# BENCHMARKING STOCHASTIC APPROXIMATION ALGORITHMS FOR FAIRNESS-CONSTRAINED TRAINING OF DEEP NEURAL NETWORKS

#### **Anonymous authors**

Paper under double-blind review

# **ABSTRACT**

The ability to train Deep Neural Networks (DNNs) with constraints is instrumental in improving the fairness of modern machine-learning models. Many algorithms have been analysed in recent years, and yet there is no standard, widely accepted method for the constrained training of DNNs. In this paper, we provide a challenging benchmark of real-world large-scale fairness-constrained learning tasks, built on top of the US Census (Folktables, Ding et al. (2021)). We point out the theoretical challenges of such tasks and review the main approaches in stochastic approximation algorithms. Finally, we demonstrate the use of the benchmark by implementing and comparing three recently proposed, but as-of-yet unimplemented, algorithms both in terms of optimization performance, and fairness improvement. We will release the code of the benchmark as a Python package after peer-review.

# 1 Introduction

There has been a considerable interest in detecting and mitigating bias in artificial intelligence (AI) systems, recently. Multiple legislative frameworks, including the AI Act in the European Union, require the bias to be removed, but there is no agreement on what the correct definition of bias is or how to remove it. A natural translation of the requirement of removing bias into the formulation of training of deep neural network (DNN) utilizes constraints bounding the difference in empirical risk across multiple subgroups (Chen et al., 2018; Nandwani et al., 2019; Ravi et al., 2019). Over the past five years, there have been numerous algorithms proposed to solve convex and non-convex empirical-risk minimization (ERM) problems subject to constraints bounding the absolute value of empirical risk (Fang et al., 2024; Berahas et al., 2021; Curtis et al., 2024a; Oztoprak et al., 2023; Berahas et al., 2023; Na et al., 2023a;b; Bollapragada et al., 2023; Curtis et al., 2024b; Shi et al., 2022; Facchinei & Kungurtsev, 2023; Huang et al., 2025; Huang & Lin, 2023). Numerous other algorithms of this kind could be construed, based on a number of design choices, including:

- sampling techniques for the ERM objective and the constraints, either the same or different;
- use of first-order or higher-order derivatives, possibly in quasi-Newton methods;
- use of globalization strategies such as filters or line search;
- use of "true" globalization strategies including random initial points and random restarts in order to reach global minimizers.

Nevertheless, there is no single toolkit implementing the algorithms, which would allow for their easy comparison, and there is no benchmark to test the combinations of design choices on.

In this paper, we consider the constrained ERM problem:

$$\min_{x \in \mathbb{R}^n} \mathbb{E}[f(x,\xi)] \quad \text{s.t.} \quad \mathbb{E}[c(x,\zeta)] \le 0, \tag{1}$$

where  $\xi$  and  $\zeta$  are random variables. Further, we provide an automated way of constructing the ERM formulations out of a computation graph of a neural network defined by PyTorch or TensorFlow, the choice of the constraints (see Table 1), and a definition of the protected subgroups to apply the constraints to. Specifically, we provide means of utilizing the US Census data via the Python package Folktables, together with definitions of up to 5.7 billion protected subgroups. This presents a challenging benchmark in stochastic approximation for the constrained training of deep neural networks.

Table 1: Particular formulations of the constraint function c to enforce fairness.

Model	Our formulation
Demographic Parity Dwork et al. (2012) Equal opportunity Hardt et al. (2016) Equalized odds Hardt et al. (2016)	$\begin{split}  \mathbb{E}_{\mathcal{D}[\operatorname{group} A]}[\ell(f_{\theta}(X), Y)] - \mathbb{E}_{\mathcal{D}[\operatorname{group} B]}[\ell(f_{\theta}(X), Y)]  &\leq \delta \\  \mathbb{E}_{\mathcal{D}[\operatorname{group} A, Y = +]}[\ell(f_{\theta}(X), Y)] - \mathbb{E}_{\mathcal{D}[\operatorname{group} B, Y = +]}[\ell(f_{\theta}(X), Y)]  &\leq \delta \\ \sum_{v \in \{+, -\}}  \mathbb{E}_{\mathcal{D}[\operatorname{group} A, Y = v]}[\ell(f_{\theta}(X), Y)] - \mathbb{E}_{\mathcal{D}[\operatorname{group} B, Y = v]}[\ell(f_{\theta}(X), Y)]  &\leq \delta \end{split}$

#### Our contributions. The contributions of this paper are:

- a literature review of algorithms subject to handling (1);
- a toolbox that (i) implements four algorithms applicable in real-world situations, and (ii) provides an easy-to-use benchmark on real-world fairness problems;
- numerical experiments that compare these algorithms on a real-world dataset, and a comparison with alternative approaches to fairness.

**Paper structure.** The rest of the paper is organized as follows. Section 2 reviews related works and presents the relevant notions of fairness. Section 3 introduces the algorithms. Section 4 reports on our experiments. Section 5 concludes.

# 2 RELATED WORK, AND BACKGROUND IN FAIRNESS

In the literature on fairness, one distinguishes among pre-processing, in-processing, and post-processing. Pre-processing methods focus on modifying the training data to mitigate biases (Tawakuli & Engel, 2024; Du et al., 2021). In-processing methods enforce fairness during the training process by modifying the learning algorithm itself (Wan et al., 2023). Post-processing methods adjust the model's predictions after training (Kim et al., 2019). The constrained ERM approach (1) belongs to the class of in-processing methods.

In-processing methods include several approaches. One trend consists in jointly learning a predictor function and an adversarial agent that aims to reconstitute the subgroups from the predictor (Adel et al., 2019; Louppe et al., 2017; Madras et al., 2018; Edwards & Storkey, 2016). Another approach consists in adding "penalization" terms to the empirical risk term. These additional penalization terms, commonly referred to as regularizers, promote models that are a compromise between fitting the training data, and optimizing a fairness metric. Differentiable regularizers include, among others, HSIC (Li et al., 2022), Fairret (Buyl et al., 2024), or Prejudice Remover (Kamishima et al., 2012).

Closer to our setting, Cotter et al. (2019) consider minimizing the empirical risk subject to the so-called rate constraints based on the model's prediction rates on different datasets. These rates, derived from a dataset, give rise to non-convex, non-smooth, and large-scale inequality constraints akin to (1). Cotter et al. (2019) argue that hard constraints, although leading to a more difficult optimization problem, offer advantages over using a weighted sum of multiple penalization terms. Indeed, while the choice of weights for the penalization terms may depend on the dataset, specifying one constraint for each goal is easier for practitioners. In addition, a penalization-based model provides a predictor that balances minimizing the data-fit term and penalties in an opaque way, whereas a constraint-based model allows for a clearer understanding of the model design: minimizing the data-fit term subject to "hard" fairness constraints. Rate constraints differ from those in (1) in that they are piecewise-constant, rendering first-order methods unsuitable for solving them. We refer to the recent work of Ramirez et al. (2025) for a detailed argument on why constraining ERM problems is preferable to penalizing the ERM with multiple terms.

Major toolboxes for evaluating the fairness of models or for training models with fairness guarantees include AIF360 (Bellamy et al., 2018) and FairLearn (Bird et al., 2020). Delaney et al. (2024) compute the Pareto front of accuracy and fairness metrics for high-capacity models, and Buyl et al. (2024) provides differentiable fairness-inducing penalization terms. We also note the recent Cooper toolbox, closest to our setting, that focuses on Lagrangian-based methods (Gallego-Posada et al., 2025).

Le Quy et al. (2022) provides a detailed survey of fairness-oriented datasets, and Ding et al. (2021) derives new datasets. The benchmark of Han et al. (2023) reviews the existence of biases in prominent datasets, finding that "not all widely used fairness datasets stably exhibit fairness issues", and assesses

the performance of a range of in-processing methods in addressing biases, focusing on differentiable minimization only. Other benchmarks of fairness methods include Defrance et al. (2024); Fabris et al. (2022); Pessach & Shmueli (2022); Chen et al. (2024). Statistical aspects of the fairness-constrained Empirical Risk Minimization have only been considered recently; see e.g. Chamon et al. (2022).

