# Provably Training Neural Network Classifiers under Fairness Constraints

You-Lin Chen

Department of Statistics, University of Chicago

and

Zhaoran Wang

Departments of Industrial Engineering and Management Sciences,
Northwestern University

and

Mladen Kolar

Booth School of Business, University of Chicago

January 1, 2021

**Abstract**

Training a classifier under fairness constraints has gotten increasing attention in the machine learning community thanks to moral, legal, and business reasons. However, several recent works addressing algorithmic fairness have only focused on simple models such as logistic regression or support vector machines due to non-convex and non-differentiable fairness criteria across protected groups, such as race or gender. Neural networks, the most widely used models for classification nowadays, are precluded and lack theoretical guarantees. This paper aims to fill this missing but crucial part of the literature of algorithmic fairness for neural networks. In particular, we show that overparametrized neural networks could meet the fairness constraints. The key ingredient of building a fair neural network classifier is establishing no-regret analysis for neural networks in the overparameterization regime, which may be of independent interest in the online learning of neural networks and related applications.

# 1  Introduction

In many real-world machine learning problems, practitioners are not only interested in the performance of their models, but also meeting societal and legal goals, while taking advantage of side information, prior knowledge, and unlabeled data. For example, in classification, fairness metrics with respect to certain protected groups, such as gender or ethnicity (Chouldechova, 2017), are used to correct biased training data and improve the accuracy (Blum and Stangl, 2019); F-measure, G-mean, and H-mean (Daskalaki et al., 2006; Lawrence et al., 2012; Kennedy et al., 2009) are used with class-imbalanced data to prevent trivial solutions; the classifier churn (Fard et al., 2016) is used to improve stability. All the metrics mentioned above are involved non-linear functions of a classifier's prediction rates on some subgroups or sub-population of data. One could cast this learning task of training a classification model satisfying the above metrics as a constrained optimization problem. Several challenges arise due to non-convex, non-differentiable, and data-dependent prediction rates. Recent works (Zafar et al., 2017; Donini et al., 2018; Oneto et al., 2019) have addressed challenges of algorithmic fairness partially, but only focused on simple models such as logistic regression or support vector machines. Neural networks, the most widely used models for classification nowadays, are precluded and lack theoretical guarantees.

Neural networks enjoy empirical success in practice, but their theoretical properties are less well understood due to their complexity and non-convexity. Nevertheless, many recent works (Cai et al., 2019; Fan et al., 2020; Du et al., 2019b) establish convergence results for neural networks. For example, Allen-Zhu et al. (2019a) and Arora et al. (2019) studied two-layer neural networks in the overparametrization regime, where neural networks can be approximated by linear models with the neural tangent kernel (Jacot et al., 2018; Lee et al., 2019; Alemohammad et al., 2020). This phenomenon of local linearization provides

2

a powerful tool to circumvent the obstacle of the non-convexity of neural networks and establish proof of convergence.

In this paper, we consider the problem of training a neural network classifiers under fairness constraints through the lens of the neural tangent kernel. Our results can also apply to any constraint involving prediction rates. We cast the training problem as a constrained optimization problem in which constraints are non-convex, non-differentiable, and data-dependent. Since neural network models are also non-convex, we face the problem of non-convexity arising from both the constraints and the model. To establish the rate of convergence under this challenging setting, we employ a game-theoretic framework proposed in Narasimhan et al. (2019). They solve optimization problems with non-differentiable and non-convex constraints by optimizing the Lagrangian and iteratively updating model parameters, Lagrange multipliers, and auxiliary variables. While their game-theoretic framework is rather general, it cannot be directly used to establish a fair neural network classifier and study its convergence due to the non-convexity of neural networks.

Our work fills a missing, but crucial, part of the literature of algorithmic fairness for neural networks, and the contribution is threefold. First, we prove the first convergence result for neural network classifiers under fairness constraints. Second, a no-regret guarantee for neural networks in overparameterization regime is provided, which may be of independent interest for online learning of neural networks and related applications. Finally, unlike the prior works, our approach does not require any best response function or oracle for optimization, which leads to a simple and computationally efficient procedure. Notably, projected stochastic gradient descent is employed for updating all parameters. Since stochastic gradient descent is the workhorse for training a neural network in practice, our results and proposed procedure are closer to the practical scenario and useful for

3

large-scale datasets.

## 1.1 Related Work

**Algorithmic fairness.** There is a growing body of literature studying different approaches to fulfill algorithmic fairness in machine learning (Hardt et al., 2016; Kilbertus et al., 2017; Salimi et al., 2019; Blum and Lykouris, 2019; Bechavod et al., 2019; Celis et al., 2019). Zafar et al. (2017) included the correlation between the decision boundary and sensitive attributes as a constraint on the learned classifier. Donini et al. (2018) and Oneto et al. (2019) designed constrained optimization problems that enforce the learned classifiers to have similar errors on the positive class independently of subgroups.

Several methodologies closely related to our work deal with non-convex constrained optimization problems and are further applied to achieve algorithmic fairness. Ma et al. (2019) studied constrained optimization where both the objective and constraints are weakly convex. Chen et al. (2017) and Agarwal et al. (2018) used a Lagrangian-based approach with optimization oracle and found a distribution over solutions rather than a pure equilibrium. Cotter et al. (2019b) replaced the constraints with differentiable surrogates and cast Lagrangian as a non-zero-sum two-player game. Their theoretical result guarantees feasibility concerning the original constraints instead of the proxy constraints.

**Neural tangent kernel.** There is a considerable body of literature analyzing deep supervised learning with overparameterized neural networks (Zou et al., 2018; Neyshabur et al., 2018b; Li and Liang, 2018; Du et al., 2019a). For example, Chizat et al. (2019) and Allen-Zhu et al. (2019b) showed that neural networks approximate a subset of the reproducing kernel Hilbert space induced by some kernels. We refer Fan et al. (2019) for a recent review. In comparison with prior works, we extend approximation results of neural

4

networks to the setting of online learning. Based on our new technique, we show that neural networks can be provably trained under fairness constraints.

## 1.2 Notations and Organization

For a continuous function $h : \mathbb{R}^n \to \mathbb{R}$, we denote the domain of $h$ by $\operatorname{dom} h$ and the gradient of $h$ at a point $x \in \operatorname{dom} h$ by $\nabla h(x)$. Given $x = (x_1, \ldots, x_n)$, $\nabla_{x_i} h(x)$ denote the partial derivative of $h$ corresponding to the coordinate $i$. In a stochastic setting, $\widehat{\nabla} h$ denotes an unbiased estimator of the gradient. $\mathcal{N}(\mu, \Sigma)$ is the multivariate normal distribution with the mean $\mu$ and the covariance matrix $\Sigma$, categorical$(E, P)$ is the categorical distribution such that $\mathbb{P}(e_i) = p_i$ and $\sum_i p_i = 1$ where $e_i \in E$ and $p_i \in P$, and Uniform$(E)$ is the discrete uniform distribution such that $\operatorname{Uniform}(E) = \operatorname{categorical}(E, \{1/n, \ldots, 1/n\})$ where $n$ is the number of elements in $E$. $\| \cdot \|_2, \| \cdot \|_\infty$ are defined as 2-norm and infinity norm, respectively. For a general norm $\| \cdot \|$, $\| \cdot \|_*$ is its dual norm. For a subset $S \subset \mathbb{R}^n$, $\Pi_S$ is defined as a projection to $S$ with respect to Euclidean norm $\| \cdot \|_2$. Given some sequences $a_n, b_n$, we write $a_n = O(b_n)$ if there exists a positive real number $C$ such that $a_n \le C b_n$ for all $n$.

The rest of the paper is organized as follows. Section 2 reviews essential background. We set up out problem and discuss the main results of the paper in Section 3 and Section 4, respectively. Section 5 presents numerical experiments on real-world datasets. Conclusion is provided in Section 6.

# 2 Background

This section reviews the necessary concepts of a game-theoretic framework for optimizing non-convex constrained optimization problems. The general framework introduced first can

be used to train models with fairness constraints in Section 2.1. Then we discuss online learning as well as its no-regret analysis in 2.2. Online learning is leveraged to find the equilibrium of a game, while no-regret analysis is the crucial ingredient to establish the convergence analysis. In the end, we demonstrate how to combine two concepts and meet the fairness requirement for a linear model in Section 2.3.

## 2.1 A general framework for fairness constraints

This section introduces a general framework for optimizing classification models with fairness constraints. We first set up our notations and consider a binary classification problem. Let $\mathcal{X} \subset \mathbb{R}^d$ be an instance space and $\mathcal{Z} = \{-1, 1\}$ be the label space. Denote a prediction model as $y(\theta; x) : \Theta \times \mathcal{X} \to \mathbb{R}$. Given a model $y(\theta; x)$ and some distribution $\mathcal{D}$ over $\mathcal{X} \times \mathcal{Z}$, the classification rate on $\mathcal{D}$ is defined as

$$R(\theta, \mathcal{D}) = \mathbb{E}_{(x,z) \sim \mathcal{D}} \, \mathbb{I}\{z = \text{sign}(y(\theta; x))\}, \tag{1}$$

where $\mathbb{I}$ is the indicator function and sign is the sign function. The goal of classification is minimizing the classification rate as small as possible. Note that $\mathbb{I}$ and sign are non-convex and non-differentiable and $0 \leq R_k(\theta) \leq 1$. Given a sub-population $A \subset X$ and a distribution $\mathcal{D}_A$ over $A$, $R(\theta, \mathcal{D}_A)$ denotes the classification rate on the sub-population $A$. The classifier's rate in (1) can be used to represent a large set of goals commonly used in practice, such as fairness, true-positive rate, recall, and precision, in terms of different populations.

In general, we have $K$ distributions $\mathcal{D}_1, \ldots, \mathcal{D}_K$ over $\mathcal{X} \times \mathcal{Z}$ and $J$ constraint functions $\phi_j : \mathbb{R}^K \to \mathbb{R}$ for $j = 1, \ldots, J$. Define rates $R_k(\theta) := R(\theta, \mathcal{D}_k)$. We aim to optimize the

following problem:

$$\min_{\theta} \quad g(\theta) = \mathbb{E}_{(x,z)\sim\mathcal{D}}\ell(y(\theta;x),z)$$

$$\text{s.t.} \quad h_j(\theta) = \phi_j(R_1(\theta),\ldots,R_K(\theta)) \leq 0 \qquad (2)$$

$$\text{for } j = 1,\ldots,J,$$

where $\ell$ is a convex loss function. For example, common convex loss functions include logistic or hinge loss. We denote TPR, TNR, FPR, and FNR as true-positive rate, true-negative rate, false positve rate, and false negative rate, respectively. Constraint functions $\phi_j$ in the optimization problem (2) can be used to represent a number of classifier metrics, such as the G-mean $1 - \sqrt{\text{TPR} \times \text{TNR}}$, H-mean $1 - 2/(1/\text{TPR} + 1/\text{TNR})$, and Q-mean $1 - \sqrt{\text{FPR}^2 + \text{FNR}^2}$. The optimization problem (2) can also be used to solve the unconstrained optimization problem

$$\min_{\theta} \psi(R_1(\theta),\ldots,R_K(\theta))$$

by introducing an auxiliary variable $\xi$ and reformulating the problem as

$$\min_{\theta,\xi} \quad \xi$$

$$\text{s.t.} \quad \psi(R_1(\theta),\ldots,R_K(\theta)) - \xi \leq 0. \qquad (3)$$

For example, we might have that $\psi(R_1(\theta),\ldots,R_K(\theta)) = \sum_{i=1}^{K} w_i R_i(\theta)$ and we aim to optimize the classifiers' rate with different weights—consider a scenario where the consequence of a false negative is more severe than that of a false positive, then we will put more weight on TPR and less weight on TNR.

Solving (2) requires addressing a number of difficult challenges. First, the constraint is data-dependent and may be computationally challenging to check. If the size of a dataset is huge, it is intractable to access the whole dataset to compute the classification rates

$R(\theta), R_k(\theta)$. Second, $h_j(\theta)$ is neither convex nor sub-differentiable with respect to $\theta$ for $j = 1, \ldots, J$, so (2) is not a convex problem. Moreover, we use neural networks to parametrize the prediction model, which are also non-convex.

To overcome the first and second challenge, Narasimhan et al. (2019) developed a general framework for solving a large class of learning problems with non-linear functions of classification rates as constraints, which partially addresses the above difficulties. We first introduce their framework and later build on their framework to achieve algorithmic fairness using neural networks and stochastic projected gradient descent. By introducing auxiliary variables $\xi = (\xi_1, \ldots, \xi_K)^\top$, we can rewrite (2) as

$$
\begin{aligned}
\min_{\theta \in \Theta, \xi \in \Xi} \quad & g(\theta) \\
\text{s.t.} \quad & \phi_j(\xi_1, \ldots, \xi_K) \leq 0, \ \text{for } j = 1, \ldots, J; \\
& R_k(\theta) \leq \xi_k, \ \text{for } k = 1, \ldots, K.
\end{aligned}
\tag{4}
$$

Note that we will assume $\phi_j$ is monotonically increasing in each argument so (2) and (4) are equivalent. The corresponding Lagrangian is

$$
\mathcal{L}(\theta, \xi, \lambda) = g(\theta) + \sum_{j=1}^{J} \lambda_j \phi_j(\xi) - \sum_{k=1}^{K} \lambda_{J+k} \xi_k + \sum_{k=1}^{K} \lambda_{J+k} R_k(\theta),
\tag{5}
$$

where $\lambda = (\lambda_1, \ldots, \lambda_J, \lambda_{J+1}, \ldots, \lambda_{K+J})^\top$ is the vector of the Lagrange multipliers with the $K$ auxiliary variables and $J$ constraints.

