H$  as a distance generating function. Since  $Z$  is already closed, by Lemma C.8, for any  $w^* \in Z$ ,

$$\begin{aligned} \left\langle g_i^{(t)}, w^{(t)} - w^* \right\rangle &= \frac{1}{\eta} \left\langle \nabla H(w^{(t)}) - \nabla H(\tilde{w}^{(t+1)}), w^{(t)} - w^* \right\rangle \\ &= \frac{1}{\eta} \left( KL(w^*||w^{(t)}) + KL(w^{(t)}||\tilde{w}^{(t+1)}) - KL(w^*||\tilde{w}^{(t+1)}) \right). \end{aligned}$$

Generalized Pythagorean Theorem, e.g. Theorem 7.7 in [40] gives,

$$KL(w^*||\tilde{w}^{(t+1)}) \geq KL(w^*||w^{(t+1)}) + KL(w^{(t+1)}||\tilde{w}^{(t+1)}).$$

Then we can bound,

$$\begin{aligned}
\eta \sum_{t=1}^T \langle g_i^{(t)}, w^{(t)} - w^* \rangle &= \sum_{t=1}^T \text{KL}(w^* || w^{(t)}) + \text{KL}(w^{(t)} || \tilde{w}^{(t+1)}) - \text{KL}(w^* || \tilde{w}^{(t+1)}) \\
&\stackrel{\text{(By Pythagorean)}}{\leq} \sum_{t=1}^T \text{KL}(w^* || w^{(t)}) + \text{KL}(w^{(t)} || \tilde{w}^{(t+1)}) \\
&\quad - \left( \text{KL}(w^* || w^{(t+1)}) + \text{KL}(w^{(t+1)} || \tilde{w}^{(t+1)}) \right) \\
&= \sum_{t=1}^T \text{KL}(w^* || w^{(t)}) - \text{KL}(w^* || w^{(t+1)}) \\
&\quad + \left( \text{KL}(w^{(t)} || \tilde{w}^{(t+1)}) - \text{KL}(w^{(t+1)} || \tilde{w}^{(t+1)}) \right) \\
&\stackrel{\text{(By telescoping)}}{\leq} \text{KL}(w^* || w^{(0)}) + \sum_{t=1}^T \text{KL}(w^{(t)} || \tilde{w}^{(t+1)}) - \text{KL}(w^{(t+1)} || \tilde{w}^{(t+1)}).
\end{aligned} \tag{17}$$

To bound the second term, we again apply the law of cosines, this time in reverse order, recovering,

$$\text{KL}(w^{(t)} || \tilde{w}^{(t+1)}) - \text{KL}(w^{(t+1)} || \tilde{w}^{(t+1)}) = \eta \langle g^{(t)}, w^{(t)} - w^{(t+1)} \rangle - \text{KL}(w^{(t+1)} || w^{(t)}).$$

As  $w^{(t+1)}, w^{(t)} \in Z$ , by Pinsker's inequality (Lemma C.7),

$$\begin{aligned}
\text{KL}(w^{(t)} || \tilde{w}^{(t+1)}) - \text{KL}(w^{(t+1)} || \tilde{w}^{(t+1)}) &\leq \eta \langle g^{(t)}, w^{(t)} - w^{(t+1)} \rangle - \frac{1}{2} \|w^{(t+1)} - w^{(t)}\|_1^2 \\
&\leq \eta \|g^{(t)}\|_\infty \|w^{(t)} - w^{(t+1)}\|_1 - \frac{1}{2} \|w^{(t+1)} - w^{(t)}\|_1^2 \\
&\leq \eta C \|w^{(t)} - w^{(t+1)}\|_1 - \frac{1}{2} \|w^{(t+1)} - w^{(t)}\|_1^2 \\
&\leq \frac{\eta^2 C^2}{2},
\end{aligned} \tag{18}$$

where the final inequality follows from maximizing the quadratic  $\eta C z - \frac{z^2}{2}$ , attained at  $z = \|w^{(t)} - w^{(t+1)}\|_1 = C\eta$ . The claim follows by plugging Equation 18 into Equation 17.  $\square$

Exponentiated gradient descent is a special case of mirror descent in the Euclidian space  $\mathcal{E} = \mathbb{R}^n$  equipped with an L1-norm  $\|\cdot\|_1$ , over the probability simplex  $Z = \Delta^n$ , and using entropy as a distance generating function. The following lemma generalizes Lemma C.4 to more general Euclidian spaces, choices of convex compact subsets, and strongly-convex distance generating functions. As the proof closely mirrors that of Lemma C.4, we defer interested readers to Beck and Teboulle [2].