The template problem (1) encompasses fairness-enforcing approaches that find applications in highrisk domains, such as credit scoring, hiring processes, medicine and healthcare (Chen et al., 2023), ranking and recommendation (Pitoura et al., 2022), but also in forecasting the observations of linear dynamical systems (Zhou et al., 2023b), or in two-sided economic markets (Zhou et al., 2023a). In addition, solving (1) is of interest in other fields, such as compression of neural networks (Chen et al., 2018), improving statistical performance of neural networks (Nandwani et al., 2019; Ravi et al., 2019), or the training of neural networks with constraints on the Lipschitz bound (Pauli et al., 2021). We note that all the aforementioned methodologies feature large-scale constraints.

**Deep neural networks (DNNs).** Consider a dataset of N observations  $\mathcal{D} = \{(X_i, Y_i), i = 1, ..., N\}$ . We seek some function  $f_{\theta}$  such that  $f_{\theta}(X_i) \approx Y_i$ . A typical formulation of this task is the following regression problem:

$$\min_{\theta \in \mathbb{R}^n} \frac{1}{N} \sum_{i=1}^N \ell(f_{\theta}(X_i), Y_i) + \mathcal{R}(\theta). \tag{2}$$

Here,  $\ell: \mathbb{R} \times \mathbb{R} \to \mathbb{R}$  is a loss function, such as the logistic loss  $\ell(y;z) = \log(1+e^{-yz})$ , the hinge loss  $\ell(y;z) = \max\{0,1-yz\}$ , the absolute deviation loss  $\ell(y;z) = |y-z|$ , or the square loss  $\ell(y;z) = \frac{1}{2}(y-z)^2$ . The term  $\mathcal{R}$  is a regularizer, and  $f_\theta$  is a deep neural network (DNN) of depth L with parameters  $\theta$ . The DNN  $f_\theta$  is defined recursively, for some input X, as

$$a_0 = X,$$
  $a_i = \rho_i(V_i(\theta)a_{i-1}), \text{ for every } i = 1, ..., L,$   $f_{\theta}(X) = a_L,$  (3)

where  $V_i(\cdot)$  are linear maps into the space of matrices, and  $\rho_i$  are activation functions applied coordinate-wise, such as ReLU  $\max(0,t)$ , quadratics  $t^2$ , hinge losses  $\max\{0,t\}$ , and SoftPlus  $\log(1+e^t)$ . A dataset  $\mathcal D$  is described by attributes (or features), such as age, income, gender, etc. The attribute which the DNN is trained to predict is called the class attribute. We denote the class attribute by Y, whereas the predicted value given by the DNN is denoted by  $\hat{Y}$ . Both Y and  $\hat{Y}$  are binary and take values in  $\{+,-\}$ .

Fairness-aware learning applied to DNNs. The goal of this approach is to reduce discriminatory behavior in the predictions of a DNN across different demographic groups (e.g., male vs. female). The demographic groups are also reffered to as subgroups. The attributes such as race or gender which must be handled cautiously are called protected. We denote by S the protected attribute,  $s_1, \ldots, s_m$  its possible values, and  $\mathcal{D}[s_i]$  the observations in  $\mathcal{D}$  such that  $S = s_i$ . A way to impose fairness on the learned predictor is to equip (2) with suitable constraints. Some possible constraint choices are shown in Table 1. Choosing loss difference bound as the constraint, denoting  $\ell^{s_i}(\theta) = \frac{1}{|\mathcal{D}[s_i]|} \sum_{X,Y \in \mathcal{D}[s_i]} \ell(f_{\theta}(X),Y)$  for  $i=1,\ldots,m$ , and setting  $\delta > 0$  yields formulation:

$$\min_{\theta \in \mathbb{R}^n} \quad \frac{1}{N} \sum_{i=1}^N \ell(f_{\theta}(X_i), Y_i) + \mathcal{R}(\theta)$$
s.t. 
$$-\delta \le \ell^{s_i}(\theta) - \frac{1}{m} \sum_{j=1}^m \ell^{s_j}(\theta) \le \delta, \quad i = 1, \dots, m.$$
(4)

Bounding the distance between subgroup losses yields m constraints. Other formulations are possible, such as bounding the distance between every pair of subgroup, providing simpler individual constraints, but in greater number (m(m+1)/2). Formulation (4) extends to several protected attributes by adding the corresponding set of equations; we omit this direct generalisation for clarity.

**Fairness metrics.** There exist tens of fairness metrics (Verma & Rubin, 2018). However, Barocas et al. (2023, Ch. 3) pointed out that most fairness metrics are combinations of three elementary fairness criteria: independence, separation, and sufficiency. These criteria cannot be minimized simultaneously, and there is a trade-off between attaining the elementary fairness metrics and the

Table 2: Three elementary notions of fairness

Independence Separation Sufficiency  $\frac{1}{m}\sum_{i=1}^{m}\left|P_{i}^{\text{ind}}-\frac{1}{m}\sum_{j}P_{j}^{\text{ind}}\right| \quad \frac{1}{2}\sum_{v\in\{+,-\}}\frac{1}{m}\sum_{i=1}^{m}\left|P_{i,v}^{\text{Sp}}-\frac{1}{m}\sum_{j}P_{j,v}^{\text{Sp}}\right| \quad \frac{1}{2}\sum_{v\in\{+,-\}}\frac{1}{m}\sum_{i=1}^{m}\left|P_{i,v}^{\text{Sf}}-\frac{1}{m}\sum_{j}P_{j,v}^{\text{Sf}}\right|$ 

prediction accuracy, i.e., the probability that the predicted value is equal to the actual value. Thus, we seek an optimal trade-off between attaining the fairness metrics and minimizing the prediction inaccuracy. We next recall the definitions of these three relevent metrics, following Barocas et al. (2023), and provide the formula for computing them in Table 2.

**Independence (Ind)** This fairness criterion requires the prediction  $\hat{Y}$  to be statistically independent of the protected attribute S. Equivalent definitions of independence for a binary classifier  $\hat{Y}$  are referred to as statistical parity (SP), demographic parity, and group fairness. Independence is the simplest criterion to work with, both mathematically and algorithmically. In a binary classification task, independence implies the equality of  $P_i^{\text{ind}} = P(\hat{Y} = + \mid S = s_i)$  for all  $i = 1, \ldots, m$ .

**Separation** (Sp) Unlike independence, the separation criterion requires the prediction  $\hat{Y}$  to be statistically independent of the protected attribute S, given the true label Y. The separation criterion also appears under the name Equalized odds (EO). In a binary classification task, the separation criterion requires that all groups experience the same true negative rate and the same true positive rate. Formally, we require the equality of  $P_{i,v}^{\mathrm{Sp}} = P(\hat{Y} = + \mid S = s_i, Y = v)$  for every  $i = 1, \ldots, m$ , and  $v \in \{+, -\}$ .

**Sufficiency (Sf)** The sufficiency criterion is satisfied if the true label Y is statistically independent of the protected attribute S, given the prediction  $\hat{Y}$ . In a binary classification task, the sufficiency criterion requires a parity of positive and negative predictive values across the groups. Formally, we require the equality of  $P_{v,s}^{\mathrm{Sf}} = P(Y = + \mid \hat{Y} = v, S = s)$ , for every  $i = 1, \ldots, m$ , and  $v \in \{+, -\}$ .

# 3 ALGORITHMS

We recall that we consider the optimization problem

$$\min_{x \in \mathbb{R}^n} F(x) \quad \text{s.t.} \quad C(x) \le 0, \tag{5}$$

where the functions  $F: \mathbb{R}^n \to \mathbb{R}$  and  $C: \mathbb{R}^n \to \mathbb{R}^m$  are defined as expectations of functions f and c, which depend on random variables  $\xi$  and  $\zeta$ , respectively. Solving (5) has the following challenges:

- large-scale objective and constraint functions, which require sampling schemes,
- the necessity of incorporating inequality constraints, not merely equality constraints (see fairness formulations in Table 1),
- the necessity to cope with the nonconvexity and nonsmoothness of F and C, due to the presence
  of neural networks.

In this section, we identify the algorithms that address these challenges most precisely. However, we note that there exists currently no algorithm with guarantees for such a general setting.