Since the functions $R_k(\theta)$ are not differentiable nor continuous, we introduce differentiable convex surrogate functions $\widetilde{R}_k$ that are upper bounds on the rates: $R_k(\theta) \leq \widetilde{R}_k(\theta)$ for all $\theta$. Despite using these surrogates, the key observation is that only $R_k(\theta)$ is needed for optimizing the Lagrangian $\mathcal{L}$ w.r.t. $\theta$. Thus, we will be able to have guarantees regarding

8

the true feasible sets. More specifically, denoting

$$\mathcal{L}_1(\xi, \lambda) = \sum_{j=1}^{J} \lambda_j \phi_j(\xi) - \sum_{k=1}^{K} \lambda_{J+k} \xi_k, \quad \widetilde{\mathcal{L}}_2(\theta, \lambda) = g(\theta) + \sum_{k=1}^{K} \lambda_{J+k} \widetilde{R}_k(\theta). \tag{6}$$

we solve (4) by alternatively minimizing $\widetilde{\mathcal{L}}_2$ with respect to $\theta$, $\mathcal{L}_1$ with respect to $\xi$, and maximizing $\mathcal{L}$ with respect to $\lambda$. This procedure can be interpreted as a game in which $\theta$-player, $\xi$-player, and $\lambda$-player have utility functions $-\widetilde{\mathcal{L}}_2$, $-\mathcal{L}_1$, and $\mathcal{L}$, respectively. To specify the domain of every variable, define

$$\begin{aligned}
\widetilde{\Theta} &= \left\{ \theta \in \Theta : (\widetilde{R}_1(\theta), \dots, \widetilde{R}_K(\theta)) \in \bigcup_{j=1}^{J} \operatorname{dom} \phi_j \right\}, \\
\Xi &= \left\{ \xi = (\xi_1, \dots, \xi_K) \in \bigcup_{j=1}^{J} \operatorname{dom} \phi_j : 0 \leq \xi_k \leq \max_{\theta \in \Theta} \widetilde{R}_k(\theta) \right\}, \\
\Lambda &= \left\{ \lambda = (\lambda_1, \dots, \lambda_J, \lambda_{J+1}, \dots, \lambda_{K+J})^\top : \|\lambda\|_\infty \leq \kappa \right\}.
\end{aligned} \tag{7}$$

The following proposition relates the equilibrium of the game and the solution of (4).

**Proposition 1.** *Suppose*

1. *$\Theta$ is compact and $\sup_{\theta \in \Theta} |g(\theta)| \leq C$;*

2. *$\phi_j$ are strictly jointly convex, monotonically increasing in each argument, and L-Lipschitz w.r.t. the infinity norm for all $j = 1, \dots, K$.*

*If $\theta_1, \dots, \theta_T$, $\xi_1, \dots, \xi_T$, and $\lambda_1, \dots, \lambda_T$ comprise an approximate coarse-correlated equilibrium, that is, satisfying,*

$$\frac{1}{T} \sum_{t=1}^{T} \mathcal{L}_1(\xi_t, \lambda_t) \leq \min_{\xi \in \Xi} \frac{1}{T} \sum_{t=1}^{T} \mathcal{L}_1(\xi, \lambda_t) + \epsilon_\xi, \tag{8}$$

9

$$\frac{1}{T}\sum_{t=1}^{T}\widetilde{\mathcal{L}}_2(\theta_t, \lambda_t) \leq \min_{\theta \in \Theta} \frac{1}{T}\sum_{t=1}^{T}\widetilde{\mathcal{L}}_2(\theta, \lambda_t) + \epsilon_\theta, \tag{9}$$

$$\frac{1}{T}\sum_{t=1}^{T}\mathcal{L}(\theta_t, \xi_t, \lambda_t) \geq \max_{\lambda \in \Lambda} \frac{1}{T}\sum_{t=1}^{T}\mathcal{L}(\theta_t, \xi_t, \lambda) - \epsilon_\lambda, \tag{10}$$

*then we have*

$$\frac{1}{T}\sum_{t=1}^{T}g(\theta_t) - \min_{\theta \in \widetilde{\Theta} \cap F} g(\theta) \leq \epsilon_\theta + \epsilon_\xi + \epsilon_\lambda \tag{11}$$

*and for all $j$,*

$$\phi_j\left(\frac{1}{T}\sum_{t=1}^{T}(R_1(\theta_t), \ldots, R_K(\theta_t))\right) \leq (L+1)(2C + \epsilon_\theta + \epsilon_\xi + \epsilon_\lambda)/\kappa. \tag{12}$$

*where $F = \{\theta : \max_j \phi_j(\widetilde{R}_1(\theta), \ldots, \widetilde{R}_K(\theta)) \leq 0\}$ is the feasible region and $\kappa$ is the upper bound of $\|\lambda\|_\infty$.*

Note that we have not discussed how to approximate coarse-correlated equilibrium yet. Proposition 1 is saying that if there exist approximate coarse-correlated equilibrium (8), (9), (10) and allow a mixed strategy or a stochastic model, we get a near-optimal and near-feasible solution to (2). In particular, if we uniformly pick an index $i \in \{1, \ldots, T\}$ and use $y(\theta_i; x)$ to make a prediction, then the corresponding random classifier's rates are

$$\mathbb{E}_i(R_1(\theta_i), \ldots, R_K(\theta_i)) = \frac{1}{T}\sum_{t=1}^{T}(R_1(\theta_t), \ldots, R_K(\theta_t))$$

and its training loss is $\mathbb{E}g(\theta_i) = (1/T)\sum_{t=1}^{T}g(\theta_t)$. The conclusions of Proposition 1 in equations (11) and (12) imply that this random classifier achieves a near-optimal and near-feasible solution. The general description of random classifier is in Algorithm 1. Such a stochastic model is necessary since the pure equilibrium may not exist in general. In

---

**Algorithm 1:** Stochastic classifier

---

**Input:** A new data $x$, a prediction model $y(\theta; x)$, a subset $S = \{\theta_s\} \subset \{\theta_t\}_{t=1}^T$ and

the probabilities $P = \{p_s\}$

Sample $\theta_s$ from categorical$(S, P)$ ;

**Output:** $y(\theta_s; x)$

---

the following section, we will introduce an online optimization procedure used to find an approximate coarse-correlated equilibrium.

We end this section by providing a proof of Proposition 1. The proof builds on ideas in Narasimhan et al. (2019), however, unlike the prior work, our approach does not require any best response function or an oracle. This leads to a simple and computationally efficient procedure. Importantly, we only need to find an approximate equilibrium for $\xi$ instead of the best response, which is needed in the framework of Narasimhan et al. (2019). As a result, we can use projected stochastic gradient descent to alternatively update $\xi, \theta, \lambda$ for large-scale datasets. Our technique is more suitable for neural network models, since stochastic gradient descent is the workhorse used to train neural networks.

We start by giving the following auxiliary results.

**Proposition 2.** *Suppose conditions in Propsotion 1 hold. For all $\lambda \geq 0$, we have*

$$\min_{\theta \in \Theta, \xi \in \Xi} \widetilde{\mathcal{L}}(\theta, \xi, \lambda) \leq \min_{\theta \in \widetilde{\Theta}: \max_j \phi_j(\widetilde{R}_1(\theta), \ldots, \widetilde{R}_K(\theta)) \leq 0} g(\theta).$$

*Proof.* For any $\lambda \geq 0$,

$$\min_{\theta \in \Theta, \xi \in \Xi} \widetilde{\mathcal{L}}(\theta, \xi, \lambda) \leq \min_{\theta \in \widetilde{\Theta}, \xi \in \Xi: \phi_j(\xi) \leq 0, \widetilde{R}_k(\theta) \leq \xi_k \forall j, k} \widetilde{\mathcal{L}}(\theta, \xi, \lambda)$$

$$\leq \min_{\theta \in \widetilde{\Theta}, \xi \in \Xi: \phi_j(\xi) \leq 0, \widetilde{R}_k(\theta) \leq \xi_k \forall j, k} g(\theta)$$

11

$$= \min_{\theta \in \widetilde{\Theta}: \max_j \phi_j(\widetilde{R}_1(\theta),\ldots,\widetilde{R}_K(\theta)) \leq 0} g(\theta),$$

where the first inequality follows from $\widetilde{\Theta} \subset \Theta$ and the equality follows from monotonicity of $\phi_j$. $\square$

**Lemma 3** (Narasimhan et al. (2019)). *Suppose conditions in Propsotion 1 hold. Then, for* $\xi$, $\Delta = (\Delta_1, \ldots, \Delta_K) \geq 0$, *we have*

$$\phi(\xi + \Delta) \leq \phi(\xi) + L \max_k \max\{\Delta_k, 0\}.$$

Now we are ready to prove Proposition 1.

*Proof of Proposition 1.* We prove Proposition 1 by showing two guarantees: optimality in (11) and feasibility in (12). We are going to utilize properties of coarse-correlated equilibrium in this proof.

Define $\bar{\lambda} = \frac{1}{T}\sum_{t=1}^{T} \lambda_t$, $\bar{\xi} = \frac{1}{T}\sum_{t=1}^{T} \xi_t$, and $\widetilde{\mathcal{L}} = \mathcal{L}_1 + \widetilde{\mathcal{L}}_2$.

**Optimality.** From (8) and (9), we know

$$\frac{1}{T}\sum_{t=1}^{T} \widetilde{\mathcal{L}}(\theta_t, \xi_t, \lambda_t) \leq \min_{\xi \in \Xi} \frac{1}{T}\sum_{t=1}^{T} \mathcal{L}_1(\xi, \lambda_t) + \epsilon_\xi + \min_{\theta \in \Theta} \frac{1}{T}\sum_{t=1}^{T} \widetilde{\mathcal{L}}_2(\theta, \lambda_t) + \epsilon_\theta$$

$$= \min_{\theta \in \Theta, \xi \in \Xi} \widetilde{\mathcal{L}}(\theta, \xi, \bar{\lambda}) + \epsilon_\theta + \epsilon_\xi$$

$$\leq \min_{\theta \in \widetilde{\Theta}: \max_j \phi_j(\widetilde{R}_1(\theta),\ldots,\widetilde{R}_K(\theta)) \leq 0} g(\theta) + \epsilon_\theta + \epsilon_\xi,$$

where the equality holds due to the linearity of $\widetilde{\mathcal{L}}$ w.r.t. $\lambda$ and the last step follows by $\bar{\lambda} > 0$ and Proposition 2. Using (10), we have for any $\lambda' \in \Lambda$

$$\frac{1}{T}\sum_{t=1}^{T} \widetilde{\mathcal{L}}(\theta_t, \xi_t, \lambda') \leq \max_{\lambda \in \Lambda} \frac{1}{T}\sum_{t=1}^{T} \widetilde{\mathcal{L}}(\theta_t, \xi_t, \lambda)$$

$$\leq \min_{\theta \in \widetilde{\Theta}: \phi_j(\widetilde{R}_1(\theta),\ldots,\widetilde{R}_K(\theta)) \leq 0 \ \forall j} g(\theta) + \epsilon_\theta + \epsilon_\xi + \epsilon_\lambda. \tag{13}$$

The optimality follows by setting $\lambda' = 0$.

**Feasibility.** Letting $j' \in \mathrm{argmax}_j\, \phi_j(\bar{\xi})$ and setting $\lambda_{j'} = \kappa$ and $\lambda_j = 0$ for $j \neq j'$, from (13) we get

$$\frac{1}{T}\sum_{t=1}^{T} g(\theta_t) + \frac{\kappa}{T}\sum_{t=1}^{T}\phi_{j'}(\xi_t) \leq \min_{\theta\in\widetilde{\Theta}:\phi_j(\widetilde{R}_1(\theta),\ldots,\widetilde{R}_K(\theta))\leq 0\forall j} g(\theta) + \epsilon_\theta + \epsilon_\xi + \epsilon_\lambda.$$

Using convexity of $\phi_j$ and the boundedness of $g$, we have

$$\max_j \phi_j(\bar{\xi}) \leq \frac{1}{T}\sum_{t=1}^{T}\phi_{j'}(\xi_t) \leq (2C + \epsilon_\theta + \epsilon_\xi + \epsilon_\lambda)/\kappa. \tag{14}$$

Similarly, letting $k' \in \mathrm{argmax}_k \left(\frac{1}{T}R_k(\theta_t) - \bar{\xi}\right)$ and setting $\lambda_{k'} = \kappa$ and $\lambda_k = 0$ for $k \neq k'$, from (13) we obtain

$$\frac{1}{T}\sum_{t=1}^{T} g(\theta_t) + \frac{\kappa}{T}\left(\sum_{t=1}^{T} R_{k'}(\theta_t) - \bar{\xi}\right) \leq \min_{\theta\in\widetilde{\Theta}:\phi_j(\widetilde{R}_1(\theta),\ldots,\widetilde{R}_K(\theta))\leq 0\forall j} g(\theta) + \epsilon_\theta + \epsilon_\xi + \epsilon_\lambda,$$

which implies

$$\max_k \max\left\{\frac{1}{T}\sum_{t=1}^{T} R_k(\theta_t) - \bar{\xi}, 0\right\} \leq (2C + \epsilon_\theta + \epsilon_\xi + \epsilon_\lambda)/\kappa. \tag{15}$$

Combining two inequality (14), (15) and the Lemma 3, we know

$$\max_j \phi_j\left(\frac{1}{T}\sum_{t=1}^{T}(R_1(\theta_t),\ldots,R_K(\theta_t))\right) \leq \max_j \phi_j(\bar{\xi}) + L\max_k \max\left\{\frac{1}{T}\sum_{t=1}^{T} R_k(\theta_t) - \bar{\xi}, 0\right\}$$

$$\leq (1 + L)(2C + \epsilon_\theta + \epsilon_\xi + \epsilon_\lambda)/\kappa,$$

which completes this proof. $\qquad\square$

## 2.2   Online learning

This section briefly presents convex online optimization problems which is the foundation of our analysis. Given a sequence of convex functions $f_1, f_2, \ldots, f_T : \Theta \to \mathbb{R}$, at each round

$t$, the learner chooses a point $\theta_t$ based on past information up to time t and incurs a loss $f_t(\theta_t)$. The goal of online learning is to control the learner's average regret in hindsight:

$$\frac{1}{T}\sum_{t=1}^{T} f_t(\theta_t) - \min_{\theta} \frac{1}{T}\sum_{t=1}^{T} f_t(\theta).$$

For example, one possible way to minimize the regret is to have the learner implement the best response strategy, that is, $\theta = \arg\min_\theta f_t(\theta)$. This strategy gives trivially negative average regret, but may be expensive in practice.

