

# Promoting Fairness in Learned Models by Learning to Active Learn under Parity Constraints

Amr Sharaf  
amrsharaf@microsoft.com  
Microsoft  
Redmond, Washington, USA

Hal Daumé III  
University of Maryland, Microsoft  
Research  
College Park, USA  
me@ha13.name

Renkun Ni  
University of Maryland  
College Park, USA

## ABSTRACT

Machine learning models can have consequential effects when used to automate decisions, and disparities between groups of people in the error rates of those decisions can lead to harms suffered more by some groups than others. Past algorithmic approaches aim to enforce parity across groups given a fixed set of training data; instead, we ask: what if we can gather more data to mitigate disparities? We develop a meta-learning algorithm for parity-constrained active learning that learns a policy to decide which labels to query so as to maximize accuracy subject to parity constraints. To optimize the active learning policy, our proposed algorithm formulates the parity-constrained active learning task as a bi-level optimization problem. The inner level corresponds to training a classifier on a subset of labeled examples. The outer level corresponds to updating the selection policy choosing this subset to achieve a desired fairness and accuracy behavior on the trained classifier. To solve this constrained bi-level optimization problem, we employ the Forward-Backward Splitting optimization method. Empirically, across several parity metrics and classification tasks, our approach outperforms alternatives by a large margin.

## KEYWORDS

active learning, meta-learning

### ACM Reference Format:

Amr Sharaf, Hal Daumé III, and Renkun Ni. 2022. Promoting Fairness in Learned Models by Learning to Active Learn under Parity Constraints. In *2022 ACM Conference on Fairness, Accountability, and Transparency (FAccT '22)*, June 21–24, 2022, Seoul, Republic of Korea. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/3531146.3534632>

## 1 INTRODUCTION

Machine learning models often lead to harms due to disparities in behavior across social groups: for example, an automated hiring system may be more likely to recommend hiring people of privileged races, genders, or age groups [23, 48]. These disparities are typically caused by biases in historic data (society is biased);

a substantial literature exists around “de-biasing” methods for algorithms, predictions, or models. Such approaches always assume that the training data is fixed, leading to a false choice between efficacy (e.g., accuracy, AUC) and “fairness” (typically measured by a metric of parity across subgroups [9, 30]). This is in stark contrast to how machine learning *practitioners* address disparities in model performance: they *collect more data* that’s more relevant or representative of the populations of interest [27, 46]. This disconnect leads to a mismatch between sources of bias, and the algorithmic interventions developed to mitigate them [52].

We consider a different trade-off: given a pre-existing dataset, which may have been collected in a highly biased manner, how can we optimize an efficacy *vs annotation cost* trade-off under a target parity constraint? We call this problem *parity-constrained active learning*, where a maximal allowed disparity (according to any of a number of different measures, see Table 1) is enforced during a data collection process. We specifically consider the case where some “starting” dataset has already been collected, distinguishing our procedure from more standard active learning settings in which we typically start from no data ([41], see §2). The goal then is to collect as little data as is needed to keep accuracy high while maintaining a constraint on parity (as a measure of fairness). As an example, consider disparities in pedestrian detection by skin tone [51]: A pedestrian detector is trained based on a dataset of 100k images, but an analysis shows that it performs significantly better at detecting people with light skin than people with dark skin. Our goal is to label few *additional* samples while achieving a high accuracy under a constraint on the disparity between these groups.

We propose to solve the parity-constrained active learning problem using a meta-learning approach, very much in the style of recent work on meta-learning for active learning [6, 18, 33]. Our algorithm, PARITY-CONSTRAINED META ACTIVE LEARNING (PANDA; see §3), uses data to learn a selection policy that picks which examples to have labeled. The data on which it learns this selection policy is the pre-existing (possibly biased!) dataset from which it continues learning.

To learn this selection policy, PANDA simulates many parity-constrained active learning tasks on this pre-existing dataset, to learn the selection policy. PANDA formulates the parity-constrained active learning task as a bi-level optimization problem. The inner level corresponds to training a classifier on a subset of labeled examples. The outer level corresponds to updating the selection policy choosing this subset to achieve a desired fairness and accuracy behavior on the trained classifier. To solve this constrained bi-level optimization problem, PANDA employs the *Forward-Backward Splitting* (FBS, [12, 24, 35]) optimization method (also known as the

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

FAccT '22, June 21–24, 2022, Seoul, Republic of Korea

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9352-2/22/06...\$15.00

<https://doi.org/10.1145/3531146.3534632>

proximal gradient method). Despite its apparent simplicity, FBS can handle non-differentiable objectives with possibly non-convex constraints while maintaining the simplicity of gradient-descent methods.

Through exhaustive empirical experiments (§4), we show:

- (1) PANDA is effective: it outperforms alternative active learning algorithms by a large margin under the same setting while enforcing the desired behavior on fairness.
- (2) PANDA is general-purpose: it learns the selection policy end-to-end and can handle a wide set of non-differentiable and non-convex constraints on fairness parity using Gumbel-Softmax reparameterization [25, 28, 36] and differentiable approximations.
- (3) PANDA is powerful: it employs a Transformer model architecture [45] to represent the learned selection policy. This architecture has achieved state-of-the-art performance in many tasks including language modeling [14], machine translation [15], and unsupervised pre-training [16].

## 2 BACKGROUND AND RELATED WORK

Concerns about biased or disparate treatment of groups or individuals by computer systems has been raised since the 1990s [21]. Machine learning and other statistical techniques provide ample opportunity for pre-existing societal bias to be incorporated into computer systems through data, leading to a burgeoning of research studying disparities in machine learning [1, 13, i.a.]. Arguably, because society is biased, societal data will be biased, and therefore, if unchecked, any machine learning model trained on such data will inherit its biases.