**Recalls and notation.** We denote the projection of a point x onto a set  $\mathcal{X}$  by  $\operatorname{proj}_{\mathcal{X}}(x) = \arg\min_{v \in \mathcal{X}} \|x - v\|^2$ . We denote by  $N \sim \mathcal{G}(p_0)$  sampling a random variable from the geometric distribution with a parameter  $p_0$ , i.e., the probability that N = n equals  $(1 - p_0)^n p_0$  for  $n \geq 0$ . We distinguish between the random variable  $\xi$  associated with the objective function and the random variable  $\zeta$  associated with the constraint function. Their probability distributions are denoted by  $\mathcal{P}_{\xi}$  and  $\mathcal{P}_{\zeta}$ . For an integer  $J \in \mathbb{N}$ , a set  $\{\xi_j\}_{j=1}^J$  of independent and identically distributed random

variables  $\xi_1, \dots, \xi_J \stackrel{iid}{\sim} \mathcal{P}_{\xi}$  is called a mini-batch. Inspired by Na et al. (2023a), we use the following notation for the stochastic estimates computed from a mini-batch of size J:

$$\overline{\nabla}^{J} f(x) = \frac{1}{J} \sum_{j=1}^{J} \nabla f(x, \xi_j), \quad \overline{c}^{J}(x) = \frac{1}{J} \sum_{j=1}^{J} c(x, \zeta_j), \quad \overline{\nabla}^{J} c(x) = \frac{1}{J} \sum_{j=1}^{J} \nabla c(x, \zeta_j). \quad (6)$$

	Objective function $F$				Constraint function $C$						
Algorithm	stochastic	weakly convex	$\mathcal{C}^1$ with Lipschitz $ abla F$	tame loc. Lipschitz	stochastic	C(x) = 0	$C(x)=0 \text{ and } C(x) \leq 0$	linear	weakly convex	$\mathcal{C}^1$ with Lipschitz $ abla C$	tame loc. Lipschitz
SGD	1	<b>(/</b> )	(✔)	1							
Berahas et al. (2023) Fang et al. (2024) Curtis et al. (2024a)	1	_	1	_	-	1	_	-	_	1	_
Na et al. (2023a)	1	_	$\checkmark(\mathcal{C}^3)$	_	-	✓	_	-	_	$\checkmark(C^3)$	_
Shi et al. (2022) Curtis et al. (2024b)	/	_	<b>√</b>	_	-	$(\checkmark)$	/	_	_	<b>✓</b>	_
Na et al. (2023b)	1	-	$\mathcal{I}(\mathcal{C}^2)$	_	-	<b>(/</b> )	✓	_	-	$\mathcal{I}(\mathcal{C}^2)$	-
Bollapragada et al. (2023)	<b>/</b>	-	<b>√</b> (+ cvx)	_	-	<b>/</b>	_	/	_	-	_
Oztoprak et al. (2023)	1	_	<b>V</b>	_	/	<b>/</b>	-	_	_	/	_
SSL-ALM Huang et al. (2025)	<b>'</b> ,	_	<i>V</i>	_	\ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \	( <b>V</b> )	<b>\</b>	/	_	_	_
Stoch. Ghost Facchinei & Kungurtsev (2023)	~	_	✓	_	· •	(V)	_	I —	_	✓	_

## 3.1 REVIEW OF METHODS FOR CONSTRAINED ERM

We compare recent constrained optimization algorithms considering a stochastic objective function in Table 3. We note that most of them do not consider the case of stochastic constraints. Among those which do consider stochastic constraints, only three admit inequality constraints. Moreover, with the exception of Huang & Lin (2023), all the algorithms in Table 3 assume F to be at least  $C^1$ , which makes addressing the challenge of nonsmoothness of F infeasible. Davis et al. (2018) leads us to conclude that assuming the objective and constraint functions to be tame and locally Lipschitz is a suitable requirement for solving (5) with theoretical guarantees of convergence. At this point, however, no such algorithm exists, to the best of our knowledge.

Consequently, we consider the practical performance of the algorithms that address the challenges of solving (5) most closely: Stochastic Ghost, SSL-ALM, and Stochastic Switching Subgradient.

#### 3.2 STOCHASTIC GHOST METHOD (STGH)

Facchinei & Kungurtsev (2023) propose the Stochastic Ghost method, that combines a deterministic method for solving (1) (Facchinei et al., 2021) with a stochastic sampling approch for nonlinear maps (Blanchet et al., 2019). The deterministic method of Facchinei et al. (2021) consists in solving subproblem (7) to obtain a direction d, and then to preform a line search. Here,  $e \in \mathbb{R}^m$  is a vector with all elements equal to one,  $\tau$  and  $\beta > 0$  are user-prescribed constants, and  $\kappa_k$  is defined as a certain convex combination of optimization subproblems related to C and  $\nabla C$ . The definition of  $\kappa_k$  enables to expand the feasibility region so that (7) is always feasible. As the problem (1) is stochastic, the subproblem (7) is modified to a stochastic version (8), using the notation in (6):

$$\min_{d} \quad \nabla F(x_k)^\top d + \frac{\tau}{2} \|d\|^2, \qquad \qquad \min_{d} \quad \overline{\nabla}^J f(x_k)^\top d + \frac{\tau}{2} \|d\|^2,$$
s.t. 
$$C(x_k) + \nabla C(x_k)^\top d \le \kappa_k e, \qquad (7) \qquad \text{s.t.} \quad \overline{c}^J(x_k) + \overline{\nabla}^J c(x_k)^\top d \le \overline{\kappa_k}^J e, \qquad (8)$$

$$\|d\|_{\infty} \le \beta, \qquad \qquad \|d\|_{\infty} \le \beta.$$

In the stochastic setting (8), an unbiased estimate  $d(x_k)$  of the line search direction d is computed using four particular mini-batches as follows. To facilitate comprehension, we denote  $X_k^J = \{X_{k,j}\}_{j=1}^J$  a mini-batch of size J with the j-th element  $X_{k,j} = (\nabla f(x_k, \xi_{k,j}), c(x_k, \zeta_{k,j}), \nabla c(x_k, \zeta_{k,j}))$ . First, we sample a random variable  $N \sim \mathcal{G}(p_0)$  from the geometric distribution. Then we sample the minibatches  $X_k^1$  and  $X_k^{2^{N+1}}$  and we partition the mini-batch  $X_k^{2^{N+1}}$  of size  $2^{N+1}$  into two mini-batches odd $(X_k^{2^{N+1}})$  and even $(X_k^{2^{N+1}})$  of size  $2^N$ . Finally, we solve (8) for each of the four mini-batches,

denoting by  $d(x_k; X_k^J)$  the solution of (8) for the corresponding mini-batch  $X_k^J$ . We obtain

$$d(x_k) = \frac{d(x_k; X_k^{2^{N+1}}) - \frac{1}{2} \left( d(x_k; \operatorname{odd}(X_k^{2^{N+1}})) + d(x_k; \operatorname{even}(X_k^{2^{N+1}})) \right)}{(1 - p_0)^N p_0} + d(x_k; X_k^1).$$
(9)

An update between the iterations  $x_k$  and  $x_{k+1}$  is then computed as  $x_{k+1} = x_k + \alpha_k d(x_k)$ , where the deterministic stepsize  $\alpha_k$  should be square-summable  $\sum_{k=1}^{\infty} \alpha_k^2 < \infty$  but not summable  $\sum_{k=1}^{\infty} \alpha_k = \infty$ . For more details, see Algorithm 1 (Appendix A).

#### 3.3 STOCHASTIC SMOOTHED AND LINEARIZED AL METHOD (SSL-ALM)

The Stochastic Smoothed and Linearized AL Method (SSL-ALM) was described in Huang et al. (2025) for optimization problems with stochastic linear constraints. Although problem (1) has nonlinear inequality constraints, we use the SSL-ALM due to the lack of algorithms in the literature dealing with stochastic non-linear constraints; see Table 3. The transition between equality and inequality constraints is handled with slack variables. Following the structure of Huang et al. (2025), we minimize over the set  $\mathcal{X} = \mathbb{R}^n \times \mathbb{R}^m_{\geq 0}$ . The method is based on the augmented Lagrangian (AL) function  $L_\rho(x,y) = F(x) + y^\top C(x) + \frac{\rho}{2} \|C(x)\|^2$ ; see e.g., (Bertsekas & Rheinboldt, 2014). Adding a smoothing term with an additional variable  $z \in \mathbb{R}^n$  yields the proximal AL function

$$K_{\rho,\mu}(x,y,z) = L_{\rho}(x,y) + \frac{\mu}{2} ||x-z||^2.$$

The SSL-ALM method was originally proposed in Huang et al. (2025) where it is interpreted as an inexact gradient descent step on the Moreau envelope. An important property of the Moreau envelope is that its stationary points coincide with those of the original function.