Another popular approach to control the regret is the online mirror descent (Srebro et al., 2011; Shalev-Shwartz, 2012). The mirror descent is given by the following updates

$$\zeta_{t+1} = \nabla h^*(\nabla h(\theta_t) - \eta \nabla f_t(\theta_t)), \quad \theta_{t+1} = \arg\min_{\theta \in \Theta} D_h(\theta, \zeta_{t+1}),$$

where $h$ is 1-strongly convex with respect to $\|\cdot\|$, $h^*$ is the convex conjugate of $h$, $D_h$ is the Bregman divergence with respect to $h$, defined by

$$D_h(x,y) = h(x) - h(y) - \nabla h(y)^\top (x - y),$$

and $\eta$ is the stepsize. The following well-known lemma provides a no-regret guarantee for online mirror descent.

**Lemma 4** (Srebro et al. (2011)). *Suppose that $\Theta$ is convex, $\sup_{t \leq T}\|\nabla f_t(\theta_t)\|_* < B$, and $\sup_{\theta \in \Theta} h(\theta) < M$. Then*

$$\frac{1}{T}\sum_{t=1}^{T} f_t(\theta_t) - \frac{1}{T}\sum_{t=1}^{T} f_t(\theta) \leq \frac{\eta B}{2} + \frac{M}{T\eta}.$$

*for any $\theta \in \Theta$.*

Unfortunately, Lemma 4 cannot be directly applied to our problem due to non-convexity arisen from neural networks. To analyze this more involved setting, we extend Lemma 4 to

14

a stochastic setting with biased. In particular, instead of accessing true gradients $\nabla f_t(\theta_t)$, we use stochastic mirror descent:

$$\mu_{t+1} = \nabla h^*(\nabla h(\theta_t) - \eta \zeta_t), \quad \theta_{t+1} = \arg\min_{\theta \in \Theta} D_h(\theta, \mu_{t+1}), \tag{16}$$

where $\zeta_t$ is a biased estimate of the gradient with the bias term $\beta_t(\theta_t)$, i.e.,

$$\mathbb{E}[\zeta_t \mid \theta_t] = \nabla f_t(\theta_t) + \beta_t(\theta_t).$$

**Proposition 5.** *Suppose that $\Theta$ is convex, $\sup_t \|\zeta_t\|_* < B$, and $\sup_{\theta \in \Theta} h(\theta) < M$. Given iterate $\theta_t$ updated by stochastic mirror descent in (16) with the bias term $\beta_t(\theta_t)$ , we have with probability $1 - \delta$*

$$\frac{1}{T}\sum_{t=1}^{T} f_t(\theta_t) - \frac{1}{T}\sum_{t=1}^{T} f_t(\theta) \le \frac{\eta B}{2} + \frac{M}{T\eta} + 8B\sqrt{\frac{M \ln(1/\delta)}{T}} + \frac{2\sqrt{2M}}{T}\sum_{t=1}^{T} \|\beta_t(\theta_t)\|_* \tag{17}$$

*for any $\theta \in \Theta$.*

*Proof.* Define $\widehat{f}_t(\theta) = f_t(\theta_t) + \zeta_t^\top(\theta - \theta_t)$. Since $\widehat{f}_t(\theta)$ is convex and $\nabla \widehat{f}_t(\theta) = \zeta_t$, applying Lemma 4, we have

$$\frac{1}{T}\sum_{t=1}^{T} \widehat{f}_t(\theta_t) - \frac{1}{T}\sum_{t=1}^{T} \widehat{f}_t(\theta) \le \frac{\eta B}{2} + \frac{M}{T\eta},$$

where we have used $\sup_t \|\zeta_t\|_* < B$. The result above yields

$$\frac{1}{T}\sum_{t=1}^{T} f_t(\theta_t) - \frac{1}{T}\sum_{t=1}^{T} f_t(\theta)$$

$$\le \frac{\eta B}{2} + \frac{M}{T\eta} + \frac{1}{T}\sum_{t=1}^{T} \widehat{f}_t(\theta) - \frac{1}{T}\sum_{t=1}^{T} f_t(\theta)$$

$$\le \frac{\eta B}{2} + \frac{M}{T\eta} + \frac{1}{T}\sum_{t=1}^{T} (\zeta_t - \nabla f_t(\theta_t))^\top (\theta - \theta_t)$$

15

$$= \frac{\eta B}{2} + \frac{M}{T\eta} + \frac{1}{T} \sum_{t=1}^{T} (\zeta_t - \nabla f_t(\theta_t) - \beta_t(\theta_t))^\top (\theta - \theta_t) + \frac{1}{T} \sum_{t=1}^{T} \beta_t(\theta_t)^\top (\theta - \theta_t),$$

where the last inequality follows from the convexity of $f_t$. By Holder's inequality and using the fact that $D_h$ is 1-strongly convex w.r.t. $\|\cdot\|$, we get for all $t$

$$\|\theta - \theta_t\| \leq 2 \sup_{\theta, \theta'} \sqrt{2D_h(\theta, \theta')} \leq 2\sqrt{2M}, \tag{18}$$

and

$$|(\zeta_t - \nabla f_t(\theta_t) - \beta_t(\theta_t))^\top (\theta - \theta_t)| \leq \|\zeta_t - \nabla f_t(\theta_t) - \beta_t(\theta_t)\|_* \cdot 2\sqrt{2M}$$
$$\leq (\|\zeta_t\|_* + \|\nabla f_t(\theta_t) + \beta_t(\theta_t)\|_*) \cdot 2\sqrt{2M}$$
$$\leq 4B\sqrt{2M},$$

since

$$\sup_t \|\nabla f_t(\theta_t) + \beta_t(\theta_t)\|_* = \sup_t \|\mathbb{E}[\zeta_t|\theta_t]\|_* \leq \mathbb{E}[\sup_t \|\zeta_t\|_* |\theta_t] \leq B,$$

from Jensen's inequality. Furthermore, since $\mathbb{E}[(\zeta_t - \nabla f_t(\theta_t) - \beta_t(\theta_t))^\top (\theta - \theta_t)|\theta_t] = 0$, we have that $(\zeta_t - \nabla f_t(\theta_t) - \beta_t(\theta_t))^\top (\theta - \theta_t)$ is a bounded martingale difference. Therefore, by Hoeffding-Azuma inequality, we get

$$\mathbb{P}\left\{ \sum_{t=1}^{T} (\zeta_t - \nabla f_t(\theta_t))^\top (\theta - \theta_t) \geq \varepsilon \right\} \leq \exp\left( \frac{-\varepsilon^2}{64TB^2M} \right).$$

The proposition follows by setting $\varepsilon = 8B\sqrt{TM \log 1/\delta}$. □

To see why the extension is useful, suppose that, instead of obtaining convex losses $f_t$, we receive its approximation $\widehat{f}_t$, which may not be convex. Then we could rewrite the regret as

16

$$\frac{1}{T} \sum_{t=1}^{T} \widehat{f}_t(\theta_t) - \frac{1}{T} \sum_{t=1}^{T} \widehat{f}_t(\theta)$$

$$= \underbrace{\frac{1}{T} \sum_{t=1}^{T} (\widehat{f}_t(\theta_t) - f_t(\theta_t))}_{(I)} + \underbrace{\frac{1}{T} \sum_{t=1}^{T} (f_t(\theta_t) - f_t(\theta))}_{(II)} + \underbrace{\frac{1}{T} \sum_{t=1}^{T} (f_t(\theta) - \widehat{f}_t(\theta))}_{(III)}, \quad (19)$$

for any $\theta$. The first $(I)$ and third term $(III)$ can be bounded by the approximation error between $f_t$ and $\widehat{f}_t$. Since $f_t$ is convex, we could use the technique of online learning to bound the second term $(II)$. However, this approach induces a bias $\beta_t(\theta) = \nabla \widehat{f}_t - \nabla f_t$ because $\theta_t$ is updated by $\nabla \widehat{f}_t$ instead of $\nabla f_t$. The bias $\beta_t(\theta)$ can be controlled by the approximation error between $\nabla \widehat{f}_t$ and $\nabla f_t$. The no-regret analysis is completed by the fact that (19) holds for any $\theta \in \Theta$, including the optimum parameter in hindsight. We will apply this strategy to our setting in which $\widehat{f}_t$, $f_t$ are objective functions induced by a neural network and its linearization, respectively, in Section 4.2. The theory of neural tangent kernel provides a tool to control the approximation of $\widehat{f}_t - f_t$ and $\nabla \widehat{f}_t - \nabla f_t$ in Section 4.1.

We end this section by noting that a constrained optimization problem can be cast as a game in which each player uses an online learning strategy to optimize their utility. With the no-regret guarantee, we can show that such a strategy can approach some equilibrium, resulting in an approximate solution of the constrained optimization problem. We will explain in detail the connection between the regret bound (17) and the constrained optimization problem (2) in the next section.

## 2.3  Online learning for finding equilibrium

It is well-known that there is a strong connection between optimization and game theory, where optimization methods help solve game theory problems (Sohrabi and Azgomi, 2020).

Specifically, iterates induced by an optimization method converges to some equilibrium in a game. In our problem, we require convergence to a coarse-correlated equilibrium (9), (8), (10) to achieve a near optimal and feasible solution of the constrained optimization problem (2). This section shows how the online learning algorithm can be used to find an approximate coarse-correlated equilibrium by considering a simple linear model. For simplicity, we remove the stochastic part of gradients and only discuss projected gradient descent stated as

$$
\begin{aligned}
\theta_{t+1} &= \Pi_\Theta(\theta_t - \eta \nabla_\theta \widetilde{\mathcal{L}}_2(\theta_t, \lambda_t)), \\
\xi_{t+1} &= \Pi_\Xi(\xi_t - \eta \nabla_\xi \mathcal{L}_1(\xi_t, \lambda_t)), \\
\lambda_{t+1} &= \Pi_\Lambda(\lambda_t + \eta \nabla_\lambda \mathcal{L}(\theta_t, \xi_t, \lambda_t)).
\end{aligned}
\tag{20}
$$

Assume the prediction model $y(\theta; x)$ is linear. The following proposition states $\theta, \xi, \lambda$ players, who have utilities $-\widetilde{\mathcal{L}}_2$, $-\mathcal{L}_1$, and $\mathcal{L}$, respectively, and follow strategy of projected gradient descent, can meet approximated coarse correlated equilibrium with convergence rate $O(1/\sqrt{T})$.

**Proposition 6.** *Suppose that conditions in Proposition 1 hold and the prediction model is linear such that $y(\theta; x) = \theta^\top x$ where $\theta \in \Theta$. Then $\{\theta_t, \xi_t, \theta_t\}_{t=1}^T$ computed by (20) with $\eta = 1/\sqrt{T}$ comprise an approximate coarse-correlated equilibrium (9), (8), (10) such that*

$$
\epsilon_\theta = O\left(\frac{1}{\sqrt{T}}\right), \quad \epsilon_\xi = O\left(\frac{1}{\sqrt{T}}\right), \quad \epsilon_\lambda = O\left(\frac{1}{\sqrt{T}}\right).
\tag{21}
$$

*Furthermore, the stochastic model exhibits near-optimal (11) and near-feasible (12) solution of (2).*

*Proof.* Let $f_t(\theta) = \widetilde{\mathcal{L}}_2(\theta, \lambda_t)$. Since $y(\theta; x) = \theta^\top x$ is linear and $\ell$ is convex, $g(\theta) = \mathbb{E}_{(x,z) \sim \mathcal{D}} \ell(y(\theta; x), z)$ is convex. Therefore $f_t(\theta)$ is convex for all $t$. Since $\Theta$ is compact and

18

projected gradient descent (20) is a special case of the stochastic mirror descent (16), Proposition 5 implies that $\{\theta_t\}_t$ satisfies the coarse correlated equilibrium in (9) with $\epsilon_\theta = 1/\sqrt{T}$ by setting $\eta = 1/\sqrt{T}$. Thus, now it remains to verify conditions in Proposition 5.

Since $\Theta$ is compact, we have $B = O(1)$, $\sup_t \|\lambda\|_\infty \leq \kappa$ and

$$\sup_{\theta \in \Theta} \max\{\|\nabla g(\theta)\|_2, \|\widetilde{R}_1(\theta)\|_2, \ldots, \|\widetilde{R}_K(\theta)\|_2\} = O(1),$$

yielding that

$$\sup_t \|\nabla f_t(\theta)\|_2 \leq \sup_{\theta \in \Theta} \|\nabla g(\theta)\|_2 + \|\lambda\|_\infty \sum_{k=1}^{K} \|\widetilde{R}_k(\theta)\|_2 = O(1).$$

Since the conditions in Proposition 5 are satisfied, (9) follows. The same argument can be applied to show that $\{\xi_t\}_t$ and $\{\lambda_t\}_t$ also satisfy the coarse correlated equilibrium. Refer to the proof of Theorem 11 for completed justifications for $\xi$ and $\lambda$. □

Note that projected gradient descent is one possible algorithm to obtain the regret bound. Other methods include the best response strategy and the Follow the Leader Algorithm (Huang et al., 2017). For example, Narasimhan et al. (2019) used best response strategy for some parameters and assumed that there exists an oracle for solving $\min f_t(\theta)$ for each time $t$.

This work focuses on the setting where the prediction model $y(\theta; x)$ is a neural network. As a result, $g(\theta)$ is no longer convex and the above argument fails. For the same reason we cannot directly use the framework of Narasimhan et al. (2019) in our setting. In the following section, we modify the previous argument to train a fair neural network classifier. The key insight is that even though $g$ is not convex, $\ell$ in (2) is convex and the neural network predictor $y(\theta; x)$ can be well approximated by a linear model in the overparameterization regime.

