

Agnostic Federated Learning

Mehryar Mohri

Google Research and Courant Institute of Mathematical Sciences, New York

MOHRI@GOOGLE.COM

Gary Sivek

Google Research, New York

GSIVEK@GOOGLE.COM

Ananda Theertha Suresh

Google Research, New York

THEERTHA@GOOGLE.COM

Abstract

A key learning scenario in large-scale applications is that of *federated learning*, where a centralized model is trained based on data originating from a large number of clients. We argue that, with the existing training and inference, federated models can be biased towards different clients. Instead, we propose a new framework of *agnostic federated learning*, where the centralized model is optimized for any target distribution formed by a mixture of the client distributions. We further show that this framework naturally yields a notion of fairness. We present data-dependent Rademacher complexity guarantees for learning with this objective, which guide the definition of an algorithm for agnostic federated learning. We also give a fast stochastic optimization algorithm for solving the corresponding optimization problem, for which we prove convergence bounds, assuming a convex loss function and hypothesis set. We further empirically demonstrate the benefits of our approach in several datasets. Beyond federated learning, our framework and algorithm can be of interest to other learning scenarios such as cloud computing, domain adaptation, drifting, and other contexts where the training and test distributions do not coincide.

1. Motivation

A key learning scenario in large-scale applications is that of *federated learning*. In that scenario, a centralized model is trained based on data originating from a large number of clients, which may be mobile phones, other mobile devices, or sensors (Konečný, McMahan, Yu, Richtárik, Suresh, and Bacon, 2016b; Konečný, McMahan, Ramage, and Richtárik, 2016a). The training data typically remains distributed over the clients, each with possibly unreliable or relatively slow network connections.

Federated learning raises several types of issues and has been the topic of multiple research efforts. These include systems, networking and communication bottleneck problems due to frequent exchanges between the central server and the clients. To deal with such problems, McMahan et al. (2017) suggested an averaging technique that consists of transmitting the central model to a subset of clients, training it with the data locally available, and averaging the local updates. Smith et al. (2017) proposed to further leverage the relationship between clients, assumed to be known, and cast

the problem as an instance of multi-task learning to derive local client models benefiting from other similar ones.

The optimization task in federated learning, which is a principal problem in this scenario, has also been the topic of multiple research work. That includes the design of more efficient communication strategies (Konečný, McMahan, Yu, Richtárik, Suresh, and Bacon, 2016b; Konečný, McMahan, Ramage, and Richtárik, 2016a; Suresh, Yu, Kumar, and McMahan, 2017), devising efficient distributed optimization methods benefiting from differential privacy guarantees (Agarwal, Suresh, Yu, Kumar, and McMahan, 2018), as well as recent guarantees for parallel stochastic optimization with a dependency graph (Woodworth, Wang, Smith, McMahan, and Srebro, 2018).

Another key problem in federated learning which appears more generally in distributed machine learning and other learning setups is that of *fairness*. In many instances in practice, the resulting learning models may be biased or unfair: they may discriminate against some protected groups (Bickel, Hammel, and O’Connell, 1975; Hardt, Price, Srebro, et al., 2016). As a simple example, a regression algorithm predicting a person’s salary could be using that person’s gender. This is a key problem in modern machine learning that does not seem to have been specifically studied in the context of federated learning.

While many problems related to federated learning have been extensively studied, the key objective of learning in that context seems not to have been carefully examined. We are also not aware of statistical guarantees derived for learning in this scenario. A crucial reason for such questions to emerge in this context is that the target distribution for which the centralized model is learned is unspecified. Which expected loss is federated learning seeking to minimize? Most centralized models for standard federated learning are trained on the aggregate training sample obtained from the subsamples drawn from the clients. Thus, if we denote by \mathcal{D}_k the distribution associated to client k , m_k the size of the sample available from that client and m the total sample size, intrinsically, the centralized model is trained to minimize the loss with respect to the *uniform distribution*

$$\bar{\mathcal{U}} = \sum_{k=1}^p \frac{m_k}{m} \mathcal{D}_k.$$

But why should $\bar{\mathcal{U}}$ be the target distribution of the learning model? Is $\bar{\mathcal{U}}$ the distribution that we expect to observe at test time? What guarantees can be derived for the deployed system?

Notice that, in practice, in federated learning, the probability that an individual data source participates in training depends on various factors such as whether the mobile device is connected to the internet or whether it is being charged. Thus, the training data may not truly reflect the usage of the learned model in inference. Additionally, these uncertainties may also affect the size of the sample m_k acquired from each client, which directly affects the definition of $\bar{\mathcal{U}}$.

We argue that in many common instances, the uniform distribution is not the natural objective distribution and that seeking to minimize the expected loss with respect to the specific distribution $\bar{\mathcal{U}}$ is *risky*. This is because the target distribution may be in general quite different from $\bar{\mathcal{U}}$. In many cases, that can result in a suboptimal or even a detrimental performance. For example, imagine a plausible scenario of federated learning where the learner has access to a large population of expensive mobile phones, which are most commonly adopted by software engineers or other technical users (say 70%) than other users (30%), and a small population of other mobile phones less used by non-technical users (5%) and significantly more often by other users (95%). The centralized model would then

某类数据占比更大，会导致模型向其倾斜，甚至变成单一类模型

be essentially based on the uniform distribution based on the expensive clients. But, clearly, such a model would not be adapted to the wide general target domain formed by the majority of phones with a 5%–95% population of general versus technical users. Many other realistic examples of this type can help illustrate the learning problem resulting from a mismatch between the target distribution and $\bar{\mathcal{U}}$. In fact, it is not clear why minimizing the expected loss with respect to $\bar{\mathcal{U}}$ could be beneficial for the clients, whose distributions are \mathcal{D}_k s.

Thus, we put forward a new framework of *agnostic federated learning* (AFL), where the centralized model is optimized for any possible target distribution formed by a mixture of the client distributions. Instead of optimizing the centralized model for a specific distribution, with the high risk of a mismatch with the target, we define an agnostic and more risk-averse objective. We show that, for some target mixture distributions, the cross-entropy loss of the hypothesis obtained by minimization with respect to the uniform distribution $\bar{\mathcal{U}}$ can be worse, by a constant additive term, than that of the hypothesis obtained in AFL, even if the learner has access to an infinite sample size (Section 3.2).

We further show that our AFL framework naturally yields a notion of fairness, which we refer to as *good-intent fairness* (Section 3.3). Indeed, the predictor solution of the optimization problem for our AFL framework treats all protected categories similarly. Beyond federated learning, our framework and solution also cover related problems in cloud-based learning services, where customers may not have any training data at their disposal or may not be willing to share that data with the cloud. In that case too, the server needs to train a model without access to the training data. Our framework and algorithm can also be of interest to other learning scenarios such as domain adaptation, drifting, and other contexts where the training and test distributions do not coincide.

The rest of the paper is organized as follows. In Section 2, we give an extensive discussion of related work, including connections with the broad literature of domain adaptation. In Section 3, we give a formal description of the learning scenario of federated learning and the formulation of the problem as AFL. Next, we give a detailed theoretical analysis of learning in the AFL framework, including data-dependent Rademacher complexity generalization bounds (Section 4). These bounds lead to a natural learning algorithm with a regularization term based on a *skewness term* that we define (Section 5). We also present an efficient convex optimization algorithm for solving the optimization problem defining our algorithm (Section 5.2). Our algorithm is a stochastic gradient-descent solution for minimax problems, for which we give a detailed analysis, including the proof of convergence in terms of the variances of the stochastic gradients. In Section 6, we present a series of experiments comparing our AFL algorithm and solution with existing federated learning solutions. In Section 7, we discuss several extensions of AFL.

2. Related work

Here, we briefly discuss several learning scenarios and work related to our study of federated learning.

The problem of federated learning is closely related to other learning scenarios where there is a mismatch between the source distribution and the target distribution. This includes the problem of *transfer learning* or *domain adaptation* from a single source to a known target domain (Ben-David, Blitzer, Crammer, and Pereira, 2006; Mansour, Mohri, and Rostamizadeh, 2009b; Cortes and Mohri, 2014; Cortes, Mohri, and Muñoz Medina, 2015), either through unsupervised adaptation techniques

(Gong et al., 2012; Long et al., 2015; Ganin and Lempitsky, 2015; Tzeng et al., 2015), or via lightly supervised ones (some amount of labeled data from the target domain) (Saenko et al., 2010; Yang et al., 2007; Hoffman et al., 2013; Girshick et al., 2014). This also includes previous applications in natural language processing (Dredze et al., 2007; Blitzer et al., 2007; Jiang and Zhai, 2007; Raju et al., 2018), speech recognition (Legetter and Woodland, 1995; Gauvain and Chin-Hui, 1994; Pietra et al., 1992; Rosenfeld, 1996; Jelinek, 1998; Roark and Bacchiani, 2003), and computer vision (Martínez, 2002)

A problem more closely related to that of federated learning is that of *multiple-source adaptation*, first formalized and analyzed theoretically by Mansour, Mohri, and Rostamizadeh (2009c,a) and later studied for various applications such as object recognition (Hoffman et al., 2012; Gong et al., 2013a,b). Recently, Zhang et al. (2015) studied a causal formulation of this problem for a classification scenario, using the same combination rules as Mansour et al. (2009c,a). The problem of *domain generalization* (Pan and Yang, 2010; Muandet et al., 2013; Xu et al., 2014), where knowledge from an arbitrary number of related domains is combined to perform well on a previously unseen domain is very closely related to that of federated learning, though the assumptions about the information available to the learner and the availability of unlabeled data may differ.

In the multiple-source adaptation problem studied by Mansour, Mohri, and Rostamizadeh (2009c,a) and Hoffman, Mohri, and Zhang (2018), each domain k is defined by the corresponding distribution \mathcal{D}_k and the learner has only access to a predictor h_k for each domain and no access to labeled training data drawn from these domains. The authors show that it is possible to define a predictor h whose expected loss $\mathcal{L}_{\mathcal{D}}(h)$ with respect to any distribution \mathcal{D} that is a mixture of the source domains \mathcal{D}_k is at most the maximum expected loss of the source predictors: $\max_k L_{\mathcal{D}_k}(h_{\mathcal{D}_k})$. They also provide an algorithm for determining h .

Our learning scenario differs from the one adopted in that work since we assume access to labeled training data from each domain \mathcal{D}_k . Furthermore, the predictor determined by the algorithm of Hoffman, Mohri, and Zhang (2018) belongs to a specific hypothesis set \mathcal{H}' , which is that of distribution weighted combinations of the domain predictors h_k , while, in our setup, the objective is to determine the best predictor in some global hypothesis set \mathcal{H} , which may include \mathcal{H}' as a subset, and which is not depending on some domain-specific predictors.

Our optimization solution also differs from the work of Farnia and Tse (2016) and Lee and Raginsky (2017) on local minimax results, where samples are drawn from a single source \mathcal{D} , and where the generalization error is minimized over a set of locally ambiguous distributions $\widehat{\mathcal{D}}$, where $\widehat{\mathcal{D}}$ is the empirical distribution. The authors propose this metric for statistical robustness. In our work, we obtain samples from p unknown distributions, and the set of distributions \mathcal{D}_λ over which we optimize the expected loss is fixed and independent of samples. Furthermore, the source distributions can differ arbitrarily and need not be close to each other. In reverse, we note that our stochastic algorithm can be used to minimize the loss functions proposed in (Farnia and Tse, 2016; Lee and Raginsky, 2017).

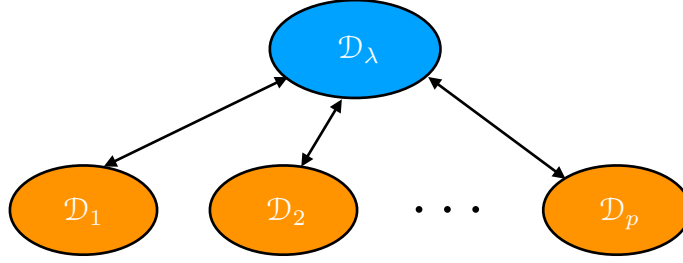


Figure 1: Illustration of the agnostic federated learning scenario.

3. Learning scenario

In this section, we introduce the learning scenario of agnostic federated learning we consider. Next, we first argue that the uniform solution commonly adopted in standard federated learning may not be an adequate solution, thereby further justifying our agnostic model. Second, we show the benefit of our model in fairness learning.