The strength of this method is that, as opposed to the Stochastic Ghost method, it does not use large mini-batch sizes. In each iteration, we sample  $\xi \stackrel{iid}{\sim} \mathcal{P}_{\xi}$  to evaluate the objective and  $\zeta_1$ ,  $\zeta_2 \stackrel{iid}{\sim} \mathcal{P}_{\zeta}$  to evaluate the constraint function and its Jacobian matrix, respectively. The function

$$G(x, y, z; \xi, \zeta_1, \zeta_2) = \nabla f(x, \xi) + \nabla c(x, \zeta_1)^\top y + \rho \nabla c(x, \zeta_1)^\top c(x, \zeta_2) + \mu(x - z)$$
 (10)

is defined so that, in iteration k,  $\mathbb{E}_{\xi,\zeta_1,\zeta_2}[G(x_k,y_{k+1},z_k;\xi,\zeta_1,\zeta_2)] = \nabla K_{\rho,\mu}(x_k,y_{k+1},z_k)$ . Denoting  $\eta,\tau$ , and  $\beta$  positive parameters, the update is

$$y_{k+1} = y_k + \eta c(x, \zeta_1),$$
  

$$x_{k+1} = \text{proj}_{\mathcal{X}}(x_k - \tau G(x_k, y_{k+1}, z_k; \xi, \zeta_1, \zeta_2)),$$
  

$$z_{k+1} = z_k + \beta(x_k - z_k).$$
(11)

For more details, see Algorithm 2 (Appendix A).

# 3.4 STOCHASTIC SWITCHING SUBGRADIENT METHOD (SSW)

The Stochastic Switching Subgradient method was described in Huang & Lin (2023) for optimization over a closed convex set  $\mathcal{X} \subset \mathbb{R}^d$  which is easy to project on. It allows for weakly-convex, possibly nonsmooth, objective and constraint functions. They consider subgradients instead of gradients.

The algorithm relies on a prescribed sequence of infeasibility tolerances  $\epsilon_k$  and of stepsizes  $\eta_k^f$  and  $\eta_k^c$ . At iteration k, we sample  $\zeta_1, \ldots, \zeta_J \overset{iid}{\sim} \mathcal{P}_{\zeta}$  to compute  $\overline{c}^J(x_k)$ . If  $\overline{c}^J(x_k)$  is smaller than  $\epsilon_k$ , we sample  $\xi \overset{iid}{\sim} \mathcal{P}_{\xi}$  and update using a stochastic estimate  $S^f(x_k, \xi)$  of a subgradient of F:

$$x_{k+1} = \operatorname{proj}_{\mathcal{X}}(x_k - \eta_k^f S^f(x_k, \xi)).$$

If not, we sample  $\zeta \stackrel{iid}{\sim} \mathcal{P}_{\zeta}$  and update using a stochastic estimate  $S^c(x_k, \zeta)$  of a subgradient of C:

$$x_{k+1} = \operatorname{proj}_{\mathcal{X}}(x_k - \eta_k^c S^c(x_k, \zeta)).$$

In either case, the updates are only saved starting from a prescribed index  $k_0$  and the final output is sampled randomly from the saved updates. The algorithm presented here is slightly more general than the one presented in Huang & Lin (2023): we allow for different stepsizes for the objective and the constraint update, while the original method employs equal stepsizes  $\eta_k^f = \eta_k^c$ . For more details, see Algorithm 3 (Appendix A).

# 4 EXPERIMENTAL EVALUATION

In this section, we illustrate the presented algorithms on a real-world instance of the ACS dataset, comparing how they fare with optimization and fairness metrics.

#### 4.1 Dataset for fair ML

Ding et al. (2021) proposed a large-scale dataset for fair Machine Learning, based on the ACS PUMS data sample (American Community Survey Public Use Microdata Sample). The ACS survey is sent annually to approximately 3.5 million US households in order to gather information on features such as ancestry, citizenship, education, employment, or income. Therefore, it has the potential to give rise to large-scale learning and optimization problems.

We use the ACSIncome dataset over the state of Oklahoma, and choose the binary classification task of predicting whether an individual's income is over \$50,000. The dataset contains 9 features and 17,917 data points, and may be accessed via the Python package Folktables. We choose race (**RAC1P**) as the protected attribute. In the original dataset, it is a categorical variable with 9 values. For the purposes of this experiment, we binarize it to obtain the non-protected group of "white" people and the protected group of "non-white" people. The dataset is split randomly into train (80%, 14,333 points) and test (20%, 3,584 points) subsets and it is stratified with respect to the protected attribute, i.e., the proportion of "white" and "non-white" samples in the training and test sets is equivalent to that in the full dataset (30.8% of positive labels in group "white", 20.7% in the group "non-white"). The protected attributes are then removed from the data so that the model cannot learn from them directly. The data is normalized using Scikit-Learn StandardScaler.

Note that ACSIncome is a real-world dataset for which ERM-based predictors without fairness safeguards are known to learn biases (Han et al., 2023). Accordingly, Table 4 (line 1) shows that an ERM predictor without fairness safeguards has poor fairness metrics; see also Figure 4.

#### 4.2 EXPERIMENTS

**Numerical setup.** Experiments are conducted on an Asus Zenbook UX535 laptop with AMD Ryzen 7 5800H CPU, and 16GB RAM, using Python with the PyTorch package (Paszke et al., 2019).

**Problems.** We consider the constrained ERM problem (4) without any regularization  $\mathcal{R}=0$ , and, as baselines, the ERM problem (2) without any regularization,  $\mathcal{R}=0$ , and with a fairness inducing regularizer  $\mathcal{R}$  that promotes small difference in accuracy between groups, provided by the Fairret library (Buyl et al., 2024). In all problems, we take as loss function the Binary Cross Entropy with Logits Loss

$$\ell(f_{\theta}(X_i), Y_i) = -Y_i \cdot \log \sigma(f_{\theta}(X_i)) - (1 - Y_i) \cdot \log(1 - \sigma(f_{\theta}(X_i))), \tag{12}$$

where  $\sigma(z) = \frac{1}{1+e^{-z}}$  is the sigmoid function, and the prediction function  $f_{\theta}$  is a neural network with 2 interconnected hidden layers of sizes 64 and 32 and ReLU activation, with a total of 194 parameters.

Algorithms and parameters. We assess the performance of four algorithms for solving the constrained problem (4): (1) Stochastic Ghost (StGh) (Sec. 3.2 - parameters  $p_0=0.4$ ,  $\alpha_0=0.05$ ,  $\rho=0.8$ ,  $\tau=1$ ,  $\beta=10$ ,  $\lambda=0.5$ ,  $\hat{\alpha}=0.05$ , (2) SSL-ALM (Sec. 3.3 - parameters  $\mu=2.0$ ,  $\rho=1.0$ ,  $\tau=0.01$ ,  $\eta=0.05$ ,  $\beta=0.5$ ,  $M_y=10$ ), (3) plain Augmented Lagrangian Method ALM (Sec. 3.3, smoothing term removed  $\mu=0$ , otherwise the same setting as SSL-ALM), and (4) Stochastic Switching Subgradient (SSw) (Sec. 3.4 -  $\eta_k^f=0.5$ ,  $\eta_k^c=0.05$ ,  $\epsilon_k=10^{-4}$  if k<500,  $\epsilon_k=0.97\epsilon_{k-1}$  for every  $k\geq 500$  at each epoch). We also provide the behavior of SGD for solving the ERM problem, both with no fairness safeguards (SGD), and with fairness regularization on accuracy as provided by the Fairret library (Buyl et al., 2024) (SGD-Fairret). These methods serve as baselines. When estimating the constraints, we sample an equal number of data points for every subgroup.

**Optimization performance.** Figures 1 and 2 present the evolution of loss and constraint values over the train and test datasets for the four algorithms addressing the constrained problem (columns 1–4),

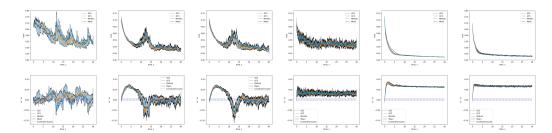


Figure 1: Train loss and constraint values (first and second row) over time (s) on the ACS Income dataset for each algorithm. From left to right: StGh, SSL-ALM, ALM, SSw, SGD, SGD-Fairret.