19

# 3    Classification with Neural Networks

This section formalizes our problem setup in Section 3.1 and gives the proposed optimization algorithm in Section 3.2. We introduce a shrinking procedure for randomized models in Section 3.3. Specifically, the shrinking procedure can be used for selecting a smaller subset of candidate parameters generating random predictions. This approach is particularly beneficial for a neural network since its size of parameters is large compared to the dimension of data $d$.

## 3.1    Problem Setup

For the rest of this paper, we use $x, y, z$ to indicate the random variable from data distribution $\mathcal{D}$, an instance of the random variable, the classification model, or the variables of functions, when the meaning is clear from the context.

Consider a prediction model $y(\cdot, \cdot) : \Theta \times \mathcal{X} \to \mathbb{R}$ parameterized by a two-layer neural network

$$y(\theta; x) = \frac{1}{\sqrt{m}} \sum_{i=1}^{m} b_i \sigma(a_i^\top x), \tag{22}$$

where $\mathcal{X} \subset \mathbb{R}^d$, $m$ is the width of the neural network, $b_i \in \{-1, 1\}$ are the output weights, $\sigma(u) = \max(0, u)$ is the rectified linear unit (ReLU) activation, and $\theta = (a_1^\top, \ldots, a_m^\top)^\top \in \Theta \subset \mathbb{R}^{dm}$ are input weights.

The neural network is initialized by the following scheme:

$$b_i \sim \text{Uniform}(\{-1, 1\}), \quad a_i \sim \mathcal{N}(0, I_d/d). \tag{23}$$

where $I_d$ is identity matrix with size $d$. During the training, $\theta$ is restricted to an $L_2$-ball

20

centered at the initialization $\theta(0) = (a_1(0), \ldots, a_m(0))$, that is, $\Theta = \{\theta : \|\theta - \theta(0)\|_2 \leq D\}$, and $b_i$ is fixed for simplicity and omitted from $\theta$. Without loss of generality, we assume $D \geq 1$. Such a setup is common in the literature of theory of deep learning (Oymak and Soltanolkotabi, 2019; Arora et al., 2019; Allen-Zhu et al., 2019a,c; Allen-Zhu and Li, 2019). Note that it is possible to extend the theory for the case where $b$ is optimized.

As $m \to \infty$, the class of functions defined in (22) approximates a subset of the reproducing kernel Hilbert space induced by the kernel

$$K(x, y) = \mathbb{E}_{a \sim \mathcal{N}(0, I/d)} [\mathbb{I}\{a^\top x > 0\} \, \mathbb{I}\{a^\top y > 0\} x^\top y].$$

This function class is sufficiently rich, if the width $m$ and the radius $D$ are sufficiently large (Arora et al., 2019).

With some abuse of notation, given a population distribution $\mathcal{D}$ and sub-populations $\mathcal{D}_1, \ldots, \mathcal{D}_K$ over $\mathcal{X} \times \mathcal{Z}$, we let

$$g(\theta) = \mathbb{E}_{x,z \sim \mathcal{D}} \ell(y(\theta; x), z),$$

where $\ell(\cdot, \cdot) : \mathbb{R} \times \mathcal{Z} \to \mathbb{R}$ is a convex loss function with respect to the first element and

$$R_k(\theta) = \mathbb{E}_{(x,z) \sim \mathcal{D}_k} C_k(y(\theta; x), z),$$
$$\widetilde{R}_k(\theta) = \mathbb{E}_{(x,z) \sim \mathcal{D}_k} \widetilde{C}_k(y(\theta; x), z)).$$

Here $\widetilde{C}_k(\cdot, \cdot) : \mathbb{R} \times \mathcal{Z}$ is a differentiable convex surrogate of $C_k$. This setting slightly generalizes the one described in Section 2.1, where we define $C_k$ as the zero-one loss and $\widetilde{C}_k$ as the hinge loss.

Note that one could use put the constrains as a penalty to the loss such as

$$\min_{\theta \in \Theta, \xi \in \Xi} \quad g(\theta) + \sum_{j=1}^{J} \lambda_j \phi_j(R_1(\theta), \ldots, R_K(\theta)). \tag{24}$$

21

Then $R_k$ can be replaced by its surrogate, and stochastic gradient descent can be applied to avoid accessing whole datasets. Such an approach may seem to solve the main challenges of fairness constraints, however, this formulation requires tuning penalty parameters $\lambda$ carefully. Large penalty parameters may hurt the training loss and result in a bad classification rate. Moreover, including fairness as a hard constraint, rather than a penalty, implies a theoretical guarantee in (12). This is an important property for practical application. For example, in the scenario in which the "80/20 Rule" must be satisfied for the prediction model $y(\theta; x)$ by law (Hardt et al., 2016), where the condition "80/20 Rule" can be stated as $R_1(\theta) < 0.8R_2(\theta)$ for some $R_1, R_2$. Another example is that race and gender discrimination should be prohibited in the classification task. The formula (24) using penalty only "encourage" models to meet algorithmic fairness, which is not enough for the practical requirement above. Again, increasing penalty can alleviate the violation of algorithmic fairness, while it reduces the prediction accuracy. On the other hand, implementing a fair classifier by imposing a hard constraint on the fairness, enjoys such a guarantee in (12) where we know the convergence rate and how the prediction model violates the fairness constraints.

## 3.2   Optimization procedure

This section describes the optimization procedure for solving the constrained problem (4). Recall that we can cast (4) as a game, in which the $\theta$-player, $\xi$-player, and $\lambda$-player have utility functions $-\widetilde{\mathcal{L}}_2$, $-\mathcal{L}_1$, and $\mathcal{L}$, respectively. Each player employs stochastic algorithms to maximize their utility and avoid accessing the whole data set. Since this is not a zero-sum game, i.e., $\mathcal{L} - \mathcal{L}_1 - \widetilde{\mathcal{L}}_2 \neq 0$. As a result, the pure Nash equilibrium may not exist.

---
**Algorithm 2:** Stochastic projected gradient descent
---
**Input:** $\theta_0, \xi_0, \lambda_0, \eta_\theta, \eta_\xi, \eta_\lambda$

\# Stochastic projected gradient descent ;

**for** $t = 1$ **to** $T$ **do**

    Sample $k$ from Uniform($\{1, \ldots, K\}$);

    Sample $j$ from Uniform($\{1, \ldots, J\}$);

    Sample $x'_{t+1}, z'_{t+1} \sim \mathcal{D}_k$ and $x_{t+1}, z_{t+1} \sim \mathcal{D}$;

    Compute $\widehat{\nabla}_\theta \widetilde{\mathcal{L}}_2(\theta_t, \lambda_t), \widehat{\nabla}_\xi \mathcal{L}_1(\xi_t, \lambda_t), \widehat{\nabla}_\lambda \mathcal{L}(\theta_t, \xi_t, \lambda_t)$ by (26);

    $\theta_{t+1} = \Pi_\Theta(\theta_t - \eta_\theta \widehat{\nabla}_\theta \widetilde{\mathcal{L}}_2(\theta_t, \lambda_t))$ ;

    $\xi_{t+1} = \Pi_\Xi(\xi_t - \eta_\xi \widehat{\nabla}_\xi \mathcal{L}_1(\xi_t, \lambda_t))$ ;

    $\lambda_{t+1} = \Pi_\Lambda(\lambda_t + \eta_\lambda \widehat{\nabla}_\lambda \mathcal{L}(\theta_t, \xi_t, \lambda_t))$ ;

**end**
---

In particular, projected stochastic gradient descent for each player is defined by

$$
\begin{aligned}
\theta_{t+1} &= \Pi_\Theta(\theta_t - \eta_\theta \widehat{\nabla}_\theta \widetilde{\mathcal{L}}_2(\theta_t, \lambda_t)), \\
\xi_{t+1} &= \Pi_\Xi(\xi_t - \eta_\xi \widehat{\nabla}_\xi \mathcal{L}_1(\xi_t, \lambda_t)), \\
\lambda_{t+1} &= \Pi_\Lambda(\lambda_t + \eta_\lambda \widehat{\nabla}_\lambda \mathcal{L}(\theta_t, \xi_t, \lambda_t)),
\end{aligned}
\tag{25}
$$

where $\widehat{\nabla}$ denotes unbiased estimators for the gradients and the Lagrangian $\mathcal{L}(\theta, \xi, \lambda)$ is

$$
\mathcal{L}(\theta, \xi, \lambda) = \underbrace{g(\theta) + \sum_{k=1}^{K} \lambda_{J+k} R_k(\theta)}_{\mathcal{L}_2(\theta, \lambda)} + \underbrace{\sum_{j=1}^{J} \lambda_j \phi_j(\xi) - \sum_{k=1}^{K} \lambda_{J+k} \xi_k}_{\mathcal{L}_1(\xi, \lambda)}.
$$

and the surrogate $\widetilde{\mathcal{L}}_2(\theta, \lambda)$ for $\mathcal{L}_2(\theta, \lambda)$

$$
\widetilde{\mathcal{L}}_2(\theta, \lambda) = g(\theta) + \sum_{k=1}^{K} \lambda_{J+k} \widetilde{R}_k(\theta).
$$

23

If $J$ or $K$ is large, it may be intractable to sample all of sub-population $D_1, \ldots, D_k$. We use following sampling scheme to get the unbiased estimator of gradients. Given $k, j$ such that $k, j$ are uniformly chosen from $\{1, \ldots, K\}$ and $\{1, \ldots, J\}$, respectively, we sample $x'_{t+1}, z'_{t+1} \sim \mathcal{D}_k$ and $x_{t+1}, z_{t+1}$ to evaluate $g(\theta)$ and $R_k$, i.e.,

$$
\begin{aligned}
\widehat{\nabla}_\theta \widetilde{\mathcal{L}}_2(\theta, \lambda) =& \nabla_\theta \ell(y(\theta; x_{t+1}), z_{t+1}) + K\lambda_{J+k} \nabla_\theta \widetilde{C}_k(y(\theta; x'_{t+1}), z'_{t+1}), \\
\widehat{\nabla}_\xi \mathcal{L}_1(\xi, \lambda) =& J\lambda_j \nabla_{\xi_k} \phi_j(\xi) e^k - K\lambda_{J+k} e^k, \\
\widehat{\nabla}_\lambda \mathcal{L}(\theta, \xi, \lambda) =& K(C_k(y(\theta; x'_{t+1}), z'_{t+1}) - \xi_k) e^{J+k} + J\phi_j(\xi) e^j,
\end{aligned}
\tag{26}
$$

for each time $t$ where $e^i = (e_1, \ldots, e_{J+K})$ is a unit vector such that $e_i = 1$ and $e_{i'} = 0$ for all $i \neq i'$. It is not hard to see that $\nabla_\theta \widetilde{\mathcal{L}}_2(\theta, \lambda)$, $\nabla_\xi \mathcal{L}_1(\xi, \lambda)$, $\nabla_\lambda \mathcal{L}(\theta, \xi, \lambda)$ are unbiased, i.e.,

$$
\begin{aligned}
\nabla_\theta \widetilde{\mathcal{L}}_2(\theta, \lambda) =& \mathbb{E}_{k, (x_{t+1}, z_{t+1}), (x'_{t+1}, z'_{t+1})} \left[ \widehat{\nabla}_\theta \widetilde{\mathcal{L}}_2(\theta, \lambda) \right], \\
\nabla_\xi \mathcal{L}_1(\xi, \lambda) =& \mathbb{E}_{j,k} \left[ \widehat{\nabla}_\xi \mathcal{L}_1(\xi, \lambda) \right], \\
\nabla_\lambda \mathcal{L}(\theta, \xi, \lambda) =& \mathbb{E}_{j,k,(x'_{t+1}, z'_{t+1})} \left[ \widehat{\nabla}_\lambda \mathcal{L}(\theta, \xi, \lambda) \right].
\end{aligned}
$$

Algorithm 2 summarizes the optimization procedure.

We have seen the connection between online learning and the framework of the constrained optimization problem in Proposition 6 for linear models. Following the same strategy, we would like to establish a no-regret analysis for neural networks. Then, combining the regret bound and Proposition 1 results in a near optimal and near feasible solution. The main technical obstacle is that the model $y(\theta; x)$ parameterized by a neural network in (22) is non-convex and the argument used in Proposition 6 is not valid. The critical observation is that $\ell$ is convex and neural networks are close to linear models in the overparameterization regime. Section 4 introduces local linearization and establishes the no-regret analysis for neural networks. Subsequently, we will put all the ingredients together and show that our algorithm trains a fair neural network classifier.

24

## 3.3 Shrinking

From Proposition 1 we observe that a stochastic model is required to achieve a near-optimal and near-feasible solution. In our setting, the prediction model is parameterized by a neural networks, which may have a considerable amount of parameters and need a large number of iterations for training. As a result, storing all parameters for all iterations for randomized prediction may be intractable. To overcome this difficulty, it can be shown that a smaller mixed equilibrium exists and can be found by using a shrinkage procedure. In particular, let $\theta_1, \ldots, \theta_T \in \Theta$ be a sequence of $T$ iterates obtained by Algorithm 2. Define $c = (g(\theta_1), \ldots, g(\theta_T))^\top$ and $a_j = (h_j(\theta_1), \ldots, h_j(\theta_T))^\top$ for $j = 1, \ldots, J$. The shrinking procedure aims to solve the following linear programming problem:

$$\min_{p \in \Delta^T} c^\top p$$
$$\text{s.t. } a_j^\top p \leq \epsilon, \quad \text{for } j = 1, \ldots, J, \tag{27}$$

for some $\epsilon > 0$ where $\Delta^T$ is the $T$-dimensional simplex. The optimal $p^* = (p_1^*, \ldots, p_T^*)$ for (27) represents the final stochastic classifier, with components representing the probability of sampling $\{1, \ldots, T\}$. Note that due to the convexity of $\phi_j$, we have for all $j$

$$\phi_j \left( \sum_{t=1}^T p_t^* (R_1(\theta_t), \ldots, R_K(\theta_t)) \right) \leq a_j^\top p^* \leq \epsilon,$$

implying that the stochastic solution induced by $p^*$ is near-feasible and achieves the smallest training error among all stochastic models.