We start with some general notation and definitions used throughout the paper. Let \mathcal{X} denote the input space and \mathcal{Y} the output space. We will primarily discuss a multi-class classification problem where \mathcal{Y} is a finite set of classes, but much of our results can be extended straightforwardly to regression and other problems. The hypotheses we consider are of the form $h: \mathcal{X} \rightarrow \Delta_{\mathcal{Y}}$, where $\Delta_{\mathcal{Y}}$ stands for the simplex over \mathcal{Y} . Thus, $h(x)$ is a probability distribution over the classes or categories that can be assigned to $x \in \mathcal{X}$. We will denote by \mathcal{H} a family of such hypotheses h . We also denote by ℓ a loss function defined over $\Delta_{\mathcal{Y}} \times \mathcal{Y}$ and taking non-negative values. The loss of $h \in \mathcal{H}$ for a labeled sample $(x, y) \in \mathcal{X} \times \mathcal{Y}$ is given by $\ell(h(x), y)$. One key example in applications is the cross-entropy loss, which is defined as follows: $\ell(h(x), y) = -\log(\mathbb{P}_{y' \sim h(x)}[y' = y])$. We will denote by $\mathcal{L}_{\mathcal{D}}(h)$ the expected loss of a hypothesis h with respect to a distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$:

$$\mathcal{L}_{\mathcal{D}}(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(h(x), y)],$$

and by $h_{\mathcal{D}}$ its minimizer: $h_{\mathcal{D}} = \operatorname{argmin}_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{D}}(h)$.

3.1. Agnostic federated learning

We consider a learning scenario where the learner receives p samples S_1, \dots, S_p , with each $S_k = ((x_{k,1}, y_{k,1}), \dots, (x_{k,m_k}, y_{k,m_k})) \in (\mathcal{X} \times \mathcal{Y})^{m_k}$ of size m_k drawn i.i.d. from a different domain or distribution \mathcal{D}_k . The learner's objective is to determine a hypothesis $h \in \mathcal{H}$ that performs well on some target distribution. We will also denote by $\widehat{\mathcal{D}}_k$ the empirical distribution associated to sample S_k of size m drawn from \mathcal{D}^m .

This scenario coincides with that of *federated learning* where training is done with the *uniform distribution* over the union of all samples S_k , that is $\widehat{\mathcal{U}} = \sum_{k=1}^p \frac{m_k}{\sum_{k=1}^p m_k} \widehat{\mathcal{D}}_k$, and where the underlying assumption is that the target distribution is $\overline{\mathcal{U}} = \sum_{k=1}^p \frac{m_k}{\sum_{k=1}^p m_k} \mathcal{D}_k$. We will not adopt that assumption since it is rather restrictive and since, as discussed later, it can lead to solutions that are

disadvantageous to domain users. Instead, we will consider an *agnostic federated learning* (AFL) scenario where the target distribution can be modeled as an unknown mixture of the distributions \mathcal{D}_k , $k = 1, \dots, p$, that is $\mathcal{D}_\lambda = \sum_{k=1}^p \lambda_k \mathcal{D}_k$ for some $\lambda \in \Delta_p$. Since the mixture weight λ is unknown, here, the learner must come up with a solution that is favorable for any λ in the simplex, or any λ in a subset $\Lambda \subseteq \Delta_p$. Thus, we define the *agnostic loss* (or *agnostic risk*) $\mathcal{L}_{\mathcal{D}_\Lambda}(h)$ associated to a predictor $h \in \mathcal{H}$ as

$$\mathcal{L}_{\mathcal{D}_\Lambda}(h) = \max_{\lambda \in \Lambda} \mathcal{L}_{\mathcal{D}_\lambda}(h). \quad (1)$$

We will extend our previous definitions and denote by $h_{\mathcal{D}_\Lambda}$ the minimizer of this loss:

$$h_{\mathcal{D}_\Lambda} = \operatorname{argmin}_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{D}_\Lambda}(h).$$

In practice, the learner has access to the distributions \mathcal{D}_k only via the finite samples S_k . Thus, for any $\lambda \in \Delta_p$, instead of the mixture \mathcal{D}_λ , only the λ -mixture of empirical distributions, $\overline{\mathcal{D}}_\lambda = \sum_{k=1}^p \lambda_k \widehat{\mathcal{D}}_k$, is accessible.¹ This leads to the definition of $\mathcal{L}_{\overline{\mathcal{D}}_\Lambda}(h)$, the *agnostic empirical loss* of a hypothesis $h \in \mathcal{H}$ for a subset of the simplex Λ :

$$\mathcal{L}_{\overline{\mathcal{D}}_\Lambda}(h) = \max_{\lambda \in \Lambda} \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h).$$

We will denote by $h_{\overline{\mathcal{D}}_\Lambda}$ the minimizer of this loss: $h_{\overline{\mathcal{D}}_\Lambda} = \operatorname{argmin}_{h \in \mathcal{H}} \mathcal{L}_{\overline{\mathcal{D}}_\Lambda}(h)$. In the next section, we will present generalization bounds relating the expected and empirical agnostic losses $\mathcal{L}_{\mathcal{D}_\Lambda}(h)$ and $\mathcal{L}_{\overline{\mathcal{D}}_\Lambda}(h)$ for all $h \in \mathcal{H}$.

Notice that the domains \mathcal{D}_k discussed thus far need not coincide with the clients. In fact, when the number of clients is very large and Λ is the full simplex, $\Lambda = \Delta_p$, it is typically preferable to consider instead domains defined by clusters of clients, as discussed in Section 7. On the other hand, if p is small or Λ more restrictive, then the model may not perform well on certain domains of interest. We mitigate the effect of large p values using a suitable regularization term derived from our theory.

3.2. Comparison with federated learning

Here, we further argue that the uniform solution $h_{\overline{\mathcal{U}}}$ commonly adopted in federated learning may not provide a satisfactory performance compared with a solution of the agnostic problem. This further motivates our AFL model.

As already discussed, since the target distribution is unknown, the natural method for the learner is to select a hypothesis minimizing the agnostic loss $\mathcal{L}_{\mathcal{D}_\Lambda}$. Is the predictor minimizing the agnostic loss coinciding with the solution $h_{\overline{\mathcal{U}}}$ of standard federated learning? How poor can the performance of the standard federated learning be? We first show that the loss of $h_{\overline{\mathcal{U}}}$ can be higher than that of the optimal loss achieved by $h_{\mathcal{D}_\Lambda}$ by a constant loss, even if the number of samples tends to infinity, that is even if the learner has access to the distributions \mathcal{D}_k and uses the predictor $h_{\overline{\mathcal{U}}}$. Similar results are known for universal compression, where the goal is to compress a sequence of random variables without knowledge of the generating distribution (Grünwald, 2007).

1. Note, $\overline{\mathcal{D}}_\lambda$ is distinct from an empirical distribution $\widehat{\mathcal{D}}_\lambda$ which would be based on a sample drawn from \mathcal{D}_λ . $\overline{\mathcal{D}}_\lambda$ is based on samples drawn from \mathcal{D}_k s.

Proposition 1 *Let ℓ be the cross-entropy loss. Then, there exist Λ , \mathcal{H} , and \mathcal{D}_k , $k \in [p]$, such that the following inequality holds:*

$$\mathcal{L}_{\mathcal{D}_\Lambda}(h_{\bar{\mathcal{U}}}) \geq \mathcal{L}_{\mathcal{D}_\Lambda}(h_{\mathcal{D}_\Lambda}) + \log \frac{2}{\sqrt{3}}.$$

Proof Consider the following two distributions with support reduced to a single element $x \in \mathcal{X}$ and two classes $\mathcal{Y} = \{0, 1\}$: $\mathcal{D}_1(x, 0) = 0$, $\mathcal{D}_2(x, 1) = 1$, $\mathcal{D}_2(x, 0) = \frac{1}{2}$, and $\mathcal{D}_2(x, 1) = \frac{1}{2}$. Let $\Lambda = \{\delta_1, \delta_2\}$, where δ_k , $k = 1, 2$, denotes the Dirac measure on index k . We will consider the case where the sample sizes m_k are all equal, that is $h_{\bar{\mathcal{U}}} = \frac{1}{2}(\mathcal{D}_1 + \mathcal{D}_2)$. Let p_0 denote the probability that h assigns to class 0 and p_1 the one it assigns to class 1. Then, the cross-entropy loss of a predictor h can be expressed as follows:

$$\begin{aligned} \mathcal{L}_{\bar{\mathcal{U}}}(h) &= \mathbb{E}_{(x,y) \sim \bar{\mathcal{U}}} [-\log p_y] = \frac{1}{4} \log \frac{1}{p_0} + \frac{1}{2} \log \frac{1}{p_1} + \frac{1}{4} \log \frac{1}{p_1} \\ &= \frac{1}{4} \log \frac{1}{p_0} + \frac{3}{4} \log \frac{1}{p_1} \\ &= D\left(\left(\frac{1}{4}, \frac{3}{4}\right) \parallel (p_0, p_1)\right) + \frac{1}{4} \log \frac{4}{1} + \frac{3}{4} \log \frac{4}{3} \\ &\geq \frac{1}{4} \log \frac{4}{1} + \frac{3}{4} \log \frac{4}{3}, \end{aligned}$$

where the last inequality follows the non-negativity of the relative entropy. Furthermore, equality is achieved when $p_0 = 1 - p_1 = \frac{1}{4}$, which defines $h_{\bar{\mathcal{U}}}$, the minimizer of $\mathcal{L}_{\bar{\mathcal{U}}}(h)$. In view of that, $\mathcal{L}_{\mathcal{D}_\Lambda}(h_{\bar{\mathcal{U}}})$ is given by the following:

$$\begin{aligned} \mathcal{L}_{\mathcal{D}_\Lambda}(h_{\bar{\mathcal{U}}}) &= \max(\mathcal{L}_{\delta_1}(\bar{\mathcal{U}}), \mathcal{L}_{\delta_2}(\bar{\mathcal{U}})) \\ &= \max\left\{\log \frac{4}{3}, \frac{1}{2} \log \frac{4}{1} + \frac{1}{2} \log \frac{4}{3}\right\} \\ &= \log \frac{4}{\sqrt{3}}. \end{aligned}$$

We now compute the loss of $h_{\mathcal{D}_\Lambda}$:

$$\begin{aligned} \min_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{D}_\Lambda}(h) &= \min_{h \in \mathcal{H}} \max_{k \in [p]} \mathcal{L}_{\mathcal{D}_k}(h) \\ &= \min_{(p_0, p_1) \in \Delta_2} \max \left\{ \log \frac{1}{p_1}, \frac{1}{2} \log \frac{1}{p_0} + \frac{1}{2} \log \frac{1}{p_1} \right\} \\ &= \min_{p_1 \in [0, 1]} \max \left\{ \log \frac{1}{p_1}, \log \frac{1}{\sqrt{p_1(1-p_1)}} \right\} \\ &= \log 2, \end{aligned}$$

since $\frac{1}{2}$ is the solution of the convex optimization in p_1 , in view of $\max \left\{ \frac{1}{p_1}, \frac{1}{\sqrt{p_1(1-p_1)}} \right\} = \frac{1}{\sqrt{p_1(1-p_1)}} \leq \frac{1}{2}$ for $p_1 > \frac{1}{2}$. \blacksquare

3.3. Good-intent fairness in learning

Here, we further discuss the relationship between our model of AFL and fairness in learning.

Fairness in machine learning has received much attention in recent past (Bickel et al., 1975; Hardt et al., 2016). There is now a broad literature on the topic with a variety of definitions of the notion of fairness. In a typical scenario, there is a protected class c among p classes c_1, c_2, \dots, c_p . While there are many definitions of fairness, the main objective of a fairness algorithm is to reduce bias and ensure that the model is fair to all the p protected categories, under some definition of fairness. The most common reasons for bias in machine learning algorithms are training data bias and overfitting bias. We first provide a brief explanation and illustration for both:

- the training data is biased: consider the regression task, where the goal is to predict the salary of a person based on features such as education, location, age, gender. Let gender be the protected class. If in the training data, there is a consistent discrimination against women irrespective of their education, e.g., their salary is lower, then we can conclude that the training data is inherently biased.
- the training procedure is biased: consider an image recognition task where the protected category is race. If the model is heavily trained on images based on certain races, then the resulting model will be biased because of over-fitting.

Our model of AFL can help define a notion of good-intent fairness, where we reduce the bias in the training procedure. Furthermore, if training procedure bias exists, it naturally highlights it.

Suppose we are interested in a classification problem and there is a protected feature class c , which can be one of p values c_1, c_2, \dots, c_p . Then, we define \mathcal{D}_k as the conditional distribution with the protected class being c_k . If \mathcal{D} is the true underlying distribution, then

$$\mathcal{D}_k(x, y) = \mathcal{D}(x, y \mid c(x, y) = c_k).$$

Let $\Lambda = \{\delta_k : k \in [p]\}$ be the collection of Dirac measures over the indices k in $[p]$. With this definition, we define a *good-intent fairness* algorithm as one seeking to minimize the agnostic loss $\mathcal{L}_{\mathcal{D}_\Lambda}$. Thus, the objective of the algorithm is to minimize the maximum loss incurred on any of the underlying protective classes and hence does not overfit the data to any particular model at the cost of others. Furthermore, it does not degrade the performance of the other classes so long as it does not affect the loss of the most-sensitive protected category. We further note that our approach does not reduce bias in the training data and is useful only for mitigating the training procedure bias.