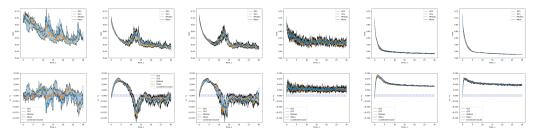


Figure 2: Test loss and constraint (first and second row) values over time (s) on the ACS Income dataset for each algorithm. From left to right: StGh, SSL-ALM, ALM, SSw, SGD, SGD-Fairret.

as well as for the two baselines: SGD without fairness (col. 5), and SGD with fairness regularization (col. 6). Each algorithm is run 10 times, and the plots display the mean, median, and quartiles values.

To a certain extent, the four algorithms (col. 1–4) succeed in minimizing the loss and satisfying the constraints on the train set. The AL-based methods (col. 2 and 3) demonstrate a better behavior compared to StGh and SSw; indeed, StGh exhibits higher variability in both loss and constraint values (col. 1), while SSw fails to satisfy the constraints within the required bounds (col. 4). We were unable to identify parameter settings for SSw that simultaneously satisfy the constraints and minimize the objective function. The ERM baselines (col. 5 and 6) exhibit lower variability in the trajectories, and minimize the loss in less time, but as expected, they do not satisfy the constraints.

The ALM and SSL-ALM schemes are the closest to satisfying the constraints on the train set. On the test set, however, they are slightly biased towards negative values. Such bias is expected on unseen data and reflects the generalization behavior of fairness-constrained estimators. This is beyond the scope of the current work; see e.g. Chamon et al. (2022).

**Fairness performance.** Figure 3 presents the distribution of predictions over both groups. The distribution of prediction without fairness guarantees (col. 5) clearly does not meet the group fairness standard. Indeed, the "non-white" group has a significantly higher likelihood than the "white" group of receiving small predicted values, and the converse holds for large predicted values. The SGD-Fairret model (col. 6) lies between the four constrained models and SGD. Among the fairness-constrained models, the ALM and SSL-ALM distributions are the closest to the distributions of SGD without fairness, which is consistent with retaining good prediction information. The four models that approximately solve the fairness formulation (col. 1–4) all have closer distributions across groups. Numerically, this is expressed in Table 4 (col. Wd), which reports the value of the Wasserstein distance between group distributions for each model.

Table 4: Fairness metrics (independence, separation, sufficiency), inaccuracy, and Wasserstein distances between groups (Wd) for the four constrained estimators and the two baselines.

	Train						Test						
Algname	Ind	Sp	Ina	Sf	Wd	Ind	Sp	Ina	Sf	Wd			
SGD	$0,094 \pm 0,004$	$0,132\pm0,007$	$0,\!201\!\pm\!0,\!001$	$0,115\pm0,006$	$0,008\pm0,000$	$0,097\pm0,006$	$0,\!176\pm0,\!016$	$0,215\pm0,002$	$0,\!171\pm\!0,\!009$	$0,008\pm0,000$			
StGh ALM SSL-ALM SSw	$\begin{array}{c} \textbf{0,048} {\pm 0,026} \\ 0,058 {\pm 0,007} \\ 0,066 {\pm 0,009} \\ 0,077 {\pm 0,029} \end{array}$	$\begin{array}{c} \textbf{0,049} {\pm 0,028} \\ 0.061 {\pm 0,016} \\ 0.071 {\pm 0,015} \\ 0.115 {\pm 0,029} \end{array}$	$\begin{array}{c} 0,273 \pm 0,024 \\ 0,240 \pm 0,012 \\ 0,233 \pm 0,017 \\ 0,224 \pm 0,017 \end{array}$	$\begin{array}{c} 0,200{\pm}0,038 \\ 0,197{\pm}0,011 \\ 0,186{\pm}0,013 \\ 0,133{\pm}0,015 \end{array}$	$\begin{array}{c} 0,002 \pm 0,001 \\ 0,003 \pm 0,000 \\ 0,003 \pm 0,001 \\ \textbf{0,001} \pm \textbf{0,001} \end{array}$	$\begin{array}{c} \textbf{0,049} {\pm 0,029} \\ 0,058 {\pm 0,012} \\ 0,066 {\pm 0,011} \\ 0,080 {\pm 0,029} \end{array}$	$\begin{array}{c} \textbf{0,096} {\pm 0,039} \\ 0,114 {\pm 0,014} \\ 0,117 {\pm 0,023} \\ 0,144 {\pm 0,050} \end{array}$	$\begin{array}{c} 0,276 \pm 0,022 \\ 0,244 \pm 0,007 \\ 0,240 \pm 0,012 \\ 0,229 \pm 0,013 \end{array}$	$\begin{array}{c} 0,211{\pm}0,033 \\ 0,221{\pm}0,017 \\ 0,215{\pm}0,022 \\ 0,175{\pm}0,031 \end{array}$	$\begin{array}{c} 0,003\pm0,002\\ 0,003\pm0,001\\ 0,004\pm0,001\\ \textbf{0,002}\pm\textbf{0,001} \end{array}$			
SGD-Fairret	$0,091 \pm 0,012$	$0,121\pm0,017$	$0,201 \pm 0,002$	0,106±0,010	$0,005\pm0,001$	0,094±0,010	$0,174\pm0,019$	$0,\!213\!\pm\!0,\!002$	$0,180\pm0,022$	0,006±0,001			

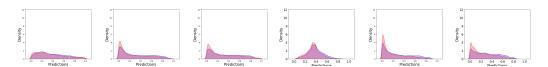


Figure 3: Distribution of predictions for each algorithm. Left to right: StGh, SSL-ALM, ALM, SSw, SGD, SGD-Fairret. Blue and red denote "white" and "non-white" groups.

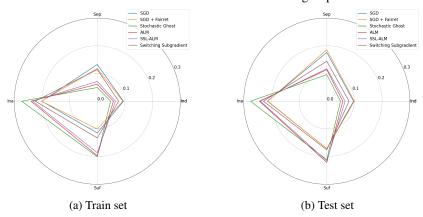


Figure 4: Average value of the three fairness metrics (independence (Ind), separation (Sp), and sufficiency (Sf)), along with inaccuracy (Ina). For all metrics, smaller values are better.

Table 4 displays the fairness metrics presented in Section 2: independence (Ind), separation (Sp), and sufficiency (Sf), along with inaccuracy (Ina). The mean value and standard deviation over 10 runs are presented for the four fairness-constrained models and the two baselines, both on train and test sets. Figure 4 presents the mean values as spider plots. For all metrics, smaller is better.

Among the four fairness-constrained models, StGh performs best in terms of independence and separation, but worst in terms of accuracy. SSw achieves fairness and accuracy metrics that have intermediate values relative to those of the unconstrained SGD model, and those of the other constrained models. This is consistent with the observation that the optimization method, with our choice of parameters, favored minimizing the objective over satisfying the constraints. The ALM and SSL-ALM methods provide the best compromise: they improve independence and separation relative to the SGD model, while moderately degrading accuracy. SGD-Fairret slightly improves sufficiency relative to the SGD model. The four models constrained in the difference of loss between subgroups have higher values of sufficiency. Similar observations hold for metrics on the test set.

For completeness, we report in Appendix B an additional experiment with one protected attribute that takes five values, and compare the optimization performance of the three algorithms for constrained minimization with two baselines.

## 5 Conclusion

To the best of our knowledge, this paper provides the first benchmark for assessing the performance of optimization methods on real-world instances of fairness constrained training of models. We highlight the challenges of this approach, namely that objective and constraints are non-convex, non-smooth, and large-scale, and review the performance of four practical algorithms.

# LIMITATIONS

Our work identifies that there is currently no algorithm with guarantees for solving the fairness constrained problem. Above all, we hope that this work, along with the Python toolbox for easy benchmarking of new optimization methods, will stimulate further interest in this topic. Also, we caution readers that the method present here is not a silver-bullet that handles all biases and ethical issues of training ML models. In particular, care must be taken that fair ML is part of a interdisciplinary pipeline that integrates the specifics of the use-case, and that it does not serve as an

excuse for pursuing Business-As-Usual policies that fail to tackle ethical issues (Balayn et al., 2023; Wachter et al., 2021).

#### REPRODUCIBILITY

Code to reproduce the experiments is provided in the Supplementary Material. This includes a readme file with instructions to reproduce experiments. Details on the computing environment are provided in Section 4.