**Lemma C.10** (Generalization of Lemma C.4 [2]). *Let  $Z$  be a convex compact subset of a Euclidean space  $\mathcal{E}$  with distance generating function  $\omega$  satisfying Definition B.9. Let  $g^{(1)}, \dots, g^{(T)} \in \mathcal{E}^*$ . Further assume  $\|g^{(t)}\|_{\mathcal{E}^*} \leq C$  for all timesteps  $t = 1, \dots, T$ . Choose step size  $\eta = \sqrt{\frac{D_Z}{T}}$  where  $D_Z$  is the Bregman radius of  $Z$ . After  $T$  iterations of online mirror descent [2], the output  $\{w\}_{t=1}^T$  satisfies,*

$$\text{Err}_V(w^{(1:T)}) \leq \frac{3C}{2} \sqrt{\frac{D_Z}{T}}.$$

## C.5 Proof of Theorem C.12

This section discusses the implications of our results for finite VC dimension problems.

First, we will use  $D_{\mathcal{X}}$  to denote the marginal distribution of a data distribution  $D$ , i.e. the distribution of  $D$  over its feature space. We also introduce the following definitions.

**Definition C.2.** The Renyi divergence  $D_{\alpha}(P||Q)$  between discrete distributions  $P, Q$  is defined by:

$$D_{\alpha}(P||Q) = \frac{1}{\alpha - 1} \log_2 \sum_{x \in \mathcal{X}} P(x) \left( \frac{P(x)}{Q(x)} \right)^{\alpha-1}$$

and between continuous distributions  $P, Q$  as:

$$D_{\alpha}(P||Q) = \frac{1}{\alpha - 1} \log_2 \int_{\mathcal{X}} P(x) \left( \frac{P(x)}{Q(x)} \right)^{\alpha-1} dx.$$

We will write  $d_{\alpha}(P||Q) := 2^{D_{\alpha}(P||Q)}$ .

**Remark C.1.** Denoting the support of  $Q$  as  $\mathcal{X}_Q$ , we can write  $\sup_{x \in \mathcal{X}_Q} \frac{P(x)}{Q(x)} = d_{\infty}(P||Q)$ .

Recall that in Theorem C.2 we describe a multi-distribution learning algorithm (Algorithm 4) with provably tight sample complexity upper bounds for convex multi-distribution learning problems (Definition B.13). We note that there is one class of multi-distribution learning problems, *non-convex finite VC multi-distribution learning*, that has been previously studied by [5, 29, 8] but does not satisfy the assumptions of convex multi-distribution learning. A *non-convex finite VC multi-distribution learning* problem is a binary-classification multi-distribution learning problem (Definition B.15) that satisfies three criteria: the hypothesis space  $\mathcal{H}$  is non-convex, of infinite size, and of finite VC dimension  $\text{VCD}(\mathcal{H}) < \infty$ . [5, 29, 8] provide upper bounds for non-convex finite VC multi-distribution learning that are identical to their upper bounds in Table 1 but replacing  $\log |\mathcal{H}|$  with  $\text{VCD}(\mathcal{H})$ .

In contrast, our Theorem C.2 upper bounds do not directly apply to non-convex finite VC multi-distribution learning. However, a similar result can be obtained by running our Algorithm 4 on a probability simplex  $\Delta \mathcal{H}'$  over some  $\varepsilon$ -covering  $\mathcal{H}'$  of  $\mathcal{H}$ . Such  $\varepsilon$ -nets of size  $n\varepsilon^{-\text{VCD}(\mathcal{H})}$  necessarily exist. For example, given an  $\varepsilon$ -net for each distribution  $D \in \mathcal{D}$ , we may take their union as the covering  $\mathcal{H}'$  and run Algorithm 2. This directly inherits a favorable upper bound from Theorem C.2.