4. Learning bounds

In this section, we present learning guarantees for agnostic federated learning. Let \mathcal{G} denote the family of the losses associated to a hypothesis set \mathcal{H} : $\mathcal{G} = \{(x, y) \mapsto \ell(h(x), y) : h \in \mathcal{H}\}$. Our learning bounds are based on the following notion of *weighted Rademacher complexity* which is defined for any hypothesis set \mathcal{H} , vector of sample sizes $\mathbf{m} = (m_1, \dots, m_p)$ and mixture weight

$\lambda \in \Delta_p$, by the following expression:

$$\mathfrak{R}_{\mathbf{m}}(\mathcal{G}, \lambda) = \mathbb{E}_{\substack{S_k \sim \mathcal{D}_k^{m_k} \\ \boldsymbol{\sigma}}} \left[\sup_{h \in \mathcal{H}} \sum_{k=1}^p \frac{\lambda_k}{m_k} \sum_{i=1}^{m_k} \sigma_{k,i} \ell(h(x_{k,i}), y_{k,i}) \right], \quad (2)$$

where $S_k = ((x_{k,1}, y_{k,1}), \dots, (x_{k,m_k}, y_{k,m_k}))$ is a sample of size m_k and $\boldsymbol{\sigma} = (\sigma_{k,i})_{k \in [p], i \in [m_k]}$ a collection of Rademacher variables, that is uniformly distributed random variables taking values in $\{-1, +1\}$. We also defined the *minimax weighted Rademacher complexity* for a subset $\Lambda \subseteq \Delta_p$ by

$$\mathfrak{R}_{\mathbf{m}}(\mathcal{G}, \Lambda) = \max_{\lambda \in \Lambda} \mathfrak{R}_{\mathbf{m}}(\mathcal{G}, \lambda). \quad (3)$$

Let $\overline{\mathbf{m}} = \frac{\mathbf{m}}{m} = (\frac{m_1}{m}, \dots, \frac{m_p}{m})$ denote the empirical distribution over Δ_p defined by the sample sizes m_k , where $m = \sum_{k=1}^p m_k$. We define the *skewness* of Λ with respect to $\overline{\mathbf{m}}$ by

$$\mathfrak{s}(\Lambda \parallel \overline{\mathbf{m}}) = \max_{\lambda \in \Lambda} \chi^2(\lambda \parallel \overline{\mathbf{m}}) + 1, \quad (4)$$

where, for any two distributions p and q in Δ_p , the chi-squared divergence $\chi^2(p \parallel q)$ is given by $\chi^2(p \parallel q) = \sum_{k=1}^p \frac{(p_k - q_k)^2}{q_k}$. We will also denote by Λ_ϵ a minimum ϵ -cover of Λ in ℓ_1 distance, that is,

$$\Lambda_\epsilon = \operatorname{argmin}_{\Lambda' \in C(\Lambda, \epsilon)} |\Lambda'|,$$

where $C(\Lambda, \epsilon)$ is a set of distributions Λ' such that for every $\lambda \in \Lambda$, there exists $\lambda' \in \Lambda'$ such that $\sum_{k=1}^p |\lambda_k - \lambda'_k| \leq \epsilon$.

Our first learning guarantee is presented in terms of $\mathfrak{R}_{\mathbf{m}}(\mathcal{G}, \Lambda)$, the skewness parameter $\mathfrak{s}(\Lambda \parallel \overline{\mathbf{m}})$ and the ϵ -cover Λ_ϵ .

Theorem 2 *Assume that the loss ℓ is bounded by $M > 0$. Fix $\epsilon > 0$ and $\mathbf{m} = (m_1, \dots, m_p)$. Then, for any $\delta > 0$, with probability at least $1 - \delta$ over the draw of samples $S_k \sim \mathcal{D}_k^{m_k}$, the following inequality holds for all $h \in \mathcal{H}$ and $\lambda \in \Lambda$:*

$$\mathcal{L}_{\mathcal{D}_\lambda}(h) \leq \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h) + 2\mathfrak{R}_{\mathbf{m}}(\mathcal{G}, \lambda) + M\epsilon + M\sqrt{\frac{\mathfrak{s}(\lambda \parallel \overline{\mathbf{m}})}{2m} \log \frac{|\Lambda_\epsilon|}{\delta}},$$

where $m = \sum_{k=1}^p m_k$.

Proof The proof is an extension of the standard proofs for Rademacher complexity generalization bounds (Koltchinskii and Panchenko, 2002; Mohri et al., 2018). Fix $\lambda \in \Lambda$. For any sample $S = S_1, \dots, S_p$, define $\Psi(S_1, \dots, S_p)$ by

$$\Psi(S_1, \dots, S_p) = \sup_{h \in \mathcal{H}} (\mathcal{L}_{\mathcal{D}_\lambda}(h) - \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h)).$$

Let $S' = (S'_1, \dots, S'_p)$ be a sample differing from $S = (S_1, \dots, S_p)$ only by point $x'_{k,i}$ in S'_k and $x_{k,i}$ in S_k . Then, since the difference of suprema over the same set is bounded by the supremum of the

differences, we can write

$$\begin{aligned}
 \Psi(S') - \Psi(S) &= \sup_{h \in \mathcal{H}} \left(\mathcal{L}_{\mathcal{D}_\lambda}(h) - \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h) \right) - \sup_{h \in \mathcal{H}} \left(\mathcal{L}_{\mathcal{D}_\lambda}(h) - \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h) \right) \\
 &\leq \sup_{h \in \mathcal{H}} \left(\mathcal{L}_{\mathcal{D}_\lambda}(h) - \mathcal{L}_{\overline{\mathcal{D}}_\lambda'}(h) \right) - \left(\mathcal{L}_{\mathcal{D}_\lambda}(h) - \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h) \right) \\
 &\leq \sup_{h \in \mathcal{H}} \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h) - \mathcal{L}_{\overline{\mathcal{D}}_\lambda'}(h) \\
 &= \sup_{h \in \mathcal{H}} \sum_{k=1}^p \frac{\lambda_k}{m_k} \sum_{i=1}^{m_k} \ell(h(x'_{k,i}), y'_{k,i}) - \sum_{k=1}^p \frac{\lambda_k}{m_k} \sum_{i=1}^{m_k} \ell(h(x_{k,i}), y_{k,i}) \\
 &= \sup_{h \in \mathcal{H}} \frac{\lambda_k}{m_k} \left[\ell(h(x'_{k,i}), y'_{k,i}) - \ell(h(x_{k,i}), y_{k,i}) \right] \\
 &\leq \frac{\lambda_k M}{m_k}.
 \end{aligned}$$

Thus, by McDiarmid's inequality, for any $\delta > 0$, the following inequality holds with probability at least $1 - \delta$ for any $h \in \mathcal{H}$:

$$\mathcal{L}_{\mathcal{D}_\lambda}(h) \leq \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h) + \mathbb{E} \left[\max_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{D}_\lambda}(h) - \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h) \right] + M \sqrt{\sum_{k=1}^p \frac{\lambda_k^2}{2m_k} \log \frac{1}{\delta}}.$$

Therefore, by the union over Λ_ϵ , with probability at least $1 - \delta$, for any $h \in \mathcal{H}$ and $\lambda \in \Lambda_\epsilon$ the following holds:

$$\mathcal{L}_{\mathcal{D}_\lambda}(h) \leq \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h) + \mathbb{E} \left[\max_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{D}_\lambda}(h) - \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h) \right] + M \sqrt{\sum_{k=1}^p \frac{\lambda_k^2}{2m_k} \log \frac{|\Lambda_\epsilon|}{\delta}}.$$

By definition of Λ_ϵ , for any $\lambda \in \Lambda$, there exists $\lambda' \in \Lambda_\epsilon$ such that $\mathcal{L}_{\mathcal{D}_\lambda}(h) \leq \mathcal{L}_{\mathcal{D}_{\lambda'}}(h) + M\epsilon$. In view of that, with probability at least $1 - \delta$, for any $h \in \mathcal{H}$ and $\lambda \in \Lambda$ the following holds:

$$\mathcal{L}_{\mathcal{D}_\lambda}(h) \leq \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h) + \mathbb{E} \left[\max_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{D}_\lambda}(h) - \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h) \right] + M\epsilon + M \sqrt{\sum_{k=1}^p \frac{\lambda_k^2}{2m_k} \log \frac{|\Lambda_\epsilon|}{\delta}}.$$

The expectation appearing on the right-hand side can be bounded following standard proofs for Rademacher complexity upper bounds (see for example (Mohri et al., 2018)), leading to

$$\mathbb{E} \left[\max_{h \in \mathcal{H}} \mathcal{L}_{\mathcal{D}_\lambda}(h) - \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h) \right] \leq \mathfrak{R}_{\mathbf{m}}(\mathcal{G}, \lambda).$$

The sum $\sum_{k=1}^p \frac{\lambda_k^2}{m_k}$ can be expressed in terms of the skewness of Λ , using the following equalities:

$$m \sum_{k=1}^p \frac{\lambda_k^2}{m_k} = \sum_{k=1}^p \frac{\lambda_k^2}{\frac{m_k}{m}} = \sum_{k=1}^p \frac{\lambda_k^2}{\frac{m_k}{m}} + \sum_{k=1}^p \frac{m_k}{m} - 2 \sum_{k=1}^p \lambda_k + 1 = \sum_{k=1}^p \frac{(\lambda_k - \frac{m_k}{m})^2}{\frac{m_k}{m}} + 1 = \chi^2(\lambda \parallel \overline{\mathbf{m}}) + 1.$$

This completes the proof. ■

It can be proven that the skewness parameter appears in a lower bound on the generalization bound. We will include that result in the final version of this paper. The theorem yields immediately upper

bounds for agnostic losses by taking the maximum over $\lambda \in \Lambda$: for any $\delta > 0$, with probability at least $1 - \delta$, for any $h \in \mathcal{H}$,

$$\begin{aligned} \mathcal{L}_{\mathcal{D}_\Lambda}(h) &\leq \max_{\lambda \in \Lambda} \left\{ \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h) + 2\mathfrak{R}_{\mathbf{m}}(\mathcal{G}, \lambda) + M\epsilon + M\sqrt{\frac{\mathfrak{s}(\lambda \parallel \overline{\mathbf{m}})}{2m} \log \frac{|\Lambda_\epsilon|}{\delta}} \right\} \\ &\leq \mathcal{L}_{\overline{\mathcal{D}}_\Lambda}(h) + \max_{\lambda \in \Lambda} \left\{ 2\mathfrak{R}_{\mathbf{m}}(\mathcal{G}, \lambda) + M\epsilon + M\sqrt{\frac{\mathfrak{s}(\lambda \parallel \overline{\mathbf{m}})}{2m} \log \frac{|\Lambda_\epsilon|}{\delta}} \right\} \\ &\leq \mathcal{L}_{\overline{\mathcal{D}}_\Lambda}(h) + 2\mathfrak{R}_{\mathbf{m}}(\mathcal{G}, \Lambda) + M\epsilon + M\sqrt{\frac{\mathfrak{s}(\Lambda \parallel \overline{\mathbf{m}})}{2m} \log \frac{|\Lambda_\epsilon|}{\delta}}. \end{aligned}$$

The following result shows that, for a family of functions taking values in $\{-1, +1\}$, the Rademacher complexity $\mathfrak{R}_{\mathbf{m}}(\mathcal{G}, \Lambda)$ can be bounded in terms of the VC-dimension and the skewness of Λ .