Cotter et al. (2019c) showed that every vertex of the feasible region and the optimal solution $p^*$ have at most $J + 1$ nonzero elements. That is, the shrinkage procedure selects at most $J + 1$ iterations for constructing randomized solution and there are only $J + 1$ non-zero elements in $p^*$, which reduces memory cost significantly. Since $J$ is much smaller

than $T$, heuristically, we only sample a small amount of iterates compared to $T$ in our experiment where $J$ is the number of constraints and $T$ is the total number of iterates.

# 4   Main Results

In this section, we show how to incorporate the game-theoretic framework, described in Section 2.1, for achieving algorithmic fairness of neural network classifiers. Specifically, our main result in Section 4.3 establishes a near-optimal and near-feasible stochastic solution for neural network classifiers with fairness constraints. The key tool used to prove this result is local linearization of overparametrized neural networks in Section 4.1 and their regret analysis through online learning in Section 4.2.

For simplicity, we consider the binary classification task and two-layer neural networks. However, with additional effort, our ideas can be extended to multi-class classification and deep neural networks.

## 4.1   Local Linearization for Neural Networks

This section presents the phenomenon of local linearization for nueral networks. Note that we assume for any $x \in \mathcal{X}$, $\|x\|_2 \leq 1$ for the rest of the paper. Define a local linearization of (22), $y^0(\theta; x)$, at the random initialization $\theta(0)$:

$$y^0(\theta; x) = \frac{1}{\sqrt{m}} \sum_{i=1}^m b_i\, \mathbb{I}\{a_i(0)^\top x > 0\} a_i^\top x := \Phi(x)^\top \theta,$$

where
$$\Phi(x) = \frac{(b_1\, \mathbb{I}\{a_1(0)^\top x > 0\} x^\top, \ldots, b_m\, \mathbb{I}\{a_m(0)^\top x > 0\} x^\top)^\top}{\sqrt{m}}.$$

Noting that $b_i$ is fixed after the initialization and $\Phi(x)$ is a feature map of $x$, we see that $y^0$ is linear w.r.t. to $\theta$. The key observation is that as the width grows, neutral networks exhibits similar behavior to linear models with the random feature $\Phi(x)$.

We need the following regularity condition on the data distribution.

**Assumption 7** (Regularity of data distribution). *For any unit vector $e$ and a constant $\zeta > 0$, there exists $c > 0$, such that $\mathbb{P}_x(|e^\top x| \leq \zeta) \leq c\zeta$.*

The regularity condition on an observation $x$ holds as long as the distribution of data has an upper-bounded density. Under the regularity conditions above, the following characterization shows that the expected approximation error of the local linearization at $\theta(0)$ vanishes to zero as $m \to \infty$.

**Proposition 8.** *Suppose that $\|x\|_2 \leq 1$ almost surely and Assumption 7 holds. Then we have that for $\theta \in \Theta = \{\theta : \|\theta - \theta(0)\|_2 \leq D\}$*

1. $\sup_x |y(\theta; x)| \leq D$ *almost surely;*

2. $\sup_x \|\nabla y(\theta; x)\|_2 \leq 1$ *almost surely;*

3. $\mathbb{E}_{x,\theta(0)}|y(\theta; x) - y^0(\theta; x)|^2 = O(D^3 m^{-1/2})$;

4. $\mathbb{E}_{x,\theta(0)}\|\nabla y(\theta; x) - \nabla y^0(\theta; x)\|_2^2 = O(Dm^{-1/2})$.

*Proof.* We establish the first statement as follows:

$$
\begin{aligned}
|y(\theta; x)|^2 &\leq \frac{1}{m}\left(\sum_{i=1}^{m} \mathbb{I}\{a_i^\top x > 0\}|a_i^\top x|\right)^2 \\
&\leq \frac{1}{m}\left(\sum_{i=1}^{m} \mathbb{I}\{a_i^\top x > 0\}\right)\left(\sum_{i=1}^{m} |a_i^\top x|^2\right)
\end{aligned}
$$

27

$$\leq \frac{1}{m} \left( \sum_{i=1}^{m} \mathbb{I}\{a_i^\top x > 0\} \right) \left( \sum_{i=1}^{m} \|a_i\|_2^2 \right)$$

$$\leq D^2,$$

where we have used that $\|x\|_2 \leq 1$. Similarly, it is not hard to see second statement that $\|\nabla y(\theta; x)\|_2 \leq 1$ since we have

$$\nabla y(\theta; x) = \frac{1}{\sqrt{m}} \left( \mathbb{I}\{a_1^\top x > 0\}x^\top, \ldots, \mathbb{I}\{a_m^\top x > 0\}x^\top \right)^\top$$

Next, we prove the local linearization for neural network models focusing on two-layer neural networks to present the main ideas in a transparent fashion. The proof closely follows ideas in Liu et al. (2019).

Given a layer $i$ s.t. $\mathbb{I}\{a_i^\top x > 0\} \neq \mathbb{I}\{a_i(0)^\top x > 0\}$, we have

$$(a_i^\top x)(a_i(0)^\top x) < 0 \quad \implies \quad |a_i(0)^\top x| \leq |(a_i(0) - a_i)^\top x| \leq \|a_i(0) - a_i\|_2, \quad (28)$$

where the last inequality follows by the fact that $\|x\|_2 \leq 1$. Thus, we have

$$|y(\theta; x) - y^0(\theta; x)| = \frac{1}{\sqrt{m}} \left| \sum_{i=1}^{m} b_i (\mathbb{I}\{a_i^\top x > 0\} - \mathbb{I}\{a_i(0)^\top x > 0\}) a_i^\top x \right|$$

$$\leq \frac{1}{\sqrt{m}} \sum_{i=1}^{m} \left| \mathbb{I}\{a_i^\top x > 0\} - \mathbb{I}\{a_i(0)^\top x > 0\} \right| (|a_i(0)^\top x| + \|a_i - a_i(0)\|_2)$$

$$\leq \frac{1}{\sqrt{m}} \sum_{i=1}^{m} \mathbb{I}\{|a_i(0)^\top x| \leq \|a_i(0) - a_i\|_2\} \times (|a_i(0)^\top x| + \|a_i - a_i(0)\|_2)$$

$$\leq \frac{1}{\sqrt{m}} \sum_{i=1}^{m} \mathbb{I}\{|a_i(0)^\top x| \leq \|a_i(0) - a_i\|_2\} \times (|a_i(0)^\top x| + \|a_i - a_i(0)\|_2)$$

$$\leq \frac{2}{\sqrt{m}} \sum_{i=1}^{m} \mathbb{I}\{|a_i(0)^\top x| \leq \|a_i(0) - a_i\|_2\} \times \|a_i - a_i(0)\|_2,$$

28

where we have used (28) to show the second inequality and $\mathbb{I}\{|c| \le d\}|c| \le \mathbb{I}\{|c| \le d\}d$ is used to show the last inequality. Applying Cauchy–Schwarz inequality twice and using Assumption 7, we get

$$
\begin{aligned}
\mathbb{E}_{x,\theta(0)}&|y(\theta;x) - y^0(\theta;x)|^2 \\
&\le \frac{4}{m}\mathbb{E}_{x,\theta(0)}\left(\sum_{i=1}^{m}\mathbb{I}\{|a_i(0)^\top x| \le \|a_i(0) - a_i\|_2\}\right)\left(\sum_{i=1}^{m}\|a_i - a_i(0)\|_2^2\right) \\
&\le \frac{4D^2}{m}\sum_{i=1}^{m}\mathbb{E}_{x,\theta(0)}\,\mathbb{I}\{|a_i(0)^\top x| \le \|a_i(0) - a_i\|_2\} \\
&\le \frac{4D^2}{m}\sum_{i=1}^{m}\mathbb{E}_{\theta(0)}\mathbb{P}\{|a_i(0)^\top x| \le \|a_i(0) - a_i\|_2\} \\
&\le \frac{4D^2}{m}\mathbb{E}_{\theta(0)}\left[\sum_{i=1}^{m}\frac{\|a_i(0) - a_i\|_2}{\|a_i(0)\|_2}\right] \\
&\le \frac{4D^2}{m}\mathbb{E}_{\theta(0)}\left(\sum_{i=1}^{m}\|a_i(0)\|_2^{-2}\right)^{1/2}\left(\sum_{i=1}^{m}\|a_i(0) - a_i\|_2^2\right)^{1/2} \\
&\le \frac{4D^3}{m}\left(\sum_{i=1}^{m}\mathbb{E}_{\theta(0)}\|a_i(0)\|_2^{-2}\right)^{1/2} \\
&= \frac{4D^3}{m^{1/2}(d-2)^{1/2}},
\end{aligned}
\tag{29}
$$

where we have used $\sum_{i=1}^{M}\|a_i(0) - a_i\|_2 \le D$ for all $i$, Jensen's inequality and the following fact. Note that if $v \sim \chi(d)$, where $\chi(d)$ is chi-squared distribution with degree of freedom $\mu$, then $1/v$ follows the inverse-chi-squared distribution and $\mathbb{E}v = 1/(d-2)$. This fact implies

$$
\mathbb{E}\|a_i(0)\|_2^{-2} = \frac{1}{d-2}.
\tag{30}
$$

Following a similar argument, we have

$$
\begin{aligned}
\mathbb{E}\|\nabla y(\theta; x) - \nabla y^0(\theta; x)\|_2^2 &= \frac{1}{m}\mathbb{E}\sum_{i=1}^{m} \mathbb{1}\{|a_i(0)^\top x| \le \|a_i(0) - a_i\|_2\}\|x\|_2 \\
&\le \frac{1}{m}\mathbb{E}\sum_{i=1}^{m} \frac{\|a_i(0) - a_i\|_2}{\|a_i(0)\|_2} \\
&\le \frac{D}{m}\mathbb{E}\left(\sum_{i=1}^{m} \|a_i(0)\|_2^{-2}\right)^{1/2} \\
&= O\left(\frac{D}{\sqrt{m}}\right).
\end{aligned}
$$

Therefore, the proof is complete. □

We note that the results of Proposition 8 hold uniformly over $\theta \in \Theta$. Also note that the randomness in the statements three and four are over the random initialization and data distribution. An immediate implication is that our main result in Section 4.3 can be easily extended to deep learning setting since similar local linearization also holds for multi-layer neutral networks (Allen-Zhu et al., 2019b; Gao et al., 2019).

## 4.2 Online Learning Problems for Neural Networks

Consider an online learning problem under a classification setting. Assume that we have a prediction model $y(\theta; x)$ for a target $z$ such that $(x, z) \sim \mathcal{D}$, where $\mathcal{D}$ is the data distribution. Instead of a fixed loss function $f$ that measures the prediction error between $y(\theta; x)$ and $z$, an adversarial environment generates a new loss function $f_t$ against our prediction model $y(\theta; x)$ for each $t$. Note that we do not assume the data distribution $\mathcal{D}$ to change over time. The adversarial environment is from the setting where different parameters are updated based on their own objective functions and against each other rather than the shift of data distribution. See Section 2.1 and Proposition 6 as an example.

Formally speaking, at each time $t$, an adversarial environment outputs a new loss functions $f_t : \mathcal{Y} \times \mathcal{Z} \rightarrow \mathbb{R}$ such that $f_t$ is convex with respect to $y$ for all $z$. An online learning algorithm for classifier $y(\theta; x)$ aims to find a sequence of parameters $\theta_t$ based on past information that controls the regret in hindsight defined by

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{x,z,\theta(0)}[f_t(y(\theta_t; x), z)] - \min_{\theta} \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{x,z,\theta(0)}[f_t(y(\theta; x), z)], \tag{31}$$

with high probability. The randomness is from the procedure for choosing $\theta_t$, which can be deterministic or stochastic. In particular, if we implement a projected gradient descent for finding $\theta_t$, then the high probability statement can be discarded as in Proposition 6. However, since stochastic gradient descent is a standard optimization algorithm for training neural networks in practice on large scale datasets, we choose projected stochastic gradient descent for this online learning problem. We will show the regret bound for stochastic projected gradient descent under our setting. The formal definition of stochastic projected gradient descent is as follows:

$$\mu_{t+1} = \theta_t - \eta \nabla f_t(y(\theta_t; x_{t+1}), z_{t+1}), \quad \theta_{t+1} = \arg\min_{\theta \in \Theta} \|\theta - \mu_{t+1}\|_2^2, \tag{32}$$

where $(x_{t+1}, z_{t+1}) \sim \mathcal{D}$ such that $\{x_t, z_t\}_t$ are independent. Even though $f_t$ is convex for all $t$, the prediction model $y(\theta; x)$ parameterized by a neural network is not convex with respect to $\theta$. Hence, the composition $f_t(y(\theta; x), z)$ is not convex with respect to $\theta$, which violates the convexity assumption in conventional online learning in Proposition 5. However, the idea is that the composition of $f_t(y, z)$ and the local linearization $y^0(\theta; x)$ is *convex*. In

31

particular, to control the regret (31), we rewrite it as

$$
\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}_{x,z,\theta(0)}[f_t(y(\theta_t;x),z)] - \frac{1}{T}\sum_{t=1}^{T}\mathbb{E}_{x,z,\theta(0)}[f_t(y(\theta;x),z)]
$$

$$
= \underbrace{\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}_{x,z,\theta(0)}[f_t(y(\theta_t;x),z)] - \frac{1}{T}\sum_{t=1}^{T}\mathbb{E}_{x,z,\theta(0)}[f_t(y^0(\theta_t;x),z)]}_{(I)}
$$

$$
\tag{33}
$$

$$
+ \underbrace{\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}_{x,z,\theta(0)}[f_t(y^0(\theta_t;x),z)] - \frac{1}{T}\sum_{t=1}^{T}\mathbb{E}_{x,z,\theta(0)}[f_t(y^0(\theta;x),z)]}_{(II)}
$$

$$
+ \underbrace{\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}_{x,z,\theta(0)}[f_t(y^0(\theta;x),z)] - \frac{1}{T}\sum_{t=1}^{T}\mathbb{E}_{x,z,\theta(0)}[f_t(y(\theta;x),z)]}_{(III)},
$$

for any $\theta \in \Theta$. The second term $(I)$ and $(III)$ can be bounded by controlling the error of linearization in Proposition 8. In fact, we show that two terms decays at the rate $O(D^{5/2}m^{-1/4})$. The second term $(II)$ in (33) can be controlled by the classical result on convex online learning, since now $y^0(\theta;x)$ is linear w.r.t. $\theta$ and $f_t(y^0(\theta_t;x),z)$ is convex w.r.t. $\theta$. However, the crucial problem is that we use the original stochastic gradient $\nabla f_t(y(\theta_t;x_{t+1}),z_{t+1})$ instead of the gradient of linearization $\nabla f_t(y^0(\theta_t;x_{t+1}),z_{t+1})$, which introduces additional error in updating. To address that, we rewrite (32) as

$$
\mu_{t+1} = \theta_t - \eta\nabla f_t(y^0(\theta_t;x_{t+1}),z_{t+1}) - \eta\Delta_t,
$$

$$
\theta_{t+1} = \arg\min_{\theta\in\Theta}\|\theta - \mu_{t+1}\|_2^2,
$$

$$
\tag{34}
$$

where

$$
\Delta_t = \nabla f_t(y(\theta_t;x_{t+1}),z_{t+1}) - \nabla f_t(y^0(\theta_t;x_{t+1}),z_{t+1}).
$$

It turns out that the noise $\Delta_t$ can also be controlled by the local linearization. In summary, we obtain the following theorem for the regret bound (31).