# REFERENCES

- Tameem Adel, Isabel Valera, Zoubin Ghahramani, and Adrian Weller. One-network adversarial fairness. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 2412–2420, 2019.
- Agathe Balayn, Mireia Yurrita, Jie Yang, and Ujwal Gadiraju. "✓ Fairness Toolkits, A Checkbox Culture?" On the Factors that Fragment Developer Practices in Handling Algorithmic Harms. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '23, pp. 482–495, New York, NY, USA, August 2023. Association for Computing Machinery. ISBN 9798400702310. doi: 10.1145/3600211.3604674.
- Solon Barocas, Moritz Hardt, and Arvind Narayanan. Fairness and Machine Learning: Limitations and Opportunities. MIT Press, 2023.
- Rachel K. E. Bellamy, Kuntal Dey, Michael Hind, Samuel C. Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilovic, Seema Nagar, Karthikeyan Natesan Ramamurthy, John Richards, Diptikalyan Saha, Prasanna Sattigeri, Moninder Singh, Kush R. Varshney, and Yunfeng Zhang. AI Fairness 360: An Extensible Toolkit for Detecting, Understanding, and Mitigating Unwanted Algorithmic Bias, October 2018.
- Albert Berahas, Frank E. Curtis, Daniel Robinson, and Baoyu Zhou. Sequential quadratic optimization for nonlinear equality constrained stochastic optimization. *SIAM Journal on Optimization*, 31: 1352–1379, 05 2021. doi: 10.1137/20M1354556.
- Albert S. Berahas, Frank E. Curtis, Michael J. O'Neill, and Daniel P. Robinson. A stochastic sequential quadratic optimization algorithm for nonlinear equality constrained optimization with rank-deficient jacobians, 2023. URL https://arxiv.org/abs/2106.13015.
- D.P. Bertsekas and W. Rheinboldt. *Constrained Optimization and Lagrange Multiplier Methods*. Computer science and applied mathematics. Academic Press, 2014. ISBN 9781483260471.
- Sarah Bird, Miro Dudík, Richard Edgar, Brandon Horn, Roman Lutz, Vanessa Milan, Mehrnoosh Sameki, Hanna Wallach, and Kathleen Walker. Fairlearn: A toolkit for assessing and improving fairness in AI. Technical Report MSR-TR-2020-32, Microsoft, May 2020.
- José H. Blanchet, Peter W. Glynn, and Yanan Pei. Unbiased multilevel monte carlo: Stochastic optimization, steady-state simulation, quantiles, and other applications. *arXiv: Statistics Theory*, 2019. URL https://api.semanticscholar.org/CorpusID:127952798.
- Raghu Bollapragada, Cem Karamanli, Brendan Keith, Boyan Lazarov, Socratis Petrides, and Jingyi Wang. An adaptive sampling augmented lagrangian method for stochastic optimization with deterministic constraints. *Computers and Mathematics with Applications*, 149:239–258, 2023. ISSN 0898-1221. doi: https://doi.org/10.1016/j.camwa.2023.09.014. URL https://www.sciencedirect.com/science/article/pii/S0898122123003991.
- Maarten Buyl, Marybeth Defrance, and Tijl De Bie. fairret: a framework for differentiable fairness regularization terms. In *International Conference on Learning Representations*, 2024.
- Luiz FO Chamon, Santiago Paternain, Miguel Calvo-Fullana, and Alejandro Ribeiro. Constrained learning with non-convex losses. *IEEE Transactions on Information Theory*, 69(3):1739–1760, 2022.

```
Changan Chen, Frederick Tung, Naveen Vedula, and Greg Mori. Constraint-aware deep neural network compression. In Computer Vision – ECCV 2018: 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part VIII, pp. 409–424, Berlin, Heidelberg, 2018. Springer-Verlag. ISBN 978-3-030-01236-6. doi: 10.1007/978-3-030-01237-3_25. URL https://doi.org/10.1007/978-3-030-01237-3_25.
```

- Richard J. Chen, Judy J. Wang, Drew F. K. Williamson, Tiffany Y. Chen, Jana Lipkova, Ming Y. Lu, Sharifa Sahai, and Faisal Mahmood. Algorithmic fairness in artificial intelligence for medicine and healthcare. *Nature Biomedical Engineering*, 7(6):719–742, Jun 2023. ISSN 2157-846X. doi: 10. 1038/s41551-023-01056-8. URL https://doi.org/10.1038/s41551-023-01056-8.
- Zhenpeng Chen, Jie M. Zhang, Max Hort, Mark Harman, and Federica Sarro. Fairness testing: A comprehensive survey and analysis of trends. *ACM Trans. Softw. Eng. Methodol.*, 33(5), June 2024. ISSN 1049-331X. doi: 10.1145/3652155. URL https://doi.org/10.1145/3652155.
- Andrew Cotter, Heinrich Jiang, Serena Wang, Taman Narayan, Seungil You, Karthik Sridharan, and Maya R. Gupta. Optimization with non-differentiable constraints with applications to fairness, recall, churn, and other goals. *Journal of Machine Learning Research*, 20(172):1–59, 2019.
- Frank E. Curtis, Michael J. O'Neill, and Daniel P. Robinson. Worst-case complexity of an sqp method for nonlinear equality constrained stochastic optimization. *Mathematical Programming*, 205(1):431–483, May 2024a. ISSN 1436-4646. doi: 10.1007/s10107-023-01981-1. URL https://doi.org/10.1007/s10107-023-01981-1.
- Frank E. Curtis, Daniel P. Robinson, and Baoyu Zhou. Sequential quadratic optimization for stochastic optimization with deterministic nonlinear inequality and equality constraints. *SIAM Journal on Optimization*, 34(4):3592–3622, 2024b. doi: 10.1137/23M1556149. URL https://doi.org/10.1137/23M1556149.
- Damek Davis, Dmitriy Drusvyatskiy, Sham Kakade, and Jason D. Lee. Stochastic subgradient method converges on tame functions, 2018. URL https://arxiv.org/abs/1804.07795.
- MaryBeth Defrance, Maarten Buyl, and Tijl De Bie. Abcfair: an adaptable benchmark approach for comparing fairness methods, 2024. URL https://arxiv.org/abs/2409.16965.
- Eoin Delaney, Zihao Fu, Sandra Wachter, Brent Mittelstadt, and Chris Russell. OxonFair: A Flexible Toolkit for Algorithmic Fairness, November 2024.
- Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. Retiring adult: New datasets for fair machine learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- Mengnan Du, Subhabrata Mukherjee, Guanchu Wang, Ruixiang Tang, Ahmed Hassan Awadallah, and Xia Hu. Fairness via representation neutralization. In *Neurips*, 2021.
- Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, ITCS '12, pp. 214–226, New York, NY, USA, 2012. Association for Computing Machinery. ISBN 9781450311151. doi: 10.1145/2090236.2090255. URL https://doi.org/10.1145/2090236.2090255.
- Harrison Edwards and Amos Storkey. Censoring representations with an adversary, 2016. URL https://arxiv.org/abs/1511.05897.
- Alessandro Fabris, Stefano Messina, Gianmaria Silvello, and Gian Antonio Susto. Algorithmic fairness datasets: the story so far. *Data Mining and Knowledge Discovery*, 36(6):2074–2152, Nov 2022. ISSN 1573-756X. doi: 10.1007/s10618-022-00854-z. URL https://doi.org/10.1007/s10618-022-00854-z.
- Francisco Facchinei and Vyacheslav Kungurtsev. Stochastic approximation for expectation objective and expectation inequality-constrained nonconvex optimization, 2023. URL https://arxiv.org/abs/2307.02943.

Francisco Facchinei, Vyacheslav Kungurtsev, Lorenzo Lampariello, and Gesualdo Scutari. Ghost penalties in nonconvex constrained optimization: Diminishing stepsizes and iteration complexity. *Mathematics of Operations Research*, 46(2):595–627, 2021. doi: 10.1287/moor.2020.1079. URL https://doi.org/10.1287/moor.2020.1079.