Lemma 3 *Let ℓ be a loss function taking values in $\{-1, +1\}$ and such that the family of losses \mathcal{G} admits VC-dimension d . Then, the following upper bound holds for the weighted Rademacher complexity of \mathcal{G} :*

$$\mathfrak{R}_{\mathbf{m}}(\mathcal{G}, \Lambda) \leq \sqrt{2\mathfrak{s}(\Lambda \parallel \overline{\mathbf{m}}) \frac{d}{m} \log \left[\frac{em}{d} \right]}.$$

Proof For any $\lambda \in \Lambda$, define the set of vectors A_λ in \mathbb{R}^m by

$$A_\lambda = \left\{ \left[\frac{\lambda_k}{m_k} \ell(h(x_{k,i}), y_{k,i}) \right]_{(k,i) \in [p] \times [m_k]} : \mathbf{x} \in \mathcal{X}^m, \mathbf{y} \in \mathcal{Y}^m \right\}.$$

For any $\mathbf{a} \in A_\lambda$, $\|\mathbf{a}\|_2 = \sqrt{\sum_{k=1}^p m_k \frac{\lambda_k^2}{m_k^2}} = \sqrt{\sum_{k=1}^p \frac{\lambda_k^2}{m_k}} \leq \sqrt{\frac{\mathfrak{s}(\Lambda \parallel \overline{\mathbf{m}})}{m}}$. Then, by Massart's lemma, for any $\lambda \in \Lambda$, the following inequalities hold:

$$\begin{aligned} \mathfrak{R}_{\mathbf{m}}(\mathcal{G}, \lambda) &= \mathbb{E}_{\substack{S_k \sim \mathcal{D}_k^{m_k} \\ \sigma}} \left[\sup_{h \in \mathcal{H}} \sum_{k=1}^p \frac{\lambda_k}{m_k} \sum_{i=1}^{m_k} \sigma_{k,i} \ell(h(x_{k,i}), y_{k,i}) \right] \\ &\leq \mathbb{E}_{\sigma} \left[\sup_{\mathbf{a} \in A} \sum_{k=1}^p \sum_{i=1}^{m_k} \sigma_{k,i} a_{k,i} \right] \\ &\leq \sqrt{\frac{\mathfrak{s}(\Lambda \parallel \overline{\mathbf{m}})}{m}} \frac{\sqrt{2 \log |A_\lambda|}}{m} \\ &= \frac{\sqrt{2\mathfrak{s}(\Lambda \parallel \overline{\mathbf{m}}) \log |A_\lambda|}}{m}. \end{aligned}$$

By Sauer's lemma, the following holds for $m \geq d$: $|A_\lambda| \leq \left(\frac{em}{d}\right)^d$. Plugging in the right-hand side in the inequality above completes the proof. \blacksquare

Both Lemma 3 and the generalization bound of Theorem 2 can thus be expressed in terms of the skewness parameter $\mathfrak{s}(\Lambda \parallel \overline{\mathbf{m}})$. Note that modulo the skewness parameter, the results look very similar to standard generalization bounds (Mohri et al., 2018). Furthermore, when Λ contains only

one distribution and is the average distribution, that is $\lambda_k = m_k/m$, then the skewness is equal to one and the results coincide with the standard guarantees in supervised learning.

Theorem 2 and Lemma 3 also provide guidelines for choosing the domains and Λ . When p is large and $\Lambda = \Delta_p$, then, the number of samples per domain could be small, the skewness parameter $\mathfrak{s}(\Lambda \parallel \bar{\mathbf{m}}) = \max_{1 \leq k \leq p} \frac{1}{m_k}$ would then be large and the generalization guarantees for the model would become weaker. We suggest some guidelines for choosing domains in Section 7. We further note that for a given p , if Λ contains distributions that are close to $\bar{\mathbf{m}}$, then the model generalizes well.

The corollary above can be straightforwardly extended to cover the case where the test samples are drawn from some distribution \mathcal{D} , instead of \mathcal{D}_λ . Define $\ell_1(\mathcal{D}, \mathcal{D}_\Lambda)$ by $\ell_1(\mathcal{D}, \mathcal{D}_\Lambda) = \min_{\lambda \in \Lambda} \ell_1(\mathcal{D}, \mathcal{D}_\lambda)$. Then, the following result holds.

Corollary 4 *Assume that the loss function ℓ is bounded by M . Then, for any $\epsilon \geq 0$ and $\delta > 0$, with probability at least $1 - \delta$, the following inequality holds for all $h \in \mathcal{H}$:*

$$\mathcal{L}_{\mathcal{D}}(h) \leq \mathcal{L}_{\overline{\mathcal{D}}_\Lambda}(h) + 2\mathfrak{R}_{\mathbf{m}}(\mathcal{G}, \Lambda) + M\ell_1(\mathcal{D}, \mathcal{D}_\Lambda) + M\epsilon + M\sqrt{\frac{\mathfrak{s}(\Lambda \parallel \bar{\mathbf{m}})}{2m} \log \frac{|\Lambda_\epsilon|}{\delta}}.$$

One straightforward choice of the parameter ϵ is $\epsilon = \frac{1}{\sqrt{m}}$, but, depending on $|\Lambda_\epsilon|$ and other terms of the bound, more favorable choices may be possible. We conclude this section by adding that alternative learning bounds can be derived for this problem, as discussed in Appendix A.

5. Algorithm

In this section, we introduce a learning algorithm for agnostic federated learning using the guarantees proven in the previous section and discuss in detail an optimization solution.

5.1. Regularization

The learning guarantees of the previous section suggest minimizing the sum of the empirical AFL term $\mathcal{L}_{\overline{\mathcal{D}}_\Lambda}(h)$, a term controlling the complexity of \mathcal{H} and a term depending on the skewness parameter. Observe that, since $\mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h)$ is linear in λ , the following equality holds:

$$\mathcal{L}_{\overline{\mathcal{D}}_\Lambda}(h) = \mathcal{L}_{\overline{\mathcal{D}}_{\text{conv}(\Lambda)}}(h), \quad (5)$$

where $\text{conv}(\Lambda)$ is the convex hull of Λ . Assume that \mathcal{H} is a vector space that can be equipped with a norm $\|\cdot\|$, as with most hypothesis sets used in learning applications. Then, given Λ and the regularization parameters $r \geq 0$ and $\gamma \geq 0$, our learning guarantees suggest minimizing the regularized loss $\mathcal{L}_{\overline{\mathcal{D}}_{\Lambda_r}}(h) + \gamma\|h\|$, where $\|\cdot\|$ is a suitable norm controlling the complexity of \mathcal{H} and where Λ_r is defined by $\Lambda_r = \{\lambda \in \text{conv}(\Lambda) : 1 + \chi^2(\lambda \parallel \bar{\mathbf{m}}) \leq r\}$. This can be equivalently formulated as the following minimization problem:

$$\min_{h \in \mathcal{H}} \max_{\lambda \in \text{conv}(\Lambda)} \mathcal{L}_{\overline{\mathcal{D}}_\lambda}(h) + \gamma\|h\| - \mu\chi^2(\lambda \parallel \bar{\mathbf{m}}), \quad (6)$$

where $\mu \geq 0$ is a hyperparameter. This defines our algorithm for AFL.

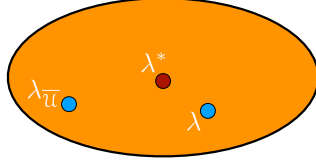


Figure 2: Illustration of the positions in Λ of λ^* , $\lambda_{\bar{u}}$, the mixture weight corresponding to the distribution \bar{u} , and an arbitrary λ . λ^* defines the least risky distribution $\bar{\mathcal{D}}_{\lambda^*}$ for which to optimize the expected loss.

Assume that ℓ is a convex function of its first argument. Then, $\mathcal{L}_{\bar{\mathcal{D}}_{\lambda}}(h)$ is a convex function of h . Since $\|h\|$ is a convex function of h for any choice of the norm, for a fixed λ , the objective $\mathcal{L}_{\bar{\mathcal{D}}_{\lambda}}(h) + \gamma\|h\| - \mu\chi^2(\lambda\|\bar{\mathbf{m}})$ is a convex function of h . The maximum over λ (taken in any set) of a family of convex functions is convex. Thus, $\max_{\lambda \in \text{conv}(\Lambda)} \mathcal{L}_{\bar{\mathcal{D}}_{\lambda}}(h) + \gamma\|h\| - \mu\chi^2(\lambda\|\bar{\mathbf{m}})$ is a convex function of h and, when the hypothesis set \mathcal{H} is a convex, (6) is a convex optimization problem. In the next subsection, we present an efficient optimization solution for this problem, for which we prove convergence guarantees.

5.2. Optimization algorithm

When the loss function ℓ is convex, the AFL minmax optimization problem above can be solved using projected gradient descent or other instances of the generic mirror descent algorithm (Nemirovski and Yudin, 1983). However, for large datasets, that is p and m large, this can be computationally costly and typically slow in practice. Juditsky, Nemirovski, and Tauvel (2011) proposed a stochastic Mirror-Prox algorithm for solving stochastic variational inequalities, which would be applicable in our context. We present a simplified version of their algorithm for the AFL problem that admits a more straightforward analysis and that is also substantially easier to implement.

Our optimization problem is over two sets of parameters, the hypothesis $h \in \mathcal{H}$ and the mixture weight $\lambda \in \Lambda$. In what follows, we will denote by $w \in \mathcal{W} \subset \mathbb{R}^N$ a vector of parameters defining a predictor h and will rewrite losses and optimization solutions only in terms of w , instead of h . We will use the following notation:

$$\mathbf{L}(w, \lambda) = \sum_{k=1}^p \lambda_k \mathbf{L}_k(w), \quad (7)$$

where $\mathbf{L}_k(w)$ stands for $\mathcal{L}_{\bar{\mathcal{D}}_k}(h)$, the empirical loss of hypothesis $h \in \mathcal{H}$ (corresponding to w) on domain k :

$$\mathbf{L}_k(w) = \frac{1}{m_k} \sum_{i=1}^{m_k} \ell(h(x_{k,i}), y_{k,i}).$$

Since the regularization terms do not make the optimization problem harder, to simplify the discussion, we will consider the unregularized version of problem (6). Thus, we will study the following problem given by the set of variables w :

$$\min_{w \in \mathcal{W}} \max_{\lambda \in \Lambda} \mathbf{L}(w, \lambda). \quad (8)$$

Observe that problem (8) admits a natural game-theoretic interpretation as a two-player game, where nature selects $\lambda \in \Lambda$ to maximize the objective, while the learner seeks $w \in \mathcal{W}$ minimizing the loss. We are interested in finding the equilibrium of this game, which is attained for some w^* , the minimizer of Equation 8 and $\lambda^* \in \Lambda$, the hardest domain mixture weights. At the equilibrium, moving w away from w^* or λ from λ^* , increases the objective function. Hence, λ^* can be viewed as the center of Λ in the manifold imposed by the loss function L , whereas \bar{U} , the empirical distribution of samples, may lie elsewhere, as illustrated by Figure 2.

By Equation (5), using the set $\text{conv}(\Lambda)$ instead of Λ does not affect the solution of the optimization problem. In view of that, in what follows, we will assume, without loss of generality, that Λ is a convex set. Observe that, since $L_k(w)$ is not an average of functions, standard stochastic gradient descent algorithms cannot be used to minimize this objective. We will present instead a new stochastic gradient-type algorithm for this problem.

Let $\nabla_w L(w, \lambda)$ denote the gradient of the loss function with respect to w and $\nabla_\lambda L(w, \lambda)$ the gradient with respect to λ . Let $\delta_w L(w, \lambda)$, and $\delta_\lambda L(w, \lambda)$ be unbiased estimates of the gradient, that is,

$$\mathbb{E}_\delta[\delta_\lambda L(w, \lambda)] = \nabla_\lambda L(w, \lambda) \text{ and } \mathbb{E}_\delta[\delta_w L(w, \lambda)] = \nabla_w L(w, \lambda).$$

We first give an optimization algorithm STOCHASTIC-AFL for the AFL problem, assuming access to such unbiased estimates. The pseudocode of the algorithm is given in Figure 3. At each step, the algorithm computes a stochastic gradient with respect to λ and w and updates the model accordingly. It then projects λ to Λ by computing a value in Λ via convex minimization. If Λ is the full simplex, then there is a near-linear time algorithm for this projection Wang and Carreira-Perpinán (2013). It then repeats the process for T steps and return the average of the weights. We provide guarantees for this algorithm in terms of the variance of the stochastic gradients when the loss function L is convex and when the set of w s, \mathcal{W} , is a compact set.

In the above analysis and in algorithm description in 3, we have ignored the regularization term. If the objective contains a regularization term such as Equation 6, then for λ_k , the regularization term yields a derivative of $-2\gamma\lambda_k/\bar{\mathbf{m}}_k$, which can be added to $\delta_\lambda L(w, \lambda)$ in Step 3 in Algorithm 3.

There are several natural candidates for the sampling method defining stochastic gradients. We highlight two techniques: PERDOMAIN GRADIENT and WEIGHTED GRADIENT. We analyze the time complexity and give bounds on the variance for both techniques in Lemmas 8 and 9 respectively.

Recently, Rakhlin and Sridharan (2013) and Daskalakis et al. (2017) gave an optimistic gradient descent algorithm for minimax optimizations. Our algorithm can also be modified to derive a stochastic optimistic algorithm, which we refer to as OPTIMISTIC-STOCHASTIC-AFL. The pseudocode of this algorithm is also given in Figure 3. However, the convergence analysis we present in the next section does not cover this algorithm.

5.3. Analysis

Throughout this section, for simplicity, we adopt the notation introduced for Equation 7. Our convergence guarantees hold under the following assumptions, which are similar to those adopted for the convergence proof of gradient descent-type algorithms.