32

**Theorem 9.** *Suppose all assumptions in Proposition 8 hold. Assume $f_t(y, z)$ is convex with respect to $y$ for all $z$, such that $|\nabla_y f_t(y_1, z) - \nabla_y f_t(y_2, z)| \leq L_f |y_1 - y_2|$, for all $t, z, y_1, y_2$. Then, with probability $1 - \delta$, we have*

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{x,z,\theta(0)}[f_t(y^0(\theta_t; x), z)] - \min_\theta \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{x,z,\theta(0)}[f_t(y^0(\theta; x), z)]$$

$$\leq O\left(\frac{\eta L_f D}{2} + \frac{D^2}{2T\eta} + \frac{L_f D^2 \sqrt{\ln(1/\delta)}}{\sqrt{T}} + \frac{L_f D^{5/2}}{m^{1/4}}\right),$$

*where $\theta_t$ is computed by (34).*

This theorem states that with regular conditions for smoothness of the gradient, we can bound the regret with extra error term controlled by the approximation error of local linearization. The theory of neural tangent kernel provides a powerful tool to overcome the difficulty from non-convexity of neural networks and to establish no-regret analysis for neural networks in overparametrized regime.

We end the section with the proof of Theorem 9.

*Proof of Theorem 9.* From the decomposition (33), it suffices to bound three terms $(I), (III), (III)$. In particular, we will show that

$$(I), (III) = O\left(\frac{L_f D^{5/2}}{m^{1/4}}\right), \quad (II) = O\left(\frac{\eta L_f D}{2} + \frac{D^2}{2T\eta} + \frac{L_f D^2 \sqrt{\ln(1/\delta)}}{\sqrt{T}} + \frac{L_f D^{5/2}}{m^{1/4}}\right).$$

and

$$(II) = O\left(\frac{\eta B_f}{2} + \frac{D^2}{2T\eta} + D\sqrt{\frac{\ln(1/\delta)}{2T}} + \frac{L_f D^{3/2}}{m^{1/4}}\right)$$

The first term $(I)$ in (33) can be bounded in the following way:

$$(I) = \mathbb{E}_{x,z,\theta(0)}\left[\frac{1}{T}\sum_{t=1}^{T}f_t(y(\theta_t;x),z) - \frac{1}{T}\sum_{t=1}^{T}f_t(y^0(\theta_t;x),z)\right]$$

$$\leq \frac{1}{T}\sum_{t=1}^{T}\mathbb{E}_{x,z,\theta(0)}|\nabla_y f_t(y(\theta_t;x),z)(y(\theta_t;x) - y^0(\theta_t;x))|$$

$$\leq \sup_{t,y,z:|y|\leq D}|\nabla_y f_t(y,z)| \times \frac{1}{T}\sum_{t=1}^{T}\mathbb{E}_{x,\theta(0)}|(y(\theta_t;x) - y^0(\theta_t;x))| \quad (35)$$

$$\leq \sup_{t,y,z:|y|\leq D}|\nabla_y f_t(y,z)| \times \frac{1}{T}\sum_{t=1}^{T}\left(\mathbb{E}_{x,\theta(0)}|(y(\theta_t;x) - y^0(\theta_t;x))|^2\right)^{1/2}$$

$$\leq \sup_{t,y,z:|y|\leq D}|\nabla_y f_t(y,z)| \times \sup_{\theta}\left(\mathbb{E}_{x,\theta(0)}|(y(\theta;x) - y^0(\theta;x))|^2\right)^{1/2}$$

$$= O\left(\frac{L_f D^{5/2}}{m^{1/4}}\right) \quad \text{almost surely,}$$

where we have used the convexity of $f_t$, $|y(\theta_t;x),z)| \leq D$ and the fact that

$$\sup_{t,y,z:|y|\leq D}|\nabla_y f_t(y,z)| \leq \sup_{t,y,z:|y|\leq D}|\nabla_y f_t(y,z) - \nabla_y f_t(0,z)| + |\nabla_y f_t(0,z)| = O(L_f D).$$

The last inequality follows by Proposition 8. The third terms $(III)$ is controlled in a similar manner.

To apply Proposition 5 for the second term $(II)$, it suffices to show $\sup_t \|\zeta_t\|_2 < O(L_f D)$, where $\zeta_t = \nabla f_t(y(\theta_t;x_{t+1}),z_{t+1})$,

$$\mathbb{E}\left[\zeta_t \mid \theta_t\right] = \mathbb{E}_{x,z,\theta(0)}[\nabla f_t(y^0(\theta_t;x),z) \mid \theta_t] + \beta_t(\theta_t), \quad (36)$$

and

$$\beta_t(\theta_t) = \mathbb{E}_{x,z,\theta(0)}[\nabla f_t(y(\theta_t;x),z) - \nabla f_t(y^0(\theta_t;x),z) \mid \theta_t]. \quad (37)$$

We have

$$\sup_t \|\zeta_t\|_2 = \sup_t \|\nabla f_t(y(\theta_t;x_{t+1}),z_{t+1})\|_2$$

$$= \sup_t \|\nabla_y f_t(y(\theta_t; x_{t+1}), z_{t+1}) \nabla_\theta y(\theta_t; x_{t+1})\|_2$$

$$\leq \sup_{t,y,z: |y| \leq D} |\nabla_y f_t(y, z)|$$

$$\leq \sup_{t,y,z: |y| \leq D} |\nabla_y f_t(y, z) - \nabla_y f_t(0, z)| + |\nabla_y f_t(0, z)|$$

$$\leq O(L_f D),$$

where we have used $|y(\theta; x)| \leq D$ and $\sup_{x,\theta} \|\nabla_\theta y(\theta; x)\|_2 \leq 1$ in Proposition 8. Applying Proposition 5, we get

$$(II) = \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{x,z,\theta(0)}[f_t(y^0(\theta_t; x), z)] - \frac{1}{T} \sum_{t=1}^{T} \mathbb{E}_{x,z,\theta(0)}[f_t(y^0(\theta; x), z)]$$

$$\leq \frac{\eta B_f}{2} + \frac{D^2}{2T\eta} + 8 B_f D \sqrt{\frac{\ln(1/\delta)}{2T}} + \frac{2D}{T} \sum_{t=1}^{T} \|\beta_t(\theta_t)\|_2$$

with probability $1 - \delta$. To bound $\frac{1}{T} \sum_{t=1}^{T} \|\beta_t(\theta_t)\|_2$, we have the decomposition:

$$\|\beta_t(\theta_t)\|_2 = \|\nabla f_t(y(\theta_t; x), z) - \nabla f_t(y^0(\theta_t; x), z)\|_2$$

$$= \mathbb{E}_{x,z,\theta(0)} \|\nabla_y f_t(y(\theta_t; x), z) \nabla_\theta y(\theta_t; x) - \nabla_y f_t(y^0(\theta_t; x), z) \nabla_\theta y^0(\theta_t; x)\|_2$$

$$\leq \mathbb{E}_{x,z,\theta(0)} |\nabla_y f_t(y(\theta_t; x), z) - \nabla_y f_t(y^0(\theta_t; x), z)| \|\nabla_\theta y(\theta_t; x)\|_2$$

$$+ \mathbb{E}_{x,z,\theta(0)} |\nabla f_t(y^0(\theta_t; x), z)| \|\nabla_\theta y(\theta_t; x) - \nabla_\theta y^0(\theta_t; x)\|_2$$

$$\leq O(L_f D) \mathbb{E}_{x,z,\theta(0)} \|\nabla_\theta y(\theta_t; x) - \nabla_\theta y^0(\theta_t; x)\|_2$$

$$= O\left(\frac{L_f D^{3/2}}{m^{1/4}}\right)$$

where we have used Jensen's inequality and the final statement in Proposition 8. Combining

all the results, we get

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}_{x,z,\theta(0)}[f_t(y^0(\theta_t;x),z)] - \min_{\theta}\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}_{x,z,\theta(0)}[f_t(y^0(\theta;x),z)]$$

$$\leq O\left(\frac{\eta L_f D}{2} + \frac{D^2}{2T\eta} + L_f D^2\sqrt{\frac{\ln(1/\delta)}{T}} + \frac{L_f D^{3/2}}{m^{1/4}} + \frac{L_f D^{5/2}}{m^{1/4}}\right)$$

$$\leq O\left(\frac{\eta L_f D}{2} + \frac{D^2}{2T\eta} + \frac{L_f D^2\sqrt{\ln(1/\delta)}}{\sqrt{T}} + \frac{L_f D^{5/2}}{m^{1/4}}\right),$$

which completes the proof. $\qquad\square$

## 4.3   Global Convergence

We are now ready to present our main theorem on global convergence. Define

$$\Theta = \{\theta : \|\theta - \theta(0)\|_2 \leq D\}, \quad \widetilde{\Theta} = \left\{\theta \in \Theta : (\widetilde{R}_1(\theta),\ldots,\widetilde{R}_K(\theta)) \in \bigcap_{j=1}^{J}\mathrm{dom}\,\phi_j,\right\},$$

$$\Lambda = \{\lambda \geq 0 : \|\lambda\|_\infty \leq \kappa\}, \quad \Xi = \left\{(\xi_1,\ldots,\xi_K) \in \bigcap_{j=1}^{J}\mathrm{dom}\,\phi_j : |\xi_k| \leq \max_{\theta\in\Theta}|\widetilde{R}_k(\theta)|\right\}. \tag{38}$$

We require the following regularity condition that allows us to use the general framework of Proposition 1 to establish a fair neural network classifier.

**Assumption 10** (Regularity of Objectives and constraints)**.** *Assume that for any* $z$, $k = 1,\ldots,K$, $j = 1,\ldots,J$

1. $\sup_{|y|\leq D}|C_k(y,z)| \leq L$;

2. $\widetilde{C}_k(y,z)$ *is differentiable and convex with respect to* $y$, $C(y,z) \leq \widetilde{C}(y,z)$ *and* $|\nabla_y\widetilde{C}_k(y_1,z) - \nabla_y\widetilde{C}_k(y_2,z)| \leq L|y_1 - y_2|$ *for all* $z$ *and* $|y_1| \leq D$, $|y_2| \leq D$;

36

3. $\ell(y, z)$ is differentiable and convex with respect to $y$ such that $|\nabla_y \ell(y_1, z) - \nabla_y \ell(y_2, z)| \leq L|y_1 - y_2|$ for all $z$ and $|y_1| \leq D$, $|y_2| \leq D$;

4. $\phi_j$ is strictly jointly convex, monotonically increasing in each argument, and L-Lipschitz w.r.t. the infinite norm such that $|\nabla_{\xi_k} \phi_j(\xi^1) - \nabla_{\xi_k} \phi_j(\xi^2)| \leq L\|\xi^1 - \xi^2\|_\infty$ for all $\xi^1, \xi^2 \in \Xi$.

Without loss of generality, we assume $\kappa \geq 1$ and $L \geq 1$. Before we state the main result in this section, we have a remark on the boundedness assumption in (38) on the spaces $\Theta$, $\Xi$ and $\Lambda$. The boundedness assumption is sufficient for applying no-regret analysis Proposition 5, which is the crucial ingredient for proving convergence to an approximate coarse-correlated equilibrium as we showed in Proposition 6. The boundedness assumption for $\Theta$ is also sufficient for local localization in Proposition 8. As a result, despite the fact that the behavior of local linearization around the initialization is not specific to overparameterized neural networks and may not fully explain their successes (Chizat et al., 2019), the boundedness assumption is natural in our setting and the projection operator is necessary. In other word, although the weights $\theta$ are indeed a measure of generalization (Neyshabur et al., 2018a) and regularization on weights can improve the generalization (Krogh and Hertz, 1991; Ba et al., 2016; Salimans and Kingma, 2016), restricting $\theta$ in $\Theta = \{\theta : \|\theta - \theta(0)\|_2 \leq D\}$ and using projected gradient descent may not be practical for training neural networks. However, since projections into some compact domains are necessary for both online learning and local linearization, projected stochastic gradient descent is a natural choose for optimization in our setting.

Our main result shows that if the classifier is parameterized by a neural network (22), and we run projected stochastic gradient descent in Algorithm 2, then we obtain a near-optimal and near feasible solution of the constrained optimization problem (4). In partic-

ular, we can build a fair neural networks classifier with provable guarantees.