**Corollary C.11** (VC Dimension Corollary of Theorem C.2). *Consider any binary classification multi-distribution learning problem  $(\mathcal{H}, \mathcal{D})$  where the hypothesis set  $\mathcal{H}$  is of finite VC dimension  $d$  and the unknown distribution set is of size  $|\mathcal{D}| = n$ . There is an algorithm that, given an  $\varepsilon$ -net of size  $\text{poly}(\varepsilon^d, \varepsilon, d, n)$  for each distribution, with probability  $1 - \delta$ , returns a distribution  $\bar{p} \in \Delta \mathcal{H}$  with,*

$$\mathbb{E}_{h \sim \bar{p}} \left[ \max_{D \in \mathcal{D}} \text{Risk}_D(h) \right] \leq \text{OPT} + \varepsilon \quad \text{and} \quad \max_{D \in \mathcal{D}} \text{Risk}_D(h_{\bar{p}}^{\text{Maj}}) \leq 2\text{OPT} + \varepsilon,$$

using a number of samples that is  $\mathcal{O}\left(\frac{d \log(dn/\varepsilon) + n \log(n/\delta)}{\varepsilon^2}\right)$ .

It is also not strictly necessary to know an  $\varepsilon$ -net in advance. Instead, one can compute a net from samples or from other information about distributions in  $\mathcal{D}$ . Theorem C.12 formalizes this claim.

**Theorem C.12.** *For any  $\mathcal{H}$  of VC dimension  $d$  and unknown set of distributions  $\mathcal{D}$  for which Assumption 1, 2 or 3 is met, there is an algorithm that, with probability  $1 - \delta$ , returns a distribution  $\bar{p} \in \Delta \mathcal{H}$  with,*

$$\mathbb{E}_{h \sim \bar{p}} \left[ \max_{D \in \mathcal{D}} \text{Risk}_D(h) \right] \leq \text{OPT} + \varepsilon \quad \text{and} \quad \max_{D \in \mathcal{D}} \text{Risk}_D(h_{\bar{p}}^{\text{Maj}}) \leq 2\text{OPT} + \varepsilon,$$

using a number of samples that is  $\mathcal{O}\left(\frac{d \log(dn/\varepsilon) + n \log(n/\delta)}{\varepsilon^2}\right)$ .

The following corollaries of Theorem 4.1 directly imply Theorem C.12.

**Corollary C.13** (Assumption 1). *Consider any binary classification multi-distribution learning problem  $(\mathcal{H}, \mathcal{D})$  where the hypothesis set  $\mathcal{H}$  is of finite VC dimension  $d$  and the unknown distribution set is of size  $|\mathcal{D}| = n$ . For  $\varepsilon \in \mathcal{O}(1/n)$ , there is an algorithm that with probability  $1 - \delta$ , returns a distribution  $\bar{p} \in \Delta\mathcal{H}$  with,*

$$\mathbb{E}_{h \sim \bar{p}} \left[ \max_{D \in \mathcal{D}} \text{Risk}_D(h) \right] \leq OPT + \varepsilon \quad \text{and} \quad \max_{D \in \mathcal{D}} \text{Risk}_D(h_{\bar{p}}^{Maj}) \leq 2OPT + \varepsilon,$$

using a number of samples that is  $\mathcal{O}\left(\frac{d \log(dn/\varepsilon) + n \log(n/\delta)}{\varepsilon^2}\right)$ .

*Proof.* By Lemma C.16, sampling  $\mathcal{O}\left(\frac{nd}{\varepsilon} \log\left(\frac{d}{\varepsilon}\right) + \frac{n}{\varepsilon} \log\left(\frac{n}{\varepsilon}\right)\right)$  datapoints provides a covering of  $\mathcal{H}$  that is simultaneously an  $\varepsilon$ -net for every  $D \in \mathcal{D}$  with probability at least  $1 - \delta$ . Moreover, by the Sauer-Shelah lemma, this net is of size  $\mathcal{O}\left(\left(\frac{\log(dn/\varepsilon) + n \log(n/\delta)}{\varepsilon^2}\right)^d\right)$ . The claim then follows from Corollary C.11, noting that since  $\varepsilon \in \mathcal{O}(1/n)$ , we only needed to sample an additional  $\mathcal{O}\left(\frac{d}{\varepsilon^2} \log\left(\frac{d}{\varepsilon}\right) + \frac{n}{\varepsilon} \log\left(\frac{n}{\varepsilon}\right)\right) \subset \mathcal{O}\left(\frac{nd}{\varepsilon} \log\left(\frac{d}{\varepsilon}\right) + \frac{n}{\varepsilon} \log\left(\frac{n}{\varepsilon}\right)\right)$  datapoints to form the cover.  $\square$