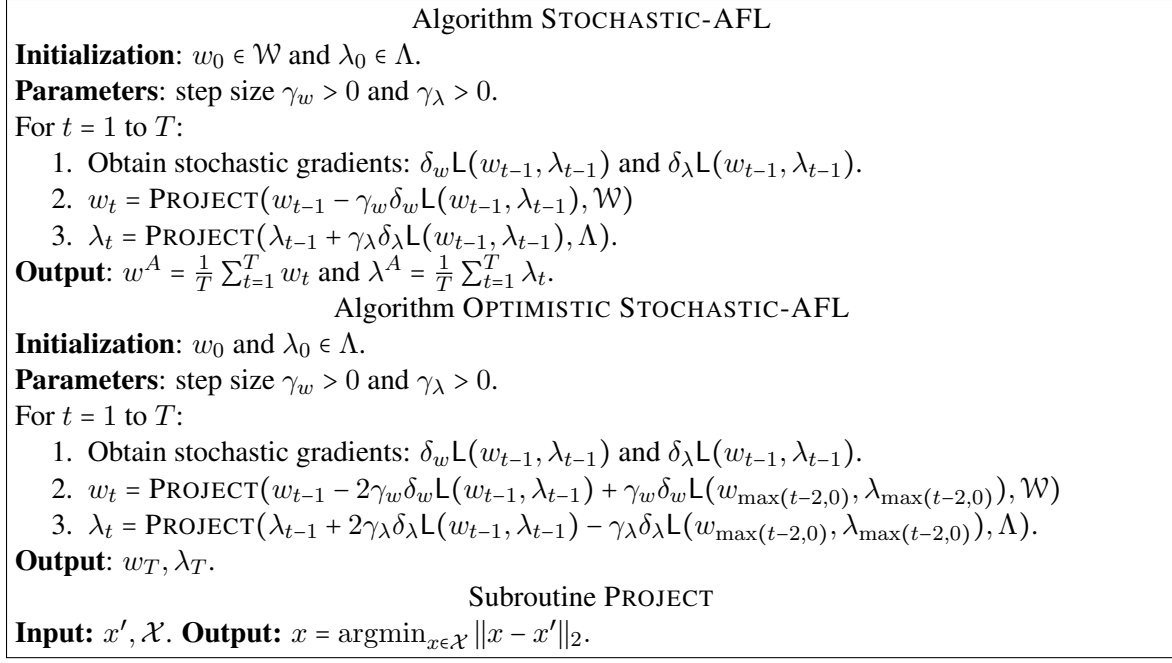


Figure 3: Pseudocodes of the STOCHASTIC-AFL and OPTIMISTIC STOCHASTIC-AFL algorithms.

Properties 1 Assume that the following properties hold for the loss function \mathcal{L} and sets \mathcal{W} and $\Lambda \subseteq \Delta_p$:

1. Convexity: $w \mapsto \mathcal{L}(w, \lambda)$ is convex for any $\lambda \in \Lambda$.
2. Compactness: $\max_{\lambda \in \Lambda} \|\lambda\|_2 \leq R_\Lambda$ and $\max_{w \in \mathcal{W}} \|w\|_2 \leq R_W$, for some $R_\Lambda > 0$ and $R_W > 0$.
3. Bounded gradients: $\|\nabla_w \mathcal{L}(w, \lambda)\|_2 \leq G_w$ and $\|\nabla_\lambda \mathcal{L}(w, \lambda)\|_2 \leq G_\lambda$ for all $w \in \mathcal{W}$ and $\lambda \in \Lambda$.
4. Stochastic variance: $\mathbb{E}[\|\delta_w \mathcal{L}(w, \lambda) - \nabla_w \mathcal{L}(w, \lambda)\|_2^2] \leq \sigma_w^2$ and $\mathbb{E}[\|\delta_\lambda \mathcal{L}(w, \lambda) - \nabla_\lambda \mathcal{L}(w, \lambda)\|_2^2] \leq \sigma_\lambda^2$ for all $w \in \mathcal{W}$ and $\lambda \in \Lambda$.
5. Time complexity: U_w denotes the time complexity of computing $\delta_w \mathcal{L}(w, \lambda)$, U_λ that of computing $\delta_\lambda \mathcal{L}(w, \lambda)$, U_p that of the projection, and d denotes the dimensionality of \mathcal{W} .

Theorem 5 Assume that the Properties 1 hold. Then, for the steps sizes $\gamma_w = \frac{2R_W}{\sqrt{T(\sigma_w^2 + G_w^2)}}$ and $\gamma_\lambda = \frac{2R_\Lambda}{\sqrt{T(\sigma_\lambda^2 + G_\lambda^2)}}$, the following guarantee holds for STOCHASTIC-AFL:

$$\mathbb{E} \left[\max_{\lambda \in \Lambda} \mathcal{L}(w^A, \lambda) - \min_{w \in \mathcal{W}} \max_{\lambda \in \Lambda} \mathcal{L}(w, \lambda) \right] \leq \frac{3R_W \sqrt{(\sigma_w^2 + G_w^2)}}{\sqrt{T}} + \frac{3R_\Lambda \sqrt{(\sigma_\lambda^2 + G_\lambda^2)}}{\sqrt{T}}.$$

and the time complexity of the algorithm is in $\mathcal{O}((U_\lambda + U_w + U_p + d + k)T)$.

Proof The time complexity of the algorithm follows the definitions of the complexity terms U_λ , U_w , and U_p the dimension d in Properties 1. To prove the convergence guarantee, we make a series of reductions. Let w^A and λ^A be a solution returned by the algorithm. First observe that since \mathcal{L}

is convex in w and linear and thus concave in λ , by the generalized von Neumann's theorem, the following holds:

$$\begin{aligned}
 \max_{\lambda \in \Lambda} \mathcal{L}(w^A, \lambda) - \min_{w \in \mathcal{W}} \max_{\lambda \in \Lambda} \mathcal{L}(w, \lambda) &= \max_{\lambda \in \Lambda} \mathcal{L}(w^A, \lambda) - \max_{\lambda \in \Lambda} \min_{w \in \mathcal{W}} \mathcal{L}(w, \lambda) && \text{(von Neumann's minimax)} \\
 &\leq \max_{\lambda \in \Lambda} \left\{ \mathcal{L}(w^A, \lambda) - \min_{w \in \mathcal{W}} \mathcal{L}(w, \lambda^A) \right\} && \text{(subadd. of max)} \\
 &= \max_{\substack{\lambda \in \Lambda \\ w \in \mathcal{W}}} \left\{ \mathcal{L}(w^A, \lambda) - \mathcal{L}(w, \lambda^A) \right\} \\
 &\leq \frac{1}{T} \max_{\substack{\lambda \in \Lambda \\ w \in \mathcal{W}}} \left\{ \sum_{t=1}^T \mathcal{L}(w_t, \lambda) - \mathcal{L}(w, \lambda_t) \right\}. && \text{(convexity in } w \text{ and lin. in } \lambda)
 \end{aligned}$$

Next, since the function is linear in λ and convex in w ,

$$\begin{aligned}
 \mathcal{L}(w_t, \lambda) - \mathcal{L}(w, \lambda_t) &= \mathcal{L}(w_t, \lambda) - \mathcal{L}(w_t, \lambda_t) + \mathcal{L}(w_t, \lambda_t) - \mathcal{L}(w, \lambda_t) \\
 &\leq (\lambda - \lambda_t) \nabla_{\lambda} \mathcal{L}(w_t, \lambda_t) + (w_t - w) \nabla_w \mathcal{L}(w_t, \lambda_t) \\
 &\leq (\lambda - \lambda_t) \delta_{\lambda} \mathcal{L}(w_t, \lambda_t) + (w_t - w) \delta_w \mathcal{L}(w_t, \lambda_t) \\
 &\quad + (\lambda - \lambda_t) (\nabla_{\lambda} \mathcal{L}(w_t, \lambda_t) - \delta_{\lambda} \mathcal{L}(w_t, \lambda_t)) + (w_t - w) (\nabla_w \mathcal{L}(w_t, \lambda_t) - \delta_w \mathcal{L}(w_t, \lambda_t)).
 \end{aligned}$$

In view of these inequalities, by the subadditivity of max, the following inequality holds:

$$\begin{aligned}
 &\max_{\substack{\lambda \in \Lambda \\ w \in \mathcal{W}}} \left\{ \sum_{t=1}^T \mathcal{L}(w_t, \lambda) - \mathcal{L}(w, \lambda_t) \right\} \\
 &\leq \max_{\substack{\lambda \in \Lambda \\ w \in \mathcal{W}}} \sum_{t=1}^T (\lambda - \lambda_t) \delta_{\lambda} \mathcal{L}(w_t, \lambda_t) + (w_t - w) \delta_w \mathcal{L}(w_t, \lambda_t) \\
 &\quad + \max_{\substack{\lambda \in \Lambda \\ w \in \mathcal{W}}} \sum_{t=1}^T \lambda (\nabla_{\lambda} \mathcal{L}(w_t, \lambda_t) - \delta_{\lambda} \mathcal{L}(w_t, \lambda_t)) - w (\nabla_w \mathcal{L}(w_t, \lambda_t) - \delta_w \mathcal{L}(w_t, \lambda_t)) \\
 &\quad + \sum_{t=1}^T \lambda_t (\nabla_{\lambda} \mathcal{L}(w_t, \lambda_t) - \delta_{\lambda} \mathcal{L}(w_t, \lambda_t)) - w_t (\nabla_w \mathcal{L}(w_t, \lambda_t) - \delta_w \mathcal{L}(w_t, \lambda_t)).
 \end{aligned}$$

We now bound each of the terms above separately. For the first term, observe that for any $w \in \mathcal{W}$,

$$\begin{aligned}
 & \sum_{t=1}^T (w_t - w) \delta_w \mathbf{L}(w_t, \lambda_t) \\
 &= \frac{1}{2\gamma_w} \sum_{t=1}^T \left(\|w_t - w\|_2^2 + \gamma_w^2 \|\delta_w \mathbf{L}(w_t, \lambda_t)\|_2^2 - \|(w_t - \gamma_w \delta_w \mathbf{L}(w_t, \lambda_t) - w)\|_2^2 \right) \\
 &\leq \frac{1}{2\gamma_w} \sum_{t=1}^T \left(\|w_t - w\|_2^2 + \gamma_w^2 \|\delta_w \mathbf{L}(w_t, \lambda_t)\|_2^2 - \|(w_{t+1} - w)\|_2^2 \right) \quad (\text{property of projection}) \\
 &= \frac{1}{2\gamma_w} \left(\|w_1 - w\|_2^2 - \|w_{T+1} - w\|_2^2 \right) + \frac{\gamma_w}{2} \sum_{t=1}^T \|\delta_w \mathbf{L}(w_t, \lambda_t)\|_2^2 \quad (\text{telescoping sum}) \\
 &\leq \frac{1}{2\gamma_w} \|w_1 - w\|_2^2 + \frac{\gamma_w}{2} \sum_{t=1}^T \|\delta_w \mathbf{L}(w_t, \lambda_t)\|_2^2 \\
 &\leq \frac{2R_{\mathcal{W}}^2}{\gamma_w} + \frac{\gamma_w}{2} \sum_{t=1}^T \|\delta_w \mathbf{L}(w_t, \lambda_t)\|_2^2 \\
 &\leq \frac{2R_{\mathcal{W}}^2}{\gamma_w} + \frac{\gamma_w}{2} \sum_{t=1}^T \|\delta_w \mathbf{L}(w_t, \lambda_t) - \nabla_w \mathbf{L}(w_t, \lambda_t) + \nabla_w \mathbf{L}(w_t, \lambda_t)\|_2^2.
 \end{aligned}$$

Since the right-hand side does not depend on w , taking the maximum of both sides over $w \in \mathcal{W}$ and the expectation yields

$$\mathbb{E} \left[\max_{w \in \mathcal{W}} \sum_{t=1}^T (w_t - w) \delta_w \mathbf{L}(w_t, \lambda_t) \right] \leq \frac{2R_{\mathcal{W}}^2}{\gamma_w} + \frac{\gamma_w T \sigma_w^2}{2} + \frac{T \gamma_w G_w^2}{2},$$

using the following identity:

$$\begin{aligned}
 & \mathbb{E} \left[\|\delta_w \mathbf{L}(w_t, \lambda_t) - \nabla_w \mathbf{L}(w_t, \lambda_t) + \nabla_w \mathbf{L}(w_t, \lambda_t)\|_2^2 \right] \\
 &= \mathbb{E} \left[\|\delta_w \mathbf{L}(w_t, \lambda_t) - \nabla_w \mathbf{L}(w_t, \lambda_t)\|_2^2 \right] - 2 \mathbb{E} \left[\delta_w \mathbf{L}(w_t, \lambda_t) - \nabla_w \mathbf{L}(w_t, \lambda_t) \right] \cdot \nabla_w \mathbf{L}(w_t, \lambda_t) + \|\nabla_w \mathbf{L}(w_t, \lambda_t)\|_2^2 \\
 &= \mathbb{E} \left[\|\delta_w \mathbf{L}(w_t, \lambda_t) - \nabla_w \mathbf{L}(w_t, \lambda_t)\|_2^2 \right] + \|\nabla_w \mathbf{L}(w_t, \lambda_t)\|_2^2.
 \end{aligned}$$