**Theorem 11.** *Suppose Assumptions 7 and 10 hold. Set*

$$\eta_\theta = \sqrt{\frac{D}{\kappa L T}}, \quad \eta_\xi = \sqrt{\frac{1}{\kappa L T}}, \quad \eta_\lambda = \sqrt{\frac{\kappa}{L^2 D T}}.$$

*Then, with probability $1 - 3\delta$, the iterates $\{\theta_t, \xi_t, \lambda_t\}_{t=1}^T$ of Algorithm 2 comprise an approximate coarse-correlated equilibrium (8), (9), (10) with*

$$
\begin{aligned}
\epsilon_\theta =& O\left(\frac{\kappa L D^{3/2}\sqrt{\ln(1/\delta)}}{\sqrt{T}} + \frac{\kappa D^{5/2}}{m^{1/4}}\right), \\
\epsilon_\xi =& O\left(\frac{\kappa D L^2 \sqrt{\ln(1/\delta)}}{\sqrt{T}}\right), \\
\epsilon_\lambda =& O\left(\frac{D L^2 \sqrt{\kappa \ln(1/\delta)}}{\sqrt{T}}\right).
\end{aligned}
\tag{39}
$$

*Moreover, with probability $1 - 3\delta$, we have*

$$\frac{1}{T}\sum_{t=1}^T g(\theta_t) - \min_{\theta \in \widetilde{\Theta} \cap F} g(\theta) \le O\left(\frac{\kappa L D^{3/2}\sqrt{\ln(1/\delta)}}{\sqrt{T}} + \frac{\kappa D^{5/2}}{m^{1/4}}\right), \tag{40}$$

*and for $j = 1, \ldots, J$,*

$$\phi_j\left(\frac{1}{T}\sum_{t=1}^T (R_1(\theta_t), \ldots, R_K(\theta_t))\right) \le O\left(\frac{L D^{3/2}\sqrt{\ln(1/\delta)}}{\sqrt{T}} + \frac{L D^2}{\kappa} + \frac{D^{5/2}}{m^{1/4}}\right), \tag{41}$$

*where $F = \{\theta : \max_j \phi_j(\widetilde{R}_1(\theta), \ldots, \widetilde{R}_K(\theta)) \le 0\}$.*

*Proof.* Since (40) and (41) immediately follow using Proposition 1, it suffices to prove the regret bounds in (39). For $\theta$, we will use the result in Theorem 9 and justify its conditions. For $\xi$ and $\lambda$, we could apply Proposition 5, since their objective functions are convex.

We first prove (9). It suffices to verify the conditions of Theorem 9. Letting

$$f_t(y, z) = \ell(y, z) + \sum_{k=1}^{K} \lambda_{t, J+k} \widetilde{C}(y, z), \tag{42}$$

by Assumption 10, we know

$$|\nabla_y f_t(y_1, x) - \nabla_y f_t(y_2, x)|$$

$$\leq |\nabla_y \ell(y_1, z) - \nabla_y \ell(y_2, z)| + \sum_{k=1}^{K} \lambda_{j+k, t} \left| \nabla_y \widetilde{C}_k(y_1, z) - \nabla_y \widetilde{C}_k(y_2, z) \right|$$

$$= O(\kappa L).$$

Therefore, using Theorem 9, with probability $1 - \delta$, we have

$$\frac{1}{T} \sum_{t=1}^{T} \widetilde{\mathcal{L}}_2(\theta_t, \lambda_t) - \min_{\theta \in \Theta} \frac{1}{T} \sum_{t=1}^{T} \widetilde{\mathcal{L}}_2(\theta, \lambda_t)$$

$$= O\left( \frac{\eta_\theta \kappa L D}{2} + \frac{D^2}{2T\eta_\theta} + \kappa L D \sqrt{\frac{\ln(1/\delta)}{T}} + \frac{\kappa L D^{5/2}}{m^{1/4}} \right),$$

which proves (9) with

$$\epsilon_\theta = O\left( \frac{D^{3/2}}{\sqrt{\kappa L T}} + \frac{\kappa L D \sqrt{\ln(1/\delta)}}{\sqrt{T}} + \frac{\kappa L D^{5/2}}{m^{1/4}} \right) \leq O\left( \frac{\kappa L D^{3/2} \sqrt{\ln(1/\delta)}}{\sqrt{T}} + \frac{\kappa D^{5/2}}{m^{1/4}} \right) \tag{43}$$

by setting $\eta_\theta = \sqrt{D/\kappa L T}$.

For $\xi$, we have

$$\sup_{\theta \in \Theta} |\widetilde{R}_k(\theta)| = |\mathbb{E}[C_k(y(\theta, x), z)]|$$

$$\leq \sup_{\theta \in \Theta} \mathbb{E}[|C_k(y(\theta, x), z) - C_k(0, z)| + |C_k(0, z)|]$$

$$\leq \sup_{|y| \leq D} \mathbb{E}[|C_k(y, z) - C_k(0, z)| + |C_k(0, z)|] \tag{44}$$

$$\leq \sup_{|y| \leq D} \mathbb{E}[|\nabla_y C_k(y, z) - \nabla_y C_k(0, z)||y| + |C_k(0, z)|]$$

$$\leq O(DL),$$

39

where we have used convexity and $\sup_{\theta \in \Theta} |y(\theta, x)| \le D$. Therefore

$$\sup_{\xi \in \Xi} |\nabla_{\xi_k} \phi_j(\xi)| \le O(DL^2),$$

yielding

$$\sup_t \|\widehat{\nabla}_\xi \mathcal{L}_1(\xi_t, \lambda_t)\|_2 \le \sup_{t, \xi \in \Xi} \|J\lambda_j \nabla_{\xi_k} \phi_j(\xi)e^k - K\lambda_{J+k}e^k\|_2 \le O(\kappa DL^2).$$

Proposition 5 implies

$$\frac{1}{T} \sum_{t=1}^{T} \mathcal{L}_1(\xi_t, \lambda_t) - \min_{\xi \in \Xi} \frac{1}{T} \sum_{t=1}^{T} \mathcal{L}_1(\xi, \lambda_t) = O\left(\frac{\eta \kappa DL^2}{2} + \frac{DL}{2T\eta} + \kappa DL^2 \sqrt{\frac{\ln(1/\delta)}{T}}\right), \quad (45)$$

with probability $1 - \delta$, so (8) follows with

$$\epsilon_\xi = O\left(\frac{\kappa DL^2 \sqrt{\ln(1/\delta)}}{\sqrt{T}}\right) \quad (46)$$

by setting $\eta_\xi = 1/\sqrt{\kappa LT}$.

For $\lambda$, using

$$\sup_t \|\widehat{\nabla}_\lambda \mathcal{L}(\theta_t, \xi_t, \lambda_t)\|_2 = \sup_{\theta \in \Theta, \xi \in \Xi} \|K(C_k(y(\theta; x'_{t+1}), z'_{t+1}) - \xi_k)e^{J+k} + J\phi_j(\xi)e^j\|_2 \le O(DL^2)$$

and $\|\lambda_t\|_\infty \le \kappa$, and Proposition 5, we get

$$\frac{1}{T} \sum_{t=1}^{T} \mathcal{L}(\theta_t, \xi_t, \lambda_t) - \min_{\lambda \in \Lambda} \frac{1}{T} \sum_{t=1}^{T} \mathcal{L}(\theta_t, \xi_t, \lambda) = O\left(\frac{DL^2 \eta_\lambda}{2} + \frac{\kappa}{2T\eta_\lambda} + DL^2 \sqrt{\frac{\kappa \ln(1/\delta)}{T}}\right),$$

with probability $1 - \delta$, which implies (10) with

$$\epsilon_\lambda = O\left(\frac{DL^2 \sqrt{\kappa \ln(1/\delta)}}{\sqrt{T}}\right),$$

by setting $\eta_\lambda = \sqrt{\kappa/(L^2 DT)}$.

40

The theorem follows by the union bound, Proposition 1, and $\sup_{\theta \in \Theta, x, z} |\ell(y(\theta, x), z)| \leq O(LD^2)$ with

$$\epsilon_\theta + \epsilon_\xi + \epsilon_\lambda = O\left(\frac{\kappa LD^{3/2}\sqrt{\ln(1/\delta)}}{\sqrt{T}} + \frac{\kappa D^{5/2}}{m^{1/4}} + \frac{\kappa DL^2\sqrt{\ln(1/\delta)}}{\sqrt{T}} + \frac{DL^2\sqrt{\kappa \ln(1/\delta)}}{\sqrt{T}}\right)$$

$$\leq O\left(\frac{\kappa LD^{3/2}\sqrt{\ln(1/\delta)}}{\sqrt{T}} + \frac{\kappa D^{5/2}}{m^{1/4}}\right).$$

$\square$

We have several remarks for Theorem 11. First of all, as the width $m$ increases, the approximation error $O(\kappa D^{5/2}/m^{1/4})$ and $O(D^{5/2}/m^{1/4})$ in (40) and (41) converges to zero. This result shows that an infinitely wide neural network could achieve optimal and near feasible solution like linear models. Next, even though our prediction model exhibits global optimality in (40) as $T$ approaches to infinity, the space over which the minimum is taken in (40) is

$$\{\theta \in \widetilde{\Theta} : \max_j \phi_j(\widetilde{R}_1(\theta), \ldots, \widetilde{R}_K(\theta)) \leq 0\},$$

instead of the original space

$$\{\theta \in \Theta : \max_j \phi_j(R_1(\theta), \ldots, R_K(\theta)) \leq 0\}.$$

This change is due to using surrogate functions $\widetilde{C}_k$ instead of the original function $C_k$. That is, we find a solution with the smallest error in a smaller space. On the other hand, even if we use surrogate functions for optimizing $\theta$, our prediction model exhibits the feasible solution in (41) in terms of the true loss function $C_k$, as we update $\xi$ and $\lambda$ by the original Lagrangian $\mathcal{L}_1$ and $\mathcal{L}$. Finally, there is a trade-off in the error bound for optimality in (40) and feasibility (41). It is not surprising that a larger $\kappa$ for penalizing the violation of constraints improves feasibility and results in a better error bound in (41). Nevertheless, increasing $\kappa$ hurts the convergence of Algorithm 2.

# 5 Experiment on COMPAS

In this section, we illustrate our algorithm and theory through analyzing a real dataset. COMPAS (Dressel and Farid, 2018) is a dataset provided by ProPublicas and containing criminal history, jail and prison time, demographics and COMPAS risk scores for defendants from Broward County. The goal of our experiments is to predict recidivism under fairness constraints. That is, we would like to ensure that some protected subgroups are treated equally by prediction models. We pre-processed the data and remove observations with some missing features and used one hot encoding for categorical features. After removing datapoints for which some of the features are missing, we have 6172 samples in COMPAS. Then data are randomly split into two groups: 70% of samples are for training and 30% samples are for testing. The classifier we use here is a 2 layer neural network with $m = 201$ hidden units. We restrict the weights in the ball with radius $D = 10$ and the Lagrange multipliers as $\|\lambda\|_\infty < \kappa = 1$. Cross-entropy loss is used as measurement of classification loss, and hinge loss is used as a differentiable convex surrogate of 0-1 loss.

In our experiment, the protected groups are two races, African-American and Caucasian, and we aim to treat each protected group equally. For comparison, we consider different methods as follows:

1. "Unconst." is the model without any constraint.

2. "$T$-Stoch." is the model with fairness constraint and trained by Algorithm 2. It has a set of $T$ parameters and uniformly pick one of those parameters to make a prediction.

3. "$J$-Stoch." is the model with fairness constraint and trained by Algorithm 2. $J$-Stoch also uses random perdition as $T$-Stochbut only use $J + 1$ parameters at most by the shrinking procedure in (27).

4. "Last" is the model with fairness constraint and trained by Algorithm 2 and uses the weight in the last iteration.

5. "Best" is the model with fairness constraint and trained by Algorithm 2 and uses the weight having the smallest loss.

First, to show the unfair treatment hiding in the classification results, we train a neural network without any constraint. The result is shown in the first line in Table 1. We can see that there is no significant difference in accuracy between black and white defendants. However, if we investigate more carefully, there is a huge difference in recall, defined as the classifier's positive classification rate, which means the classifier fails in different ways for two protected groups and tends to predict black defendants as more likely to re-offend. Note that removing the race feature cannot solve this unfair treatment and would reduce accuracy since there are others features correlated with race.

To address this issue, we add a constraint on equal opportunity

$$p^+(\mathcal{D}_B) = p^+(\mathcal{D}_W),$$

where $p^+(\mathcal{D}_\square) = \mathbb{E}_{(x,z)\sim\mathcal{D}_\square} \mathbb{I}\{y(\theta; x) > 0\}$ for $\square = B, W$ and $\mathcal{D}_B$ and $\mathcal{D}_W$ are the distributions that the target of instance is positive and the race of instance is African-American and Caucasian, respectively. We aim to optimize

$$\begin{aligned}
\min_{\theta \in \Theta} \quad & g(\theta) \\
\text{s.t.} \quad & p^+(\mathcal{D}_B) - p^+(\mathcal{D}_W) \leq 0 \\
& p^+(\mathcal{D}_W) - p^+(\mathcal{D}_B) \leq 0
\end{aligned} \tag{47}$$

using Algorithm 2. Heuristically, we only record the weights at the end of every epoch and discard the first 1000 iterations. The $T$-stochastic classifier would uniformly pick one of

|  | Train (%) | | | | Test (%) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Algo. | ACC | A(B) | A(W) | R(B) | R(W) | ACC | A(B) | A(W) | R(B) | R(W) |
| Unconst. | 74.42 | 72.94 | 74.88 | 59.07 | 39.62 | 65.60 | 64.66 | 67.63 | 50.0 | 28.87 |
| $J$-Stoch. | 75.23 | 71.80 | 77.84 | 55.60 | 55.40 | 65.98 | 64.97 | 64.87 | 48.82 | 41.42 |
| $T$-Stoch. | 72.31 | 68.14 | 74.14 | 49.86 | 51.97 | 63.44 | 62.59 | 61.34 | 46.27 | 41.00 |
| Last | 75.09 | 74.57 | 74.41 | 71.67 | 45.96 | 67.17 | 66.73 | 66.73 | 66.07 | 36.82 |
| Best | 75.99 | 73.30 | 77.97 | 61.51 | 53.68 | 66.73 | 66.01 | 66.01 | 55.29 | 40.16 |

Table 1: COMPAS Experiment Results. All numbers are represented by the percentage. ACC and A. stand for accuracy and R. stands for recall. B and W means black and white defendants, respectively.

those weights and make a prediction. The solution is further shrank using (27) to establish a $J$-stochastic classifier, which only uses 2 weights in our experiment. The results are shown in Table 1 and Figure 1.