**Corollary C.14** (Assumption 2). *Consider any binary classification multi-distribution learning problem  $(\mathcal{H}, \mathcal{D})$  where the hypothesis set  $\mathcal{H}$  is of finite VC dimension  $d$  and the unknown distribution set is of size  $|\mathcal{D}| = n$ . We say an algorithm has weak unlabeled access if the algorithm can access, for each  $D \in \mathcal{D}$ , a marginal distribution  $D'_{\mathcal{X}}$  such that  $d_{\infty}(D'_{\mathcal{X}} \| D_{\mathcal{X}}) \in \text{poly}(1/\varepsilon, d, n)$ , with probability  $1 - \delta$ . There is an algorithm that, given weak unlabeled access, with probability  $1 - \delta$ , returns a distribution  $\bar{p} \in \Delta\mathcal{H}$  with,*

$$\mathbb{E}_{h \sim \bar{p}} \left[ \max_{D \in \mathcal{D}} \text{Risk}_D(h) \right] \leq OPT + \varepsilon \quad \text{and} \quad \max_{D \in \mathcal{D}} \text{Risk}_D(h_{\bar{p}}^{Maj}) \leq 2OPT + \varepsilon,$$

using a number of samples that is  $\mathcal{O}\left(\frac{d \log(dn/\varepsilon) + n \log(n/\delta)}{\varepsilon^2}\right)$ .

*Proof.* Observe that when  $d_{\infty}(D'_{\mathcal{X}} \| D_{\mathcal{X}}) < \gamma$ ,  $D'_{\mathcal{X}}$  can be written as a mixture over  $D_{\mathcal{X}}$  with probability at least  $\frac{1}{\gamma}$  and some other distribution  $\tilde{D}_{\mathcal{X}}$  with probability at most  $1 - \frac{1}{\gamma}$ . Once again invoking uniform convergence, we observe that sampling  $\Theta\left(d_{\infty}(D'_{\mathcal{X}} \| D_{\mathcal{X}}) \frac{d \log(d/\varepsilon) + \log(1/\delta)}{\varepsilon^2}\right)$  i.i.d. samples from distribution  $D'_{\mathcal{X}}$ , with probability at least  $1 - \delta$ , yields an  $\varepsilon$ -covering on  $D$ . By Sauer-Shelah's lemma, the resulting covering  $\mathcal{H}'_D$  is of size  $\mathcal{O}\left((\text{poly}(1/\varepsilon, d, n))^d\right)$ . Repeating this procedure for each  $D \in \mathcal{D}$ , with probability at least  $1 - n\delta$ , we have an  $\varepsilon$ -covering  $\mathcal{H}'$  of  $\mathcal{D}$  of size  $|\mathcal{H}'| \in \mathcal{O}\left(n(\text{poly}(1/\varepsilon, d, n))^d\right)$ . We can then appeal directly to Theorem 4.1 for a sample complexity bound on learning  $(\mathcal{H}', \mathcal{D})$ .  $\square$

**Corollary C.15** (Assumption 3). *Consider any binary classification multi-distribution learning problem  $(\mathcal{H}, \mathcal{D})$  where the hypothesis set  $\mathcal{H}$  is of finite VC dimension  $d$  and the unknown distribution set is of size  $|\mathcal{D}| = n$ . There is an algorithm that, given access to the marginal distribution  $D_{\mathcal{X}}$  of every  $D \in \mathcal{D}$ , with probability  $1 - \delta$ , returns a distribution  $\bar{p} \in \Delta\mathcal{H}$  with,*

$$\mathbb{E}_{h \sim \bar{p}} \left[ \max_{D \in \mathcal{D}} \text{Risk}_D(h) \right] \leq OPT + \varepsilon \quad \text{and} \quad \max_{D \in \mathcal{D}} \text{Risk}_D(h_{\bar{p}}^{Maj}) \leq 2OPT + \varepsilon,$$

using a number of samples that is  $\mathcal{O}\left(\frac{d \log(1/\varepsilon) + n \log(n/\delta)}{\varepsilon^2}\right)$ .