Similarly, using the projection property, the following inequality can be shown:

$$\mathbb{E} \left[\max_{\lambda \in \Lambda} \sum_{t=1}^T (\lambda - \lambda_t) \delta_\lambda \mathbf{L}(w_t, \lambda_t) \right] \leq \frac{2R_{\Lambda}^2}{\gamma_\lambda} + \frac{\gamma_\lambda T \sigma_\lambda^2}{2} + \frac{T \gamma_\lambda G_\lambda^2}{2}.$$

For the second term, by the Cauchy-Schwarz inequality, we can write

$$\begin{aligned}
 \max_{\lambda \in \Lambda} \sum_{t=1}^T \lambda (\nabla_\lambda \mathbf{L}(w_t, \lambda_t) - \delta_\lambda \mathbf{L}(w_t, \lambda_t)) &\leq R_{\Lambda} \left\| \sum_{t=1}^T \nabla_\lambda \mathbf{L}(w_t, \lambda_t) - \delta_\lambda \mathbf{L}(w_t, \lambda_t) \right\|_2 \\
 &\leq R_{\Lambda} \sum_{t=1}^T \|\nabla_\lambda \mathbf{L}(w_t, \lambda_t) - \delta_\lambda \mathbf{L}(w_t, \lambda_t)\|_2.
 \end{aligned}$$

Taking the expectation of both sides and using Jensen's inequality yields

$$\mathbb{E} \left[\max_{\lambda \in \Lambda} \sum_{t=1}^T \lambda (\nabla_\lambda \mathbf{L}(w_t, \lambda_t) - \delta_\lambda \mathbf{L}(w_t, \lambda_t)) \right] \leq R_{\Lambda} \sqrt{T} \sigma_\lambda.$$

Similarly, we obtain the following:

$$\mathbb{E} \left[\max_{w \in \mathcal{W}} w \nabla_w \mathcal{L}(w_t, \lambda_t) - \delta_w \mathcal{L}(w_t, \lambda_t) \right] \leq R_{\mathcal{W}} \sqrt{T} \sigma_w.$$

For the third term, observe that the stochastic gradients at time t are unbiased, conditioned on λ_t , and w_t , hence,

$$\mathbb{E} \left[\sum_{t=1}^T \lambda_t (\nabla_{\lambda} \mathcal{L}(w_t, \lambda_t) - \delta_{\lambda} \mathcal{L}(w_t, \lambda_t)) - w_t (\nabla_w \mathcal{L}(w_t, \lambda_t) - \delta_w \mathcal{L}(w_t, \lambda_t)) \right] = 0.$$

Combining the upper bounds just derived gives:

$$\begin{aligned} \mathbb{E} \left[\max_{\lambda \in \Lambda} \mathcal{L}(w^A, \lambda) - \min_{w \in \mathcal{W}} \max_{\lambda \in \Lambda} \mathcal{L}(w, \lambda) \right] \\ \leq \frac{2R_{\mathcal{W}}^2}{T\gamma_w} + \frac{\gamma_w(\sigma_w^2 + G_w^2)}{2} + \frac{2R_{\Lambda}^2}{T\gamma_{\lambda}} + \frac{\gamma_{\lambda}(\sigma_{\lambda}^2 + G_{\lambda}^2)}{2} + \frac{R_{\mathcal{W}}\sigma_w}{\sqrt{T}} + \frac{R_{\Lambda}\sigma_{\lambda}}{\sqrt{T}}. \end{aligned}$$

Setting $\gamma_w = \frac{2R_{\mathcal{W}}}{\sqrt{T(\sigma_w^2 + G_w^2)}}$ and $\gamma_{\lambda} = \frac{2R_{\Lambda}}{\sqrt{T(\sigma_{\lambda}^2 + G_{\lambda}^2)}}$ to minimize this upper bound completes the proof. \blacksquare

5.4. Stochastic gradients

The convergence results of Theorem 5 depend on the variance of the stochastic gradients. Thus, before proceeding to the results, we first compute the gradients with respect to w and λ . Let $\mathcal{L}_{k,i}(w) = \ell(h(x_{k,i}, y_{k,i}))$. For any $w \in \mathcal{W}$, $\lambda \in \Lambda$ and $k \in [p]$, the gradient with respect to w is given by

$$\nabla_w \mathcal{L}(w, \lambda) = \sum_{k=1}^p \frac{\lambda_k}{m_k} \sum_{i=1}^{m_k} \nabla_w \mathcal{L}_{k,i}(w).$$

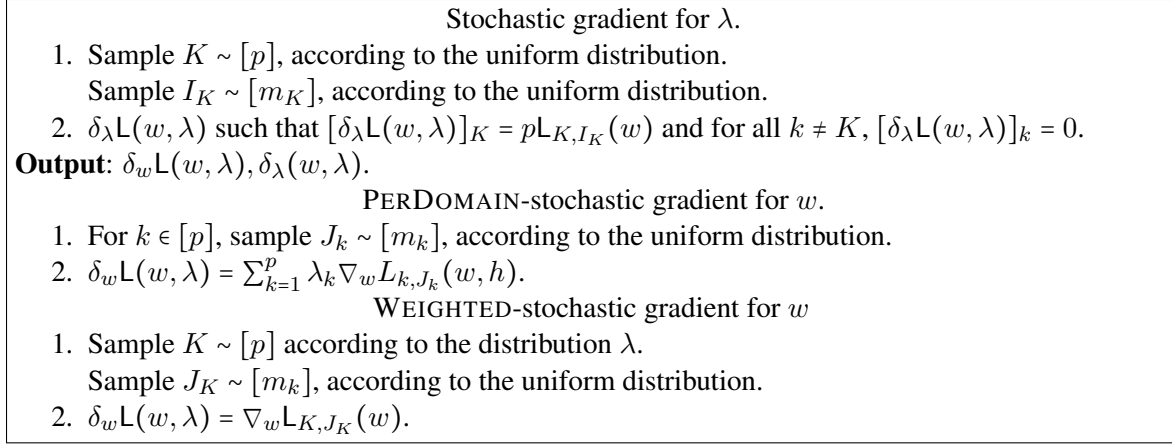
For any $w \in \mathcal{W}$, $\lambda \in \Lambda$ and $k \in [p]$, the gradient with respect to λ_k is given by

$$[\nabla_{\lambda} \mathcal{L}(w, \lambda)]_k = \frac{1}{m_k} \sum_{i=1}^{m_k} \mathcal{L}_{k,i}(w) = \mathcal{L}_k(w).$$

We first discuss the stochastic gradients for λ . Notice that the gradient for λ is independent of λ . Thus, a natural choice for the stochastic gradient with respect to λ is based on uniformly sampling a domain $K \in [p]$ and then sampling $x_{K,i}$ from domain K . This leads to the definition of the stochastic gradient $\delta_{\lambda} \mathcal{L}(w, \lambda)$ shown in Figure 4. The following lemma bounds the variance for that definition of $\delta_{\lambda} \mathcal{L}(w, \lambda)$.

Lemma 6 *The stochastic gradient $\delta_{\lambda} \mathcal{L}(w, \lambda)$ is unbiased. Further, if the loss function is bounded by M , then the following upper bound holds for the variance of $\delta_{\lambda} \mathcal{L}(w, \lambda)$:*

$$\sigma_{\lambda}^2 = \max_{w \in \mathcal{W}, \lambda \in \Lambda} \text{Var}(\delta_{\lambda} \mathcal{L}(w, \lambda)) \leq p^2 M^2.$$


 Figure 4: Definition of the stochastic gradients with respect to λ and w .

Proof The unbiasedness of $\delta_\lambda \mathbf{L}(w, \lambda)$ follows directly its definition. For the variance, observe that, for index $k \in [p]$, since the probability of not drawing domain k is $(1 - \frac{1}{p})$, the variance is given by the following

$$\begin{aligned} \text{Var}_k[\delta_\lambda \mathbf{L}(w, \lambda)] &= \left[1 - \frac{1}{p}\right] [0 - \mathbf{L}_k(w)]^2 + \frac{1}{p} \sum_{k=1}^p \frac{1}{m_k} \sum_{i=1}^{m_k} [p \mathbf{L}_{k,i}(w) - \mathbf{L}_k(w)]^2 \\ &\leq \left[1 - \frac{1}{p}\right] M^2 + \frac{1}{p} \sum_{k=1}^p \frac{1}{m_k} \sum_{i=1}^{m_k} [pM]^2 = pM^2. \end{aligned}$$

Summing over all indices from $k \in [p]$ completes the proof. ■

If the above variance is too high, then we can sample one J_k for every domain k . This is same as computing the gradient of a batch and reduces the variance by a factor of p .

The gradient with respect to w depends both on λ and w . There are two natural stochastic gradients: the PERDOMAIN-stochastic gradient and the WEIGHTED-stochastic gradient. For a PerDomain-stochastic gradient, we sample an element uniformly from $[m_k]$ for each $k \in [p]$. For the WEIGHTED-stochastic gradient, we sample a domain according to λ and sample an element out of it. To bound the variance of these two stochastic gradients, we need a few definitions.

Definition 7 *The following definitions are used:*

- the intra-domain variance with respect to w is defined as follows:

$$\sigma_I^2(w) = \max_{w \in \mathcal{W}, k \in [p]} \frac{1}{m_k} \sum_{j=1}^{m_k} [\nabla_w L_{k,j}(w) - \nabla_w L_k(w)]^2.$$

- the outer-domain variance with respect to w is defined as follows:

$$\sigma_O^2(w) = \max_{w \in \mathcal{W}, \lambda \in \Lambda} \sum_{k=1}^p \lambda_k [\nabla_w \mathbf{L}_k(w) - \nabla_w \mathbf{L}(w, \lambda)]^2.$$

- the time complexity of computing the loss and gradient with respect to w for a single sample is denoted by U .

With these definitions, we can bound the variance of both PERDOMAIN and WEIGHTED stochastic gradients.

Lemma 8 PERDOMAIN stochastic gradient is unbiased and runs in time $pU + \mathcal{O}(p \log m)$ and the variance satisfy,

$$\sigma_w^2 \leq R_\Lambda \sigma_I^2(w).$$

Proof The time complexity and the unbiasedness follow from the definitions. We now bound the variance. Since $\nabla_w \mathcal{L}_{k,J_k}$ is an unbiased estimate of $\nabla_w \mathcal{L}_k(w)$ and we have:

$$\text{Var}[\delta_w] = \sum_{k=1}^p \lambda_k^2 \text{Var}[\nabla_w \mathcal{L}_{k,J_k}(w) - \nabla_w \mathcal{L}_k(w)] \leq \sum_{k=1}^p \lambda_k^2 \sigma^2(w, I) \leq R_\Lambda \sigma_I^2(w).$$

This completes the proof. ■

Lemma 9 WEIGHTED stochastic gradient is unbiased and runs in time $U + \mathcal{O}(k + \log n)$ and the variance satisfy the following inequality:

$$\sigma_w^2 \leq \sigma_I^2(w) + \sigma_O^2(w).$$

Proof The time complexity and the unbiasedness follow from the definitions. We now bound the variance. By definition for any w, λ ,

$$\begin{aligned} \text{Var}(\delta_w) &= \sum_{k=1}^p \frac{\lambda_k}{m_k} \sum_{j=1}^{m_k} (\nabla_w \mathcal{L}_{k,j}(w) - \mathcal{L}(w, \lambda))^2 \\ &= \sum_{k=1}^p \frac{\lambda_k}{m_k} \sum_{j=1}^{m_k} (\nabla_w \mathcal{L}_{k,j}(w) - \mathcal{L}_k(w))^2 + \sum_{k=1}^p \lambda_k (\mathcal{L}_k(w) - \mathcal{L}(w, h))^2 \\ &\leq \sigma_I^2(w) + \sigma_O^2(w), \end{aligned}$$

where the second equality follows from the unbiasedness of the stochastic gradients. ■

Since $R_\Lambda \leq 1$, at first glance, the above two lemmas may suggest that PERDOMAIN stochastic is always better than WEIGHTED stochastic gradient. Note, however, that the time complexities of the algorithms is dominated by U and thus, the time complexity of PERDOMAIN-stochastic gradient is roughly k times larger than that of WEIGHTED-stochastic gradient. Hence, if k is small, it is preferable to choose the PERDOMAIN-stochastic gradient.