Stochastic classifiers perform similarly to the unconstrained classifier in term of accuracy, as well as accuracy in different protected groups. Moreover, the recall of African-American aligns with the recall of Caucasian for stochastic classifiers in the training set. For the test set, we also observe improvement in the unfair treatment significantly. This means that indeed we fairly treat different subgroups, without losing any predictive power. Same phenomenon cannot be observed for "Last" and "Best" classifiers. Thus, the randomized model indeed is useful for meeting fairness constraints in general.

We also plot the loss and recalls of two protected groups for each of the iterates in Figure 1. We observe oscillation which is caused by violating constraints alternately and suggests that no pure equilibrium exists and that a stochastic classifier may be necessary.
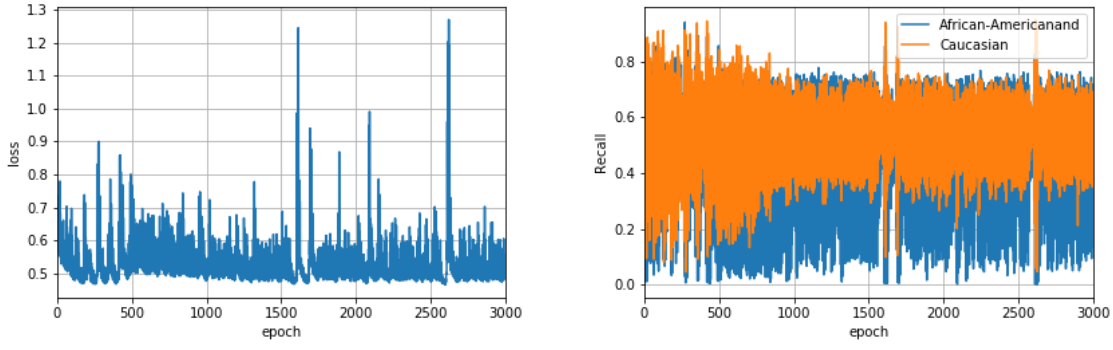
Figure 1: The plots for the loss and recall of two protected groups for each iteration during training on the COMPAS data set with equal opportunity constraints. We can see there is some oscillation in two plots. The oscillation caused by conflicting between two constraints in (47) and violating one of them alternately shows that there is no pure equilibrium to converge to.

The oscillation may commonly occur in practice, especially for optimizing non-convex and non-smooth Lagrangian.

We can understand why randomized models work for achieving the fairness constraint from this oscillation. Assume we could find two models with the same high accuracy but different behavior on recalls. One has a high recall on the black, and the other has a high recall on white defendants. If we randomly select one model to make the prediction, the average result will balance the recall on two races, which results in a fair classifier. Neural networks particularly fit this perspective since overparameterized neural networks can represent many different but high performing models corresponding to different parameters. We may find several candidate parameters with high accuracy but distinct prediction behaviors when adding various constraints. Then mixing candidate models indeed can im-

45

prove fairness via the previous argument. The game-theoretic framework can be interpreted as an automated strategy to encourage the neural network to explore different-behavior parameters.

# 6  Conclusion

This work shows how to provably train neural network models under fairness constraints using a game-theoretic framework. We modify the original game-theoretic framework to implement efficiently stochastic gradients gradient for large-scale datasets. Despite the difficulty of non-convex neural network models, we establish a no-regret bound for online learning of neural networks. Our results shed new light on the theoretical understanding of algorithmic fairness for neural networks.

One possible avenue for future work is developing an online procedure (Agrawal et al., 2014; Li and Ye, 2019) for shrinkage with theoretical guarantees, which may reduce memory costs also during training. Moreover, randomized prediction violates different fairness principles, such as "similar individuals receive similar outcomes" (Dwork et al., 2012). Cotter et al. (2019a) discussed clustering and ensemble procedures to decrease randomness while satisfying the fairness defined from an individual perspective, while it may be computationally intractable for large models. It may be interesting to find more memory and computationally efficient algorithms for neural networks.

# References

Agarwal, A., Beygelzimer, A., Dudik, M., Langford, J., and Wallach, H. (2018). A reductions approach to fair classification. In *International Conference on Machine Learning*, pages 60–69.

Agrawal, S., Wang, Z., and Ye, Y. (2014). A dynamic near-optimal algorithm for online linear programming. *Operations Research*, 62(4):876–890.

Alemohammad, S., Wang, Z., Balestriero, R., and Baraniuk, R. (2020). The recurrent neural tangent kernel. *arXiv preprint arXiv:2006.10246*.

Allen-Zhu, Z. and Li, Y. (2019). What can resnet learn efficiently, going beyond kernels? *Advances in Neural Information Processing Systems 32*, pages 9015–9025.

Allen-Zhu, Z., Li, Y., and Liang, Y. (2019a). Learning and generalization in overparameterized neural networks, going beyond two layers. In *Advances in neural information processing systems*, pages 6155–6166.

Allen-Zhu, Z., Li, Y., and Song, Z. (2019b). A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning*, pages 242–252.

Allen-Zhu, Z., Li, Y., and Song, Z. (2019c). On the convergence rate of training recurrent neural networks. *Advances in Neural Information Processing Systems 32*, pages 6673–6685.

Arora, S., Du, S. S., Hu, W., Li, Z., and Wang, R. (2019). Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. *arXiv preprint arXiv:1901.08584*.

Ba, J. L., Kiros, J. R., and Hinton, G. E. (2016). Layer normalization. *arXiv preprint arXiv:1607.06450*.

Bechavod, Y., Ligett, K., Roth, A., Waggoner, B., and Wu, S. Z. (2019). Equal opportunity in online classification with partial feedback. In *Advances in Neural Information Processing Systems*, pages 8972–8982.

Blum, A. and Lykouris, T. (2019). Advancing subgroup fairness via sleeping experts. *arXiv preprint arXiv:1909.08375*.

Blum, A. and Stangl, K. (2019). Recovering from biased data: Can fairness constraints improve accuracy? *arXiv preprint arXiv:1912.01094*.

Cai, Q., Yang, Z., Lee, J. D., and Wang, Z. (2019). Neural temporal-difference learning converges to global optima. *arXiv preprint arXiv:1905.10027*.

Celis, L. E., Huang, L., Keswani, V., and Vishnoi, N. K. (2019). Classification with fairness constraints: A meta-algorithm with provable guarantees. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 319–328.

Chen, R., Lucier, B., Singer, Y., and Syrgkanis, V. (2017). Robust optimization for non-convex objectives. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 4708–4717. Curran Associates Inc.

Chizat, L., Oyallon, E., and Bach, F. (2019). On lazy training in differentiable programming. In *Advances in Neural Information Processing Systems*, pages 2933–2943.

Chouldechova, A. (2017). Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 5(2):153–163.

Cotter, A., Gupta, M., and Narasimhan, H. (2019a). On making stochastic classifiers deterministic. *Advances in Neural Information Processing Systems*, pages 10910–10920.

Cotter, A., Jiang, H., and Sridharan, K. (2019b). Two-player games for efficient non-convex constrained optimization. In *Algorithmic Learning Theory*, pages 300–332.

Cotter, A., Jiang, H., Wang, S., Narayan, T., You, S., Sridharan, K., and Gupta, M. R. (2019c). Optimization with non-differentiable constraints with applications to fairness, recall, churn, and other goals. *Journal of Machine Learning Research*, 20(172):1–59.

Daskalaki, S., Kopanas, I., and Avouris, N. (2006). Evaluation of classifiers for an uneven class distribution problem. *Applied artificial intelligence*, 20(5):381–417.

Donini, M., Oneto, L., Ben-David, S., Shawe-Taylor, J. S., and Pontil, M. (2018). Empirical risk minimization under fairness constraints. In *Advances in Neural Information Processing Systems*, pages 2791–2801.

Dressel, J. and Farid, H. (2018). The accuracy, fairness, and limits of predicting recidivism. *Science advances*, 4(1):eaao5580.

Du, S., Lee, J., Li, H., Wang, L., and Zhai, X. (2019a). Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, pages 1675–1685.

Du, S. S., Hou, K., Salakhutdinov, R. R., Poczos, B., Wang, R., and Xu, K. (2019b). Graph neural tangent kernel: Fusing graph neural networks with graph kernels. In *Advances in Neural Information Processing Systems*, pages 5723–5733.

Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. (2012). Fairness through awareness. *Proceedings of the 3rd innovations in theoretical computer science conference*, pages 214–226.

Fan, J., Ma, C., and Zhong, Y. (2019). A selective overview of deep learning. *arXiv preprint arXiv:1904.05526*.

Fan, J., Wang, Z., Xie, Y., and Yang, Z. (2020). A theoretical analysis of deep q-learning. In *Learning for Dynamics and Control*, pages 486–489. PMLR.

Fard, M. M., Cormier, Q., Canini, K., and Gupta, M. (2016). Launch and iterate: Reducing prediction churn. In *Advances in Neural Information Processing Systems*, pages 3179–3187.

Gao, R., Cai, T., Li, H., Wang, L., Hsieh, C.-J., and Lee, J. D. (2019). Convergence of adversarial training in overparametrized networks. *arXiv preprint arXiv:1906.07916*.

Hardt, M., Price, E., and Srebro, N. (2016). Equality of opportunity in supervised learning. In *Advances in neural information processing systems*, pages 3315–3323.

Huang, R., Lattimore, T., György, A., and Szepesvári, C. (2017). Following the leader and fast rates in online linear prediction: Curved constraint sets and other regularities. *The Journal of Machine Learning Research*, 18(1):5325–5355.

Jacot, A., Gabriel, F., and Hongler, C. (2018). Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pages 8571–8580.

Kennedy, K., Mac Namee, B., and Delany, S. J. (2009). Learning without default: A study of one-class classification and the low-default portfolio problem. In *Irish Conference on Artificial Intelligence and Cognitive Science*, pages 174–187. Springer.

Kilbertus, N., Carulla, M. R., Parascandolo, G., Hardt, M., Janzing, D., and Schölkopf, B. (2017). Avoiding discrimination through causal reasoning. In *Advances in Neural Information Processing Systems*, pages 656–666.

Krogh, A. and Hertz, J. (1991). A simple weight decay can improve generalization. *Advances in neural information processing systems*, 4:950–957.

Lawrence, S., Burns, I., Back, A., Tsoi, A. C., and Giles, C. L. (2012). Neural network classification and prior class probabilities. In *Neural networks: Tricks of the trade*, pages 295–309. Springer.

Lee, J., Xiao, L., Schoenholz, S. S., Bahri, Y., Sohl-Dickstein, J., and Pennington, J. (2019). Wide neural networks of any depth evolve as linear models under gradient descent. *arXiv preprint arXiv:1902.06720*.

Li, X. and Ye, Y. (2019). Online linear programming: Dual convergence, new algorithms, and regret bounds. *arXiv preprint arXiv:1909.05499*.

Li, Y. and Liang, Y. (2018). Learning overparameterized neural networks via stochastic gradient descent on structured data. In *Advances in Neural Information Processing Systems*, pages 8157–8166.

Liu, B., Cai, Q., Yang, Z., and Wang, Z. (2019). Neural proximal/trust region policy optimization attains globally optimal policy. *arXiv preprint arXiv:1906.10306*.

Ma, R., Lin, Q., and Yang, T. (2019). Proximally constrained methods for weakly convex optimization with weakly convex constraints. *arXiv preprint arXiv:1908.01871*.

Narasimhan, H., Cotter, A., and Gupta, M. (2019). Optimizing generalized rate metrics through game equilibrium. *arXiv preprint arXiv:1909.02939*.

Neyshabur, B., Bhojanapalli, S., and Srebro, N. (2018a). A pac-bayesian approach to spectrally-normalized margin bounds for neural networks. In *International Conference on Learning Representations*.

Neyshabur, B., Li, Z., Bhojanapalli, S., LeCun, Y., and Srebro, N. (2018b). Towards understanding the role of over-parametrization in generalization of neural networks. *arXiv preprint arXiv:1805.12076*.

Oneto, L., Donini, M., and Pontil, M. (2019). General fair empirical risk minimization. *arXiv preprint arXiv:1901.10080*.

Oymak, S. and Soltanolkotabi, M. (2019). Towards moderate overparameterization: global convergence guarantees for training shallow neural networks. *arXiv preprint arXiv:1902.04674*.

Salimans, T. and Kingma, D. P. (2016). Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 901–909.

Salimi, B., Rodriguez, L., Howe, B., and Suciu, D. (2019). Capuchin: Causal database repair for algorithmic fairness. *arXiv preprint arXiv:1902.08283*.

Shalev-Shwartz, S. (2012). Online learning and online convex optimization. *Foundations and Trends in Machine Learning*, 4(2):107–194.

Sohrabi, M. K. and Azgomi, H. (2020). A survey on the combined use of optimization methods and game theory. *Archives of Computational Methods in Engineering*, 27(1):59–80.

Srebro, N., Sridharan, K., and Tewari, A. (2011). On the universality of online mirror descent. In *Advances in neural information processing systems*, pages 2645–2653.

Zafar, M. B., Valera, I., Rogriguez, M. G., and Gummadi, K. P. (2017). Fairness constraints: Mechanisms for fair classification. In *Artificial Intelligence and Statistics*, pages 962–970.

Zou, D., Cao, Y., Zhou, D., and Gu, Q. (2018). Stochastic gradient descent optimizes over-parameterized deep relu networks. *arXiv preprint arXiv:1811.08888*.