*Proof.* By uniform convergence, sampling  $\Theta\left(\frac{d \log(d/\varepsilon) + \log(1/\delta)}{\varepsilon^2}\right)$  i.i.d. samples from a distribution  $D_{\mathcal{X}}$ , with probability at least  $1 - \delta$ , yields an  $\varepsilon$ -covering on  $D$ . By Sauer-Shelah's lemma, the resulting covering  $\mathcal{H}'_D$  is of size  $\mathcal{O}\left((\log(d/\varepsilon) + \frac{1}{d} \log(1/\delta)/\varepsilon^2)^d\right)$ . Repeating this procedure for each  $D \in \mathcal{D}$ , with probability at least  $1 - n\delta$ , we have an  $\varepsilon$ -covering of  $\mathcal{D}$  of size  $\mathcal{O}\left(n(\log(d/\varepsilon) + \frac{1}{d} \log(1/\delta)/\varepsilon^2)^d\right) \in \text{poly}(\varepsilon^d, \varepsilon, d, n)$  and can appeal to Corollary C.11.  $\square$

**Lemma C.16** (Corollary 3.7 in Haussler, Welzl). *Let  $\mathcal{F}$  be a function class consisting of functions from  $\mathcal{X}$  to  $[0, 1]$  and  $\mathcal{P}$  a probability measure on  $\mathcal{X}$ . Given  $N \geq \frac{8d}{\varepsilon} \log \frac{8d}{\varepsilon} + \frac{4}{\varepsilon} \log \frac{2}{\delta}$  random samples  $\mathbf{x}$  from  $\mathcal{P}$ , with probability at least  $1 - \delta$ , the projection of  $\mathcal{F}$  on  $\mathbf{x}$  constitutes an  $\varepsilon$ -net. That is, for any  $f_1, f_2 \in \mathcal{F}$  where  $\Pr_{x \sim \mathcal{P}}(f_1(x) \neq f_2(x)) \geq \varepsilon$ ,  $\|f_1(x) - f_2(x)\|_{\mathbf{x}} > 0$ .*

---

**Algorithm 5** Resampling-based Multi-Distribution Learning (R-MDL)

---

**Input:** Parameter space  $\Theta$ , initial parameter  $\theta^{(0)} \in \Theta$ , training set  $\mathcal{X}_D$  and validation set  $\mathcal{X}'_D$  for  $D$  in  $\mathcal{D}$  where  $n := |\mathcal{D}|$ , iterations  $T$ , minibatch size  $B$ , adversary minibatch size  $B'$ , blackbox minibatch learning algorithm  $\mathcal{A}$ ;

**Initialize:**  $q^{(0)} = [1/n]^n$ ;

**for** iterations  $t = 1, 2, \dots, T$  **do**

For each  $D \in \mathcal{D}$ , sample  $B'/n$  points  $\mathcal{X}'^{(t)}_D$  from the empirical distribution of  $\mathcal{X}'_D$ . In other words, sample uniformly with replacement from  $\mathcal{X}'_D$ .

Run Hedge update on adversary using datasets  $[\ell(\theta, \mathcal{X}'^{(t)}_D)]_{D \in \mathcal{D}}$ . In other words, letting  $\exp$  be component-wise exponentiation,  $\cdot$  be component-wise multiplication, and  $\eta$  some learning rate:

$$q^{(t)} = \frac{q^{(t-1)} \cdot \exp(-\eta[\ell(\theta, \mathcal{X}'^{(t)}_D)]_{D \in \mathcal{D}})}{\|q^{(t-1)} \cdot \exp(-\eta[\ell(\theta, \mathcal{X}'^{(t)}_D)]_{D \in \mathcal{D}})\|_1};$$

Sample  $B$  datapoints from the mixture over the empirical distributions of  $\{\mathcal{X}_D\}_{D \in \mathcal{D}}$  weighted by  $q^{(t)}$ . In other words, sample with replacement from  $\mathcal{X}_D$  with probability  $q^{(t)}_D$ .