For large values of p , to do a fair comparison, we need to average p independent copies of the WEIGHTED-stochastic gradient, which we refer to as p -WEIGHTED, and compare it with the PERDOMAIN-stochastic gradient. Since the variance of average of p i.i.d. random variables is $1/p$ times the individual variance, by Lemma 9, the following holds:

$$\text{Var}(k\text{-WEIGHTED}) = \frac{\sigma_I^2(w) + \sigma_O^2(w)}{p}.$$

Table 1: Test accuracy of the train model on various domains, as a function of training loss for the adult dataset. Of all the model, domain agnostic model that minimizes $\mathcal{L}_{\mathcal{D}_\Lambda}$ has the best accuracy on the worst domain. All experiments are averaged over 50 runs.

Training loss	\mathcal{U}	doctorate	non-doctorate	\mathcal{D}_Λ
$\mathcal{L}_{\text{doctorate}}$	53.35 ± 0.91	73.58 ± 0.48	53.12 ± 0.89	53.12 ± 0.89
$\mathcal{L}_{\text{non-doctorate}}$	82.15 ± 0.09	69.46 ± 0.29	82.29 ± 0.09	69.46 ± 0.29
$\mathcal{L}_{\widehat{\mathcal{U}}}$	82.10 ± 0.09	69.61 ± 0.35	82.24 ± 0.09	69.61 ± 0.35
$\mathcal{L}_{\mathcal{D}_\Lambda}$	80.10 ± 0.39	71.53 ± 0.88	80.20 ± 0.40	71.53 ± 0.88

Table 2: Test accuracy of the train model on the different clothing classes, as a function of training loss for the Fashion MNIST dataset. Of the two models, the domain agnostic model that minimizes $\mathcal{L}_{\mathcal{D}_\Lambda}$ has the best accuracy overall and on the worst domain. All experiments are averaged over 50 runs.

Training loss	\mathcal{U}	shirt	pullover	T-shirt/top	\mathcal{D}_Λ
$\mathcal{L}_{\widehat{\mathcal{U}}}$	81.8 ± 1.3	71.2 ± 7.8	87.8 ± 6.0	86.2 ± 4.9	71.2 ± 7.8
$\mathcal{L}_{\mathcal{D}_\Lambda}$	82.3 ± 0.9	74.5 ± 6.0	87.6 ± 4.5	84.9 ± 4.4	74.5 ± 6.0

Further, observe that $R_\Lambda = \max_{\lambda \in \Lambda} \sum_{k=1}^p \lambda_k^2 \geq \frac{1}{p}$. Thus,

$$\text{Var}(\text{PERDOMAIN}) \geq \frac{\sigma_I^2(w)}{p}.$$

Hence, the right choice of the stochastic variance of w depends on the application. If all domains are roughly equally weighted, then we have $R(\Lambda) \approx \frac{1}{p}$ and the PERDOMAIN-variance is a more favorable choice. Otherwise, if $\sigma_O^2(w)$ is small, then the WEIGHTED-stochastic gradient is more favorable.

6. Experiments

To study the benefits of our AFL algorithm, we carried out experiments with three datasets. Even though our optimization convergence guarantees hold only for convex functions and stochastic gradient, we show that our domain-agnostic learning performs well for non-convex functions and variants of stochastic gradient descent such as momentum and Adagrad too.

In all the three experiments, we compare the domain agnostic model with the model trained with $\widehat{\mathcal{U}}$, the uniform distribution over the union of samples, and the models trained on individual domains. In all of these experiments, we used PERDOMAIN stochastic gradients and set $\Lambda = \Delta_p$. All algorithms were implemented in Tensorflow (Abadi et al., 2015).

Table 3: Test perplexity of the train model on various domains, as a function of training loss for the language model dataset. Of all the model, the domain agnostic model that minimizes $\mathcal{L}_{\mathcal{D}_\Lambda}$ admits the best perplexity on the worst domain.

Training loss	\mathcal{U}	document	conversation	\mathcal{D}_Λ
$\mathcal{L}_{\text{document}}$	414.96	83.97	615.75	615.75
$\mathcal{L}_{\text{conversation}}$	108.97	1138.76	61.01	1138.76
$\mathcal{L}_{\widehat{\mathcal{U}}}$	68.18	96.98	62.50	96.98
$\mathcal{L}_{\mathcal{D}_\Lambda}$	79.98	86.33	78.48	86.33

6.1. Adult dataset

The Adult dataset is a census dataset from the UCI Machine Learning Repository (Blake, 1998). It contains 32,561 training samples with numerical and categorical features, each representing a person. The task consists of predicting if the person’s income exceeds \$50,000. We split this dataset into two domains depending on whether the person had a doctorate degree or not, resulting into domains: the `doctorate` domain containing 413 examples and the `non-doctorate` domain containing 32,148 examples. We trained a logistic regression model with just the categorical features and Adagrad optimizer. The performance of the models averaged over 50 runs is reported in Table 1. The performance on \mathcal{D}_Λ of the model trained with $\widehat{\mathcal{U}}$, that is standard federated learning, is about 69.6%. In contrast, the performance of our AFL model is at least about 71.5% on *any* target distribution \mathcal{D}_Λ . The uniform average over the domains of the test accuracy of the AFL model is slightly less than that of the uniform model, but the agnostic model is less biased and performs better on \mathcal{D}_Λ . Furthermore, of the two domains, the `doctorate` domain is the harder one for predictions. For this domain, the performance of the domain agnostic model is close to the model trained only on `doctorate` data and is better than that of the model trained with the uniform distribution $\widehat{\mathcal{U}}$.

6.2. Fashion MNIST

The Fashion MNIST dataset, originally announced by Xiao et al. (2017), is an MNIST-like dataset where images are classified into 10 categories of clothing, instead of handwritten digits. The dataset includes 60,000 training images and 10,000 test images given as 28x28 arrays of grayscale pixel intensities, spread evenly among the ten categories. We first trained a simple logistic regression classifier and observed that the lowest performance was achieved for the following three categories: `t-shirt/top`, `pullover`, and `shirt`. Next, we extracted the subset of the data labeled with these three categories and split this subset into three domains, each consisting of one class of clothing. We then trained a classifier for the three classes using logistic regression and the Adam optimizer. The results are shown in Table 2. Since here the domain uniquely identifies the label, in this experiment, we did not compare against models trained on specific domains. Of the three domains or classes, the `shirt` class is the hardest one to distinguish from others. The domain-agnostic model improves the performance for `shirt` more than it degrades it on `pullover` and `shirt`, leading to both `shirt`-specific and overall accuracy improvement when compared to the model trained with the uniform distribution $\widehat{\mathcal{U}}$. Furthermore, in this experiment, note that our agnostic learning solution not only improves the loss of the worst domain, but also generalizes better and hence improves

the average test accuracy. Our AFL model achieves a performance of about %74.5 on *any* target distribution \mathcal{D}_λ , while the performance of standard federated learning can be as low as about %71.2.

6.3. Language models

Motivated by the keyboard application (Hard et al., 2018), where a single client uses a trained language model in multiple environments such as chat apps, email, and web input, we created a dataset that combines two very different types of language datasets: `conversation` and `document`. For `conversation`, we used the Cornell movie dataset that contain movie dialogues Danescu-Niculescu-Mizil and Lee (2011). This dataset contains about 300,000 sentences with an average sentence length of 8. For `documents`, we used the Penn TreeBank (PTB) dataset that contains approximately 50,000 sentences with an average sentence length of 20 Marcus et al. (1993). We created a single dataset by combining both of the above corpuses, with `conversation` and `document` as domains. We preprocessed the data to remove punctuations, capitalized the data uniformly, and computed a vocabulary of 10,000 most frequent words. We trained a two-layer LSTM model with LSTM and projection size of 512 with momentum optimizer. The performance of the models are measured by their perplexity, that is the exponent of cross-entropy loss. The results are reported in Table 3. Of the two domains, the `document` domain is the one admitting the higher perplexity. For this domain, the test perplexity of the domain agnostic model is close to that of the model trained only on `document` data and is better than that of the model trained with the uniform distribution $\hat{\mathcal{U}}$.

7. Extensions

In this section, we briefly discuss several extensions of the framework, theory and algorithms that we presented.

7.1. Domain definitions

The choice of the domains can significantly impact learnability in federated learning. In view of our learning bounds, if the number of domains, p , is large and Λ is the full simplex, $\Lambda = \Delta_p$, then the models may not generalize well. Thus, if the number of clients is very large, using each client as a domain may be a poor choice for better generalization. Ideally, each domain is represented with a sufficiently large number of samples and is relatively homogeneous or pure. This suggests using a clustering algorithm for defining the domains based on the similarity of the client distributions. Different Bregman divergences could be used to define the divergence or similarity between distributions. Thus, techniques such as those of Banerjee, Merugu, Dhillon, and Ghosh (2005) could be used to determine clusters of clients using a suitable Bregman divergence.

Client clusters can also be determined based on domain expertise. For example, in federated keyboard next word prediction (Hard et al., 2018), domains can be chosen to be the native language of the clients. If the model is used in variety of applications, domains can also be based on the application of interest. For example, the keyboard in (Hard et al., 2018) is used in chat apps, social apps, and

web inputs. Here, domains can be the app that was used. Training models agnostically ensures that the user experience is favorable in all apps.

7.2. Incorporating a prior on Λ

Agnostic federated learning as defined in (1) treats all domains equally and does not incorporate any prior knowledge of λ . Suppose we have a prior distribution $p_\Lambda(\lambda)$ over $\lambda \in \Lambda$ at our disposal, then, we can modify (1) to incorporate that prior. If the loss function ℓ is the cross-entropy loss, then the agnostic loss can be modified as follows:

$$\max_{\lambda \in \Lambda} (\mathcal{L}_{D_\lambda}(h) + \log p_\Lambda(\lambda)). \quad (9)$$

In this formulation, larger weights are assigned to more likely domains. The generalization guarantees of Theorem 2 can be appropriately modified to include these changes. Furthermore, if the prior $p_\Lambda(\lambda)$ is a log-concave function of λ , then the new objective is convex in h and concave in λ and a slight modification of our proposed algorithm can be used to determine the global minima. We note that we could also adopt a multiplicative formulation with the prior multiplying the loss, instead of the additive one with the negative log of the probability in Equation 9.

7.3. Domain features and personalization

We studied agnostic federated learning, where we learn a model that performs well on all domains. First, notice that we do not make any assumption on the hypothesis set \mathcal{H} and the hypotheses can use the domain k as a feature. Such models could be useful for applications where the target domain is known at inference time. Second, while the paper deals with learning a centralized model, the resulting model $h_{\mathcal{D}_\Lambda}$ can be combined with a personalized model, on the client's machine, to design better client-specific models. This can be done for example by learning an appropriate mixture weight $\alpha_k \in [0, 1]$ to use a mixture $\alpha_k h_{\mathcal{D}_\Lambda} + (1 - \alpha_k) h_k$ of the domain agnostic centralized model $h_{\mathcal{D}_\Lambda}$ and a client- or domain-specific model h_k .

8. Conclusion

We introduced a new framework of AFL for which we presented a detailed theoretical analysis. We also gave an algorithm for this problem benefiting from our theoretical analysis, as well as a new stochastic optimization solution needed for large-scale problems. Our experimental results suggest that our solution can lead to significant benefits in practice.

9. Acknowledgements

We thank Shankar Kumar, Rajiv Mathews, and Brendan McMahan for helpful comments and discussions.

References

- Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. URL <https://www.tensorflow.org/>. Software available from tensorflow.org.
- Naman Agarwal, Ananda Theertha Suresh, Felix X. Yu, Sanjiv Kumar, and Brendan McMahan. cpSGD: Communication-efficient and differentially-private distributed SGD. In *Proceedings of NeurIPS*, pages 7575–7586, 2018.
- Arindam Banerjee, Srujana Merugu, Inderjit S Dhillon, and Joydeep Ghosh. Clustering with Bregman divergences. *Journal of machine learning research*, 6(Oct):1705–1749, 2005.
- Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations for domain adaptation. In *NIPS*, pages 137–144, 2006.
- P. J. Bickel, E. A. Hammel, and J. W. O’Connell. Sex bias in graduate admissions: Data from Berkeley. *Science*, 187(4175):398–404, 1975. ISSN 0036-8075.
- Catherine L Blake. UCI repository of machine learning databases, Irvine, University of California. <http://www.ics.uci.edu/~mlearn/MLRepository>, 1998.
- John Blitzer, Mark Dredze, and Fernando Pereira. Biographies, Bollywood, Boom-boxes and Blenders: Domain Adaptation for Sentiment Classification. In *Proceedings of ACL 2007*, Prague, Czech Republic, 2007.
- Corinna Cortes and Mehryar Mohri. Domain adaptation and sample bias correction theory and algorithm for regression. *Theor. Comput. Sci.*, 519:103–126, 2014.
- Corinna Cortes, Mehryar Mohri, and Andres Muñoz Medina. Adaptation algorithm and theory based on generalized discrepancy. In *KDD*, pages 169–178, 2015.
- Cristian Danescu-Niculescu-Mizil and Lillian Lee. Chameleons in imagined conversations: A new approach to understanding coordination of linguistic style in dialogs. In *Proceedings of the 2nd Workshop on Cognitive Modeling and Computational Linguistics*, pages 76–87. Association for Computational Linguistics, 2011.
- Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng. Training GANs with optimism. *arXiv preprint arXiv:1711.00141*, 2017.
- Mark Dredze, John Blitzer, Pratha Pratim Talukdar, Kuzman Ganchev, Joao Graca, and Fernando Pereira. Frustratingly Hard Domain Adaptation for Parsing. In *Proceedings of CoNLL 2007*, Prague, Czech Republic, 2007.