- Yuchen Fang, Sen Na, Michael W. Mahoney, and Mladen Kolar. Fully stochastic trust-region sequential quadratic programming for equality-constrained optimization problems. *SIAM Journal on Optimization*, 34(2):2007–2037, 2024. doi: 10.1137/22M1537862. URL https://doi.org/10.1137/22M1537862.
- Jose Gallego-Posada, Juan Ramirez, Meraj Hashemizadeh, and Simon Lacoste-Julien. Cooper: A Library for Constrained Optimization in Deep Learning, April 2025.
- Xiaotian Han, Jianfeng Chi, Yu Chen, Qifan Wang, Han Zhao, Na Zou, and Xia Hu. FFB: A Fair Fairness Benchmark for In-Processing Group Fairness Methods. In *The Twelfth International Conference on Learning Representations*, October 2023.
- Moritz Hardt, Eric Price, and Nathan Srebro. Equality of opportunity in supervised learning. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, NIPS'16, pp. 3323–3331, Red Hook, NY, USA, 2016. Curran Associates Inc. ISBN 9781510838819.
- Ruichuan Huang, Jiawei Zhang, and Ahmet Alacaoglu. Stochastic smoothed primal-dual algorithms for nonconvex optimization with linear inequality constraints, 2025. URL https://arxiv.org/abs/2504.07607.
- Yankun Huang and Qihang Lin. Oracle complexity of single-loop switching subgradient methods for non-smooth weakly convex functional constrained optimization, 2023. URL https://arxiv.org/abs/2301.13314.
- Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. Fairness-Aware Classifier with Prejudice Remover Regularizer. In Peter A. Flach, Tijl De Bie, and Nello Cristianini (eds.), *Machine Learning and Knowledge Discovery in Databases*, pp. 35–50, Berlin, Heidelberg, 2012. Springer. ISBN 978-3-642-33486-3. doi: 10.1007/978-3-642-33486-3\_3.
- Michael P. Kim, Amirata Ghorbani, and James Zou. Multiaccuracy: Black-box post-processing for fairness in classification. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '19, pp. 247–254, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450363242. doi: 10.1145/3306618.3314287. URL https://doi.org/10.1145/3306618.3314287.
- Tai Le Quy, Arjun Roy, Vasileios Iosifidis, Wenbin Zhang, and Eirini Ntoutsi. A survey on datasets for fairness-aware machine learning. *WIREs Data Mining and Knowledge Discovery*, 12(3), 03 2022. ISSN 1942-4795. doi: 10.1002/widm.1452. URL https://doi.org/10.1002/widm.1452.
- Zhu Li, Adrián Pérez-Suay, Gustau Camps-Valls, and Dino Sejdinovic. Kernel dependence regularizers and Gaussian processes with applications to algorithmic fairness. *Pattern Recognition*, 132: 108922, December 2022. ISSN 0031-3203. doi: 10.1016/j.patcog.2022.108922.
- Gilles Louppe, Michael Kagan, and Kyle Cranmer. Learning to pivot with adversarial networks. *Advances in neural information processing systems*, 30, 2017.
- David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. Learning adversarially fair and transferable representations. In *International Conference on Machine Learning*, pp. 3384–3393. PMLR, 2018.
- Sen Na, Mihai Anitescu, and Mladen Kolar. An adaptive stochastic sequential quadratic programming with differentiable exact augmented lagrangians. *Mathematical Programming*, 199(1): 721–791, May 2023a. doi: 10.1007/s10107-022-01846-z. URL https://doi.org/10.1007/s10107-022-01846-z.
- Sen Na, Mihai Anitescu, and Mladen Kolar. Inequality constrained stochastic nonlinear optimization via active-set sequential quadratic programming, 2023b. URL https://arxiv.org/abs/2109.11502.

```
Yatin Nandwani, Abhishek Pathak, Mausam, and Parag Singla. A primal dual formulation for deep learning with constraints. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/cf708fc1decf0337aded484f8f4519ae-Paper.pdf.
```

- Figen Oztoprak, Richard Byrd, and Jorge Nocedal. Constrained optimization in the presence of noise. *SIAM Journal on Optimization*, 33(3):2118–2136, 2023. doi: 10.1137/21M1450999. URL https://doi.org/10.1137/21M1450999.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019. URL http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf.
- Patricia Pauli, Anne Koch, Julian Berberich, Paul Kohler, and Frank Allgöwer. Training robust neural networks using lipschitz bounds. *IEEE Control Systems Letters*, 6:121–126, 2021.
- Dana Pessach and Erez Shmueli. A review on fairness in machine learning. *ACM Comput. Surv.*, 55 (3), February 2022. ISSN 0360-0300. doi: 10.1145/3494672. URL https://doi.org/10.1145/3494672.
- Evaggelia Pitoura, Kostas Stefanidis, and Georgia Koutrika. Fairness in rankings and recommendations: an overview. *The VLDB Journal*, 31(3):431–458, May 2022. ISSN 0949-877X. doi: 10.1007/s00778-021-00697-y. URL https://doi.org/10.1007/s00778-021-00697-y.
- Juan Ramirez, Meraj Hashemizadeh, and Simon Lacoste-Julien. Position: Adopt Constraints Over Penalties in Deep Learning, July 2025.
- Sathya N. Ravi, Tuan Dinh, Vishnu Suresh Lokhande, and Vikas Singh. Explicitly imposing constraints in deep networks via conditional gradients gives improved generalization and faster convergence. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):4772–4779, Jul. 2019. doi: 10.1609/aaai.v33i01.33014772. URL https://ojs.aaai.org/index.php/AAAI/article/view/4404.
- Qiankun Shi, Xiao Wang, and Hao Wang. A momentum-based linearized augmented lagrangian method for nonconvex constrained stochastic optimization. *Optimization Online*, 2022. URL https://optimization-online.org/?p=19870.
- Amal Tawakuli and Thomas Engel. Make your data fair: A survey of data preprocessing techniques that address biases in data towards fair ai. *Journal of Engineering Research*, 2024. ISSN 2307-1877. doi: https://doi.org/10.1016/j.jer.2024.06.016. URL https://www.sciencedirect.com/science/article/pii/S2307187724001871.
- Sahil Verma and Julia Rubin. Fairness definitions explained. In *Proceedings of the International Workshop on Software Fairness*, FairWare '18, pp. 1–7, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 9781450357463. doi: 10.1145/3194770.3194776. URL https://doi.org/10.1145/3194770.3194776.
- Sandra Wachter, Brent Mittelstadt, and Chris Russell. Bias Preservation in Machine Learning: The Legality of Fairness Metrics Under EU Non-Discrimination Law, January 2021.
- Mingyang Wan, Daochen Zha, Ninghao Liu, and Na Zou. In-processing modeling techniques for machine learning fairness: A survey. *ACM Trans. Knowl. Discov. Data*, 17(3), March 2023. ISSN 1556-4681. doi: 10.1145/3551390. URL https://doi.org/10.1145/3551390.
- Quan Zhou, Jakub Mareček, and Robert Shorten. Subgroup fairness in two-sided markets. *Plos one*, 18(2):e0281443, 2023a.

Quan Zhou, Jakub Mareček, and Robert Shorten. Fairness in forecasting of observations of linear dynamical systems. *Journal of Artificial Intelligence Research*, 76:1247–1280, April 2023b. ISSN 1076-9757. doi: 10.1613/jair.1.14050. URL http://dx.doi.org/10.1613/jair.1.14050.

# A ALGORITHMS IN MORE DETAIL

In this section, we provide the pseudocodes of algorithms presented in Section 3 as Algorithms 1 to 3. Recall that we denote by  $X_k^J = \{X_{k,j}\}_{j=1}^J$  a mini-batch of size J with the j-th element

$$X_{k,j} = (\nabla f(x_k, \xi_{k,j}), c(x_k, \zeta_{k,j}), \nabla c(x_k, \zeta_{k,j})). \tag{13}$$

# B ADDITIONAL EXPERIMENT: ONE PROTECTED ATTRIBUTE WITH FIVE VALUES

**Numerical setup.** We use the same numerical setup as in section 4 (hardware and software). In this experiment, we run the algorithm for 90 seconds instead of 30.

**Dataset and problems.** We consider the dataset ACSIncome, over the state of Virginia, and choose this time Mariage as the protected attribute. This attribute takes fives values, as opposed to the binary attribute setup of section 4.

We consider three optimization problems as approaches to tackle the learning task. First, we consider the *constrained* learning problem as described in eq. (4), with m=5. Second, we consider the unconstrained, but *penalized*, problem

$$\min_{\theta \in \mathbb{R}^n} \frac{1}{N} \sum_{i=1}^N \ell(f_{\theta}(X_i), Y_i) + \mathcal{R}(\theta) + \lambda \sum_{i=1}^m \left| \ell^{s_i}(\theta) - \frac{1}{m} \sum_{j=1}^m \ell^{s_j}(\theta) \right|, \tag{14}$$

where  $\lambda=0.4$  is a penalization weight tuned on the validation set. Third, we consider the unconstrained and unpenalized problem, as described in eq. (2), for comparison.