Run minibatch update on learner, with  $\theta^{(t)} = \mathcal{A}(\theta^{(t-1)}, \mathcal{X}^{(t)})$ .

**end for**

**Return:**  $\theta^{(T)}$ .

---

## D Experiment Details

**R-MDL Algorithm** The R-MDL algorithm is defined in full in Algorithm 5. This algorithm is implemented in the Github repository [ericzhao28/multidistributionlearning](https://github.com/ericzhao28/multidistributionlearning).

**Additional Observation: R-MDL converges faster than ERM or GDRO.** The R-MDL methods in Table 2 used a fraction of the training epochs that their GDRO counterparts used. The ratio of R-MDL to GDRO training epochs is 1:3, 2:5, 1:2 on the Waterbirds, CelebA, and MultiNLI datasets respectively. This fast convergence rate is predicted by our theory, particularly Corollary 5.2. In our Figure 1, we also replicate the Figure 2 of Sagawa et al. [36], appending our additional results on R-MDL. We again see a trend of faster test error convergence (solid lines) and more uniform per-group risks by the R-MDL algorithm.

**Datasets** Our experiments were performed on three datasets: Multi-NLI, CelebA, and Waterbirds [36]. We use identical preprocessing settings and dataset splits as Sagawa et al. [36]. Our experiments, unless otherwise specified, replicate the exact hyperparameter settings adopted by Sagawa et al. [36] for their Table 2 experiments. This includes the choice of random seeds, batch sizes, learning rates, learning schedules, and regularization. We defer readers to Sagawa et al. [36] or to our public source code for replication details.

The **Multi-NLI dataset** [42] concerns the following natural language inference task: determine if one statement is entailed by, neutral with, or contradicts a given statement. This dataset is challenging because traditional ERM models are prone to spuriously correlating “contradiction” labels with the existence of negation words. The dataset is divided into 6 distributions: the Cartesian product of the label space (entailment, neutral, contradiction) and an indicator of whether the sentence contains a negation word. The label space annotations were annotated by [42] while negation labels were annotated by Sagawa et al. [36]. There are 206,175 datapoints available in the Multi-NLI dataset; the smallest distribution (entailment + negation) is represented by only 1,521 datapoints. We use a randomly shuffled 50-20-30 training-validation-testing split.

The **CelebA dataset** is a dataset of celebrity face images and a label space of potential physical attributes [23]. This dataset is challenging because traditional ERM models are prone to spuriously correlating attribute labels with demographic information such as race and gender. Following Sagawa et al. [36], we divide the dataset into 4 distributions: the Cartesian product of the blond vs dark hair attribute label (“Blond\_Hair”) with the “gender” attribute label (“Male”). Note that the authors of Liu et al. [23] limited the “gender” attribute label to binary options of male and not male. There are

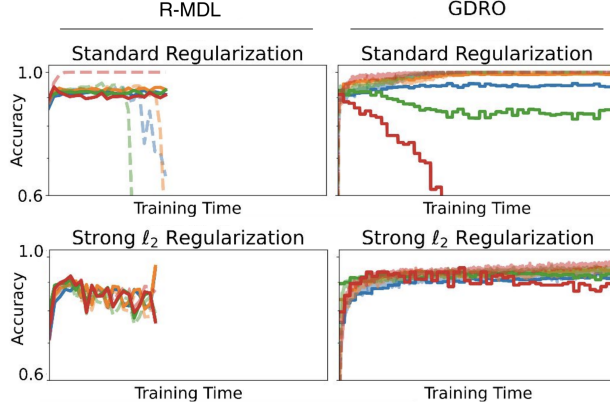


Figure 1: Training (light, dashed) and validation (dark, solid) accuracies for GDRO and R-MDL during training, plotted on a log scale. Note that R-MDL validation accuracy will be noisier than those of GDRO as we constrain R-MDL to limited samples (with replacement) from the validation set. In addition, in the left-most plot, training accuracy for all groups except the blond male group (red) dips to zero due to lack of data—this is because the blond male group (red) is the most challenging so the adversary eventually stops sampling from other groups. Under standard regularization, the red-group accuracy drops off in GDRO while R-MDL maintains a high red-group accuracy by heavily sampling from the red group, as reflected in the near-perfect red-group training error.