- Farzan Farnia and David Tse. A minimax approach to supervised learning. In *Proceedings of NIPS*, pages 4240–4248, 2016.
- Yaroslav Ganin and Victor S. Lempitsky. Unsupervised domain adaptation by backpropagation. In *ICML*, volume 37, pages 1180–1189, 2015.
- Jean-Luc Gauvain and Chin-Hui. Maximum a posteriori estimation for multivariate gaussian mixture observations of Markov chains. *IEEE Transactions on Speech and Audio Processing*, 2(2):291–298, 1994.
- Ross B. Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, pages 580–587, 2014.
- Boqing Gong, Yuan Shi, Fei Sha, and Kristen Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *CVPR*, pages 2066–2073, 2012.
- Boqing Gong, Kristen Grauman, and Fei Sha. Connecting the dots with landmarks: Discriminatively learning domain-invariant features for unsupervised domain adaptation. In *ICML*, volume 28, pages 222–230, 2013a.
- Boqing Gong, Kristen Grauman, and Fei Sha. Reshaping visual datasets for domain adaptation. In *NIPS*, pages 1286–1294, 2013b.
- Peter D. Grünwald. *The minimum description length principle*. MIT press, 2007.
- Andrew Hard, Kanishka Rao, Rajiv Mathews, Françoise Beaufays, Sean Augenstein, Hubert Eichner, Chloé Kiddon, and Daniel Ramage. Federated learning for mobile keyboard prediction. *arXiv preprint arXiv:1811.03604*, 2018.
- Moritz Hardt, Eric Price, Nati Srebro, et al. Equality of opportunity in supervised learning. In *Proceedings of NIPS*, pages 3315–3323, 2016.
- Judy Hoffman, Brian Kulis, Trevor Darrell, and Kate Saenko. Discovering latent domains for multisource domain adaptation. In *ECCV*, volume 7573, pages 702–715, 2012.
- Judy Hoffman, Erik Rodner, Jeff Donahue, Kate Saenko, and Trevor Darrell. Efficient learning of domain-invariant image representations. In *ICLR*, 2013.
- Judy Hoffman, Mehryar Mohri, and Ningshan Zhang. Algorithms and theory for multiple-source adaptation. In *Proceedings of NeurIPS*, pages 8256–8266, 2018.
- Frederick Jelinek. *Statistical Methods for Speech Recognition*. The MIT Press, 1998.
- Jing Jiang and ChengXiang Zhai. Instance Weighting for Domain Adaptation in NLP. In *Proceedings of ACL 2007*, pages 264–271, Prague, Czech Republic, 2007. Association for Computational Linguistics.
- Anatoli Juditsky, Arkadi Nemirovski, and Claire Tauvel. Solving variational inequalities with stochastic mirror-prox algorithm. *Stochastic Systems*, 1(1):17–58, 2011.

- Vladimir Koltchinskii and Dmitry Panchenko. Empirical margin distributions and bounding the generalization error of combined classifiers. *Annals of Statistics*, 30, 2002.
- Jakub Konečný, H Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated optimization: Distributed machine learning for on-device intelligence. *arXiv preprint arXiv:1610.02527*, 2016a.
- Jakub Konečný, H Brendan McMahan, Felix X Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*, 2016b.
- Jaeho Lee and Maxim Raginsky. Minimax statistical learning and domain adaptation with Wasserstein distances. *arXiv preprint arXiv:1705.07815*, 2017.
- C. J. Legetter and Phil C. Woodland. Maximum likelihood linear regression for speaker adaptation of continuous density hidden Markov models. *Computer Speech and Language*, pages 171–185, 1995.
- Jianwei Liu, Jiajia Zhou, and Xionglin Luo. Multiple source domain adaptation: A sharper bound using weighted Rademacher complexity. In *Technologies and Applications of Artificial Intelligence (TAAI), 2015 Conference on*, pages 546–553. IEEE, 2015.
- Mingsheng Long, Yue Cao, Jianmin Wang, and Michael I. Jordan. Learning transferable features with deep adaptation networks. In *ICML*, volume 37, pages 97–105, 2015.
- Yishay Mansour, Mehryar Mohri, and Afshin Rostamizadeh. Multiple source adaptation and the Rényi divergence. In *UAI*, pages 367–374, 2009a.
- Yishay Mansour, Mehryar Mohri, and Afshin Rostamizadeh. Domain adaptation: Learning bounds and algorithms. In *COLT*, 2009b.
- Yishay Mansour, Mehryar Mohri, and Afshin Rostamizadeh. Domain adaptation with multiple sources. In *NIPS*, pages 1041–1048, 2009c.
- Mitchell P Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. Building a large annotated corpus of english: The penn treebank. *Computational linguistics*, 19(2):313–330, 1993.
- Aleix M. Martínez. Recognizing imprecisely localized, partially occluded, and expression variant faces from a single sample per class. *IEEE Trans. Pattern Anal. Mach. Intell.*, 24(6):748–763, 2002.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Proceedings of AISTATS*, pages 1273–1282, 2017.
- Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. MIT Press, second edition, 2018.
- Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant feature representation. In *ICML*, volume 28, pages 10–18, 2013.

- Arkadii Semenovich Nemirovski and David Berkovich Yudin. *Problem complexity and Method Efficiency in Optimization*. Wiley, 1983.
- Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Trans. Knowl. Data Eng.*, 22(10):1345–1359, 2010.
- S. Della Pietra, V. Della Pietra, R. L. Mercer, and S. Roukos. Adaptive language modeling using minimum discriminant estimation. In *HLT '91: Proceedings of the workshop on Speech and Natural Language*, pages 103–106, Morristown, NJ, USA, 1992. Association for Computational Linguistics.
- Anirudh Raju, Behnam Hedayatnia, Linda Liu, Ankur Gandhe, Chandra Khatri, Angeliki Metallinou, Anu Venkatesh, and Ariya Rastrow. Contextual language model adaptation for conversational agents. *arXiv preprint arXiv:1806.10215*, 2018.
- Sasha Rakhlin and Karthik Sridharan. Optimization, learning, and games with predictable sequences. In *Proceedings of NIPS*, pages 3066–3074, 2013.
- Brian Roark and Michiel Bacchiani. Supervised and unsupervised PCFG adaptation to novel domains. In *Proceedings of HLT-NAACL*, 2003.
- Roni Rosenfeld. A Maximum Entropy Approach to Adaptive Statistical Language Modeling. *Computer Speech and Language*, 10:187–228, 1996.
- Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *ECCV*, volume 6314, pages 213–226, 2010.
- Virginia Smith, Chao-Kai Chiang, Maziar Sanjabi, and Ameet S. Talwalkar. Federated multi-task learning. In *Proceedings of NIPS*, pages 4427–4437, 2017.
- Ananda Theertha Suresh, Felix X Yu, Sanjiv Kumar, and H Brendan McMahan. Distributed mean estimation with limited communication. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3329–3337. JMLR. org, 2017.
- Eric Tzeng, Judy Hoffman, Trevor Darrell, and Kate Saenko. Simultaneous deep transfer across domains and tasks. In *ICCV*, pages 4068–4076, 2015.
- Weiran Wang and Miguel A Carreira-Perpinán. Projection onto the probability simplex: An efficient algorithm with a simple proof, and an application. *arXiv preprint arXiv:1309.1541*, 2013.
- Blake E. Woodworth, Jialei Wang, Adam D. Smith, Brendan McMahan, and Nati Srebro. Graph oracle models, lower bounds, and gaps for parallel stochastic optimization. In *Proceedings of NeurIPS*, pages 8505–8515, 2018.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. *CoRR*, abs/1708.07747, 2017. URL <http://arxiv.org/abs/1708.07747>.
- Zheng Xu, Wen Li, Li Niu, and Dong Xu. Exploiting low-rank structure from latent domains for domain generalization. In *ECCV*, volume 8691, pages 628–643, 2014.

Jun Yang, Rong Yan, and Alexander G. Hauptmann. Cross-domain video concept detection using adaptive svms. In *ACM Multimedia*, pages 188–197, 2007.

Kun Zhang, Mingming Gong, and Bernhard Schölkopf. Multi-source domain adaptation: A causal view. In *AAAI*, pages 3150–3157, 2015.

Appendix A. Alternative learning guarantees

An objective similar to that of AFL was considered in the context of multiple source domain adaptation by [Liu et al. \(2015\)](#). The authors presented generalization bounds for a scenario where the target is based on some specific mixture λ of the source domains. Our theoretical results differ from those of this work in two ways. First, our generalization bounds do not hold for a single mixture weight λ but for any subset Λ of the simplex. Second, the complexity terms in the bounds presented by these authors are proportional to $\sqrt{m} \max_{k \in [p]} \frac{\lambda_k}{m_k}$, while our guarantees are in terms of $\sqrt{\sum_{k=1}^p \frac{\lambda_k^2}{m_k}}$, which is strictly tighter. In particular, in the special case where $k = 2$, $\lambda_1 = \frac{1}{\sqrt{m}}$, $\lambda_2 = 1 - \lambda_1$ and $m_1 = 1$ and $m_2 = m - 1$, the bounds of [Liu et al. \(2015\)](#) are proportional to a constant and thus not informative, $\sqrt{m} \max_{k \in [p]} \frac{\lambda_k}{m_k} = 1$, while our guarantees are in terms of $\frac{1}{\sqrt{m}}$.

Our generalization error in Theorem 2 is particularly useful when Λ is a strict subset of the simple, $\Lambda \subset \Delta_p$. If $\Lambda = \Delta_p$, we can give the following alternative learning guarantee based.

Theorem 10 *For any $\delta > 0$, with probability at least $1 - \delta$ over the draw of samples $S_k \sim \mathcal{D}_k^{m_k}$, the following inequality holds for all $h \in \mathcal{H}$ and $\lambda \in \Lambda$:*

$$L_{\mathcal{D}_\lambda}(h) \leq L_{\overline{\mathcal{D}}_\lambda}(h) + \sum_{k=1}^p \left(2\lambda_k \mathfrak{R}_{m_k}^k(\mathcal{G}) + \lambda_k M \sqrt{\frac{1}{2m_k} \log \frac{p}{\delta}} \right),$$

where $\mathfrak{R}_{m_k}^k(\mathcal{G})$ is the Rademacher complexity over domain \mathcal{D}_k with m_k samples.

The proof is a direct application of known Rademacher complexity bounds ([Mohri et al., 2018](#)) and a union bound and is omitted.

To relate the generalization bounds of Theorem 2 and Theorem 10, observe that, by the sub-additivity of sup and the linearity of expectation, the following inequality holds:

$$\begin{aligned} \mathfrak{R}_{\mathbf{m}}(\mathcal{G}, \lambda) &= \mathbb{E}_{S_k \sim \mathcal{D}_k^{m_k}} \left[\sup_{h \in \mathcal{H}} \sum_{k=1}^p \frac{\lambda_k}{m_k} \sum_{i=1}^{m_k} \sigma_{k,i} \ell(h(x_{k,i}), y_{k,i}) \right] \\ &\leq \sum_{k=1}^p \frac{\lambda_k}{m_k} \mathbb{E}_{S_k \sim \mathcal{D}_k^{m_k}} \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^{m_k} \sigma_{k,i} \ell(h(x_{k,i}), y_{k,i}) \right] \\ &= \sum_{k=1}^p \lambda_k \mathfrak{R}_{m_k}^k(\mathcal{G}). \end{aligned}$$

Furthermore, by the sub-additivity of $\sqrt{\cdot}$, the following inequality holds:

$$\sqrt{\frac{\mathfrak{s}(\lambda \parallel \overline{\mathbf{m}})}{m}} = \sqrt{\sum_{k=1}^p \frac{\lambda_k^2}{m_k}} \leq \sum_{k=1}^p \sqrt{\frac{\lambda_k^2}{m_k}} = \sum_{k=1}^p \lambda_k \sqrt{\frac{1}{m_k}}.$$

Hence, up to the logarithmic factors in the second term, the guarantee of Theorem 2 is stronger than that of Theorem 10. However, Λ_ϵ can be large and exponential in p , and it is not clear which of the bounds are stronger in general. This depends on $\overline{\mathbf{m}}$ and λ . Deriving learning bounds that improve upon both of the learning bounds above remains an interesting open question.