**Algorithms and parameters.** We solve the constrained learning problem (4) with Stochastic Ghost, Switching Subgradient, and SSL-ALM. We solve the baseline penalized problem (14) and the basic unconstrained unpenalized problem (2) using SGD.

The hyperparameters for each algorithm were tuned on the validation set; we picked the values resulting in lowest loss and constraint satisfaction after 60 seconds. Our hyperparameter choices are listed below:

- SSL-ALM:  $\tau = \eta = 0.01, \beta = 0.5, \mu = 2, \rho = 1.$
- Stochastic Ghost:  $\beta=1.0, \gamma_0=0.005, \zeta=0.05, \rho=0.1, \tau=1, \lambda=0.5.$
- SSw:  $\eta^f = 0.05$  constant,  $\eta^c_k$  diminishing with  $\eta^c_0 = 0.25$ ,  $\eta^c_k = \frac{\eta^c_{k-1}}{\sqrt{k}}$ , k > 0; constraint tolerance  $\epsilon$  diminishing with  $\epsilon_0 = 0.01$ ,  $\epsilon_k = \frac{\epsilon_{k-1}}{\sqrt{k}}$ , k > 0.

**Optimization performance.** We report in Figure 5 the evolution of the mean and quartiles of the train and test values over 10 runs. SGD on the unconstrained and unpenalized problem (2) (first row) converges to a model such that the constraint are consistently above the constraint bound for three values of the protected attribute (Wid, Div, and Nev). SGD on the penalized problem (second row) manages to meet all constraints for the training set, but constraint Div on the test set is eventually violated. The three constrained methods minimize function values while keeping with the constraint bounds. The Stochastic Ghost does not reach convergence in 90 seconds, we report in fig. 6 its behavior over 180 seconds. With that time budget, we see that it manages to minimize the loss well while keeping within constraints. The performance of SGD on penalized problem and the three constrained algorithms is comparable, especially so for SSw: the objective and constraint value trajectories can hardly be distinguished. However, note that the penalized problem required consequent preliminary computations in order to tune the penalization parameter  $\lambda$ . We found that the performance of the estimator was quite sensitive to the value of  $\lambda$ . In contrast, the constrained formulation does not feature a hyperparameter, and thus does not requires tuning. The algorithms for constrained minimization do depend on hyperparameters, which control their convergence speed. Nevertheless, they converged to feasible solutions for every (reasonable) hyperparameter value we tried. This observation is consistent with the argument of Ramirez et al. (2025).

# Algorithm 1 Stochastic Ghost algorithm

810 811

812

827 828 829

830 831

832

833

834

835

836

837

838

839

840

841

843 844 845

846

847

848

849

850

851

852

853

854

855

856

858

859

861

863

```
813
               Require: Training dataset \mathcal{D}, constraint dataset \mathcal{C}, initial neural network weights x_0
814
               Require: Parameters p_0 \in (0,1), \alpha_0, \hat{\alpha}, \rho, \tau, \beta
                1: for Iteration k = 0 to K - 1 do
815
                           Sample \xi \overset{iid}{\sim} \mathcal{P}_{\xi} and \zeta \overset{iid}{\sim} \mathcal{P}_{\zeta}
Sample N \sim \mathcal{G}(p_0)
Set J = 2^{N+1}
816
                3:
817
                4:
818
                           Sample a mini-batch \{\zeta_j\}_{j=1}^J so that \zeta_1, \dots, \zeta_J \stackrel{iid}{\sim} \mathcal{P}_{\zeta}
                5:
819
820
                           Sample a mini-batch \{\xi_j\}_{j=1}^J so that \xi_1, \dots, \xi_J \stackrel{iid}{\sim} \mathcal{P}_{\xi}
                6:
821
                           Set X_k^1 and X_k^{2^{N+1}} using (13)
                7:
822
                           Compute d(x_k) from (9)
                8:
823
                9:
                           Set \alpha_k = \alpha_{k-1}(1 - \hat{\alpha}\alpha_{k-1})
824
               10:
                           Update x_{k+1} = x_k + \alpha_k d(x_k)
825
               11: end for
826
```

# Algorithm 2 Stochastic Smoothed and Linearized AL Method for solving (1)

```
Require: Training dataset \mathcal{D}, constraint dataset \mathcal{C}, initial neural network weights x_0
Require: Parameters \mu, \eta, M_y > 0, \tau, \beta, \rho \ge 0
 1: for Iteration k = 0 to K - 1 do
           Sample \xi \stackrel{iid}{\sim} \mathcal{P}_{\xi} and \zeta_1, \zeta_2 \stackrel{iid}{\sim} \mathcal{P}_{\zeta}
 2:
 3:
           y_{k+1} = y_k + \eta c(x, \zeta_1)
 4:
           if ||y_{k+1}|| \geq M_y then
 5:
                 y_{k+1} = 0
 6:
           x_{k+1} = \operatorname{proj}_{\mathcal{X}}(x_k - \tau G(x_k, y_{k+1}, z_k; \xi, \zeta_1, \zeta_2)), where G is defined in (10)
 7:
 8:
            z_{k+1} = z_k + \beta(x_k - z_k)
 9: end for
```

# Algorithm 3 Stochastic Switching Subgradient Method

```
Require: Training dataset \mathcal{D}, constraint dataset \mathcal{C}, initial neural network weights x_0 \in \mathcal{X}
Require: Total number of iterations K, sequence of tolerances of infeasibility \epsilon_k \geq 0, sequences of
      stepsizes \eta_k^I and \eta_k^c, mini-batch size J, starting index k_0 for recording outputs, I = \emptyset
 1: for Iteration k = 0 to K - 1 do
           Sample a mini-batch \{\zeta_j\}_{j=1}^J so that \zeta_1, \dots, \zeta_J \stackrel{iid}{\sim} \mathcal{P}_{\zeta}
 2:
           Set \overline{c}^J(x_k) = \frac{1}{J} \sum_{j=1}^J c(x_k, \zeta_j)
 3:
           if \overline{c}^J(x_k) < \epsilon_k then
 4:
                 Sample \xi \stackrel{iid}{\sim} \mathcal{P}_{\xi} and generate S^f(x_k, \xi)
 5:
                 Set x_{k+1} = \operatorname{proj}_{\mathcal{X}}(x_k - \eta_k^f S^f(x_k, \xi)) and, if k \ge k_0, I = I \cup \{k\}
 6:
 7:
                 Sample \zeta \stackrel{iid}{\sim} \mathcal{P}_{\zeta} and generate S^c(x_k, \zeta)
 8:
                 Set x_{k+1} = \operatorname{proj}_{\mathcal{X}}(x_k - \eta_k^c S^c(x_k, \zeta)) and, if k \ge k_0, I = I \cup \{k\}
 9:
           end if
10:
12: Output: x_{\tau} with \tau randomly sampled from I using P(\tau = k) = \frac{\eta_k}{\sum_{s \in I} \eta_s}.
```

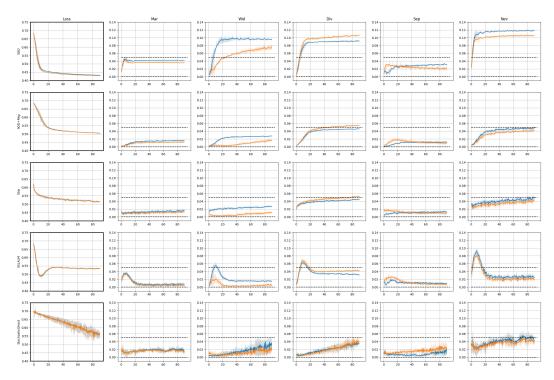


Figure 5: Train (blue) and test (orange) statistics (mean and quartiles) over time (s) on the ACS Income dataset for each algorithm: Switching Subgradient (row 1), SSL-ALM (row 2), SGD (row 3), and Stochastic Ghost (row 4). The plots depict the mean values for loss (leftmost column) and constraints (second to rightmost column) at each timestamp, rounded to the nearest 0.5 seconds, over 10 runs. The shaded area depicts the region between the first and third quartiles.

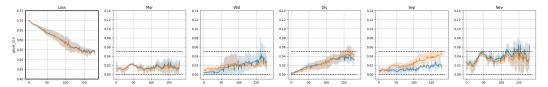


Figure 6: A longer run of the Stochastic Ghost algorithm in the setting of fig. 5: 180 seconds, and rounding each second.