162,770 datapoints available in the CelebA dataset; the smallest distribution (blond-hair + male) is represented by only 1,387 datapoints. We use the official training-testing-validation dataset split.

The **Waterbirds dataset** is a dataset by Sagawa et al. [36] curated from a larger Caltech-UCSD Birds-200-2011 (CUB) dataset [41]. It concerns the task of predicting whether a bird is of some waterbird (sub)species from an image of said bird. This dataset is challenging because traditional ERM models are prone to spuriously correlating backgrounds with foreground subjects; for instance, a model may often predict that a bird is a waterbird only because the image of the bird was taken at a beach. The dataset has 4 distributions: the Cartesian product of the waterbird vs not waterbird label with whether the background of the picture is over water. There are 4,795 datapoints available in the Waterbirds dataset; the smallest distribution (waterbirds on land) is represented by only 56 examples.

**Models** We use two classes of models in our experiments: Resnet-50 [18] and BERT [11]. We use the *torchvision* [25] implementation of the convolutional neural network Resnet-50, with a default choice of a stochastic gradient descent optimizer with momentum 0.9 and batch size 128. Batch normalization is used; data augmentation and dropout are not used. We use the *HuggingFace* [43] implementation of the language model BERT, with a default choice of an Adam optimizer with dropout and batch size 32.

**Hyperparameters** In the *Standard Regularization* experiments, we use a Resnet-50 model with an  $\ell_2$  regularization parameter of  $\lambda = 0.0001$  and a fixed learning rate of  $\alpha = 0.001$  for both Waterbirds and CelebA datasets. The ERM and Group DRO baselines are trained on CelebA for 50 epochs and Waterbirds for 300 epochs. Our multi-distribution learning method is trained on CelebA for only 20 epochs and Waterbirds for 100 epochs; this is due to the faster training error convergence of our method. For the MultiNLI dataset, we use a BERT model with a linearly decaying learning rate starting at  $\alpha_0 = 0.00002$  and no  $\ell_2$  regularization. The ERM and Group DRO baselines are trained on MultiNLI for 20 epochs. Our multi-distribution learning method is trained on MultiNLI for only 10 epochs. Our multi-distribution learning method uses adversary learning rates  $\eta_+$  of 1, 1, 0.2 on Waterbirds, CelebA and MultiNLI respectively.

In the *Strong Regularization* experiments, we follow similar settings to the *Standard Regularization* experiments. The only change is that an  $\ell_2$  regularization parameter of  $\lambda = 1$  is used for Waterbirds and an  $\ell_2$  regularization parameter of  $\lambda = 0.1$  is used for CelebA. Our multi-distribution learning method uses adversary learning rates  $\eta_+$  of 1 and 0.2 on Waterbirds and CelebA respectively.

In the *Early Stopping* experiments, we follow similar settings to the *Standard Regularization* experiments. The only change is that all CelebA and Waterbird experiments are run for a single epoch. MultiNLI experiments are run for 3 epochs. Our multi-distribution learning method uses adversary learning rates  $\eta_+$  of 1, 1, 1 on Waterbirds, CelebA and MultiNLI respectively.

The only hyperparameters we use that differ from prior literature are the number of training epochs and the adversary learning rates of our method (R-MDL). The choice of epoch was not fine-tuned, and was selected due to our observation of early training error convergence. We selected our adversary learning rate  $\eta_-$  by training our method, on each dataset, for both  $\eta_- = 1$  and  $\eta_- = 0.2$  and selecting the  $\eta_-$  yielding the highest validation-split worst-group accuracy.

**Compute** The total amount of compute run for the experiments in this section is approximately 50 GPU hours. A “n1-standard-8” machine was leased from the Cloud computing service Google Cloud; the “n1-standard-8” machine was equipped with 8 Intel Broadwell chips and 1 NVIDIA Tesla V100 GPU. The cost of these computing resources totaled approximately USD \$2 per hour, with a total cost of approximately USD \$100. All results described in this section, with the exception of existing results cited from other works, were obtained with experiments on said machine. All experiments were implemented in Python and PyTorch.