Much technical machine learning research has gone into defining what disparate treatment means formally, leading to a zoo of parity metrics [37] (see Table 1 for examples), proofs of their incompatibilities [10, 32], and analyses of how they conform to normative notions of fairness [44]. This has led to machine learning algorithms that optimize not just for accuracy, but rather for accuracy subject to a constraint on parity between groups [2].

A parallel line of research has considered the human side of analyzing disparities in machine learning systems, including visualization [8], debugging [9], and needs-finding [27, 46]. One finding of the latter is that machine learning practitioners and data scientists often have control over training data, which is not taken into account in most technical machine learning research. For instance, [27]’s results show that 78% of practitioners who had attempted to address disparities did so by trying to collect more data, despite the lack of tools that support this.

Curating more data is not a foreign concept in the machine learning research: active learning—the learning paradigm in which the learner itself chooses which examples to have labeled next—has been studied extensively over the past five decades (e.g., Angluin [4], Cohn et al. [11], Fedorov [19], Jiang and Ip [29], Settles [41]). Most active learning approaches select samples to label based on some notion of uncertainty (e.g., entropy of predictions, or margin). Most relevant to our work are recent active learning approaches based on meta-learning [6, 18]: here, instead of designing the selection strategy by hand, the selection strategy is learned based on offline, simulated active learning problems. So long as those offline

problems are sufficiently similar to the target, real, active learning problem, there is hope that the learned strategy will generalize well.

We are aware of only one paper that considers active learning in the context of fairness: Fair Active Learning (FAL) by [3]. FAL uses a handselection strategy that interpolates between an uncertainty-based selection criteria, and a “fairness” criteria that estimates the impact on disparity if the label of a given point were queried (by computing expected disparity over all possible labels). FAL selects data points to be labeled to balance a convex combination of model accuracy and parity, with the trade-off specified by a hyperparameter. Empirically, [3] showed a significant reduction in disparity while maintaining accuracy. Our setting is slightly different than FAL—we assume pre-existing data—but we compare extensively to this approach experimentally under similar conditions (§4).

## 3 PROBLEM DEFINITION AND PROPOSED APPROACH

In this section we define *parity-constrained active learning* and describe our algorithm, PANDA.

### 3.1 Problem Definition: Parity-Constrained Active Learning

We consider the following model. We have collected a dataset of  $N$  labeled examples,  $D^0 = (\mathbf{x}_n, y_n)_{n=1}^N$  over an input space  $\mathcal{X}$  (e.g., images) and output space  $\mathcal{Y}$  (e.g., pedestrian bounding boxes), and have access to a collection of  $M$ -many *unlabeled* examples,  $U = (\mathbf{x}_m)_{m=1}^M$ . Each input example  $x \in \mathcal{X}$  is associated with a unique group  $g(x)$  (e.g., skin tone). We fix a hypothesis class  $\mathcal{H} \subset \mathcal{Y}^{\mathcal{X}}$  and learning algorithm  $\mathcal{A}$  that maps a labeled sample  $D$  to a classifier  $h \in \mathcal{H}$ . Finally, we have a loss function  $\ell(y, \hat{y}) \in \mathbb{R}^{\geq 0}$  (e.g., squared error, classification error, etc.) and a target *disparity metric*,  $\Delta(h) \in \mathbb{R}^{\geq 0}$  (such as those in Table 1). The goal is to label as few examples from  $U$  as possible to learn a classifier  $h$  on the union of  $D^0$  and the newly labeled examples with high accuracy subject to a constraint that  $\Delta(h) < \rho$  for a given threshold  $\rho > 0$ . We assume access to a (small) validation set  $V$  of labeled examples (which can be taken to be a subset of  $D$ ). We will denote population expectations and disparities by  $\mathbb{E}$  and  $\Delta$ , respectively, and their estimates on a finite sample by  $\hat{\mathbb{E}}_A$  and  $\hat{\Delta}_A$ , where  $A$  is the sample.

The specific interaction model we assume is akin to standard active learning with labeling budget  $B$ :

- 1: train the initial classifier on the pre-existing dataset:  $h^0 = \mathcal{A}(D^0)$ .
- 2: **for** round  $b = 1 \dots B$  **do**
- 3:   generate categorical probability distribution  $Q = \pi(h^{b-1}, U)$  over  $U$  using policy  $\pi$ .
- 4:   sample an unlabeled example  $\mathbf{x} \sim Q$ , query its label  $y$ , and set  $D^b = D^{b-1} \cup \{(\mathbf{x}, y)\}$ .
- 5:   train/update classifier:  $h^b = \mathcal{A}(D^b)$ .
- 6: **end for**
- 7: **return** classifier  $h^B$ , it’s validation loss  $\hat{\mathbb{E}}_V \ell(y, h^B(\mathbf{x}))$  and validation disparity  $\hat{\Delta}_V(h^B)$ .

METRIC	DESCRIPTION & MATHEMATICAL DEFINITION
Demographic Parity	Predictions $h(x)$ are statistically independent of the group $g(x)$ [20]: $\Delta^{\text{DP}}(h) \triangleq \max_a   \mathbb{E}[h(x)   g(x)=a] - \mathbb{E}[h(x)]  $
Equalized Odds	Predictions $h(x)$ are independent of the group $g(x)$ given the true label $y$ [26]: $\Delta^{\text{EO}}(h) \triangleq \max_{a,y}   \mathbb{E}[h(x)   g(x)=a, Y=y] - \mathbb{E}[h(x)   Y=y]  $
Error-rate Balance	False positive and false negative error rates are equal across groups [10]: $\Delta^{\text{EB}}(h) \triangleq \max_{a,a',y}   \mathbb{E}[h(x)   g(x)=a, Y=y] - \mathbb{E}[h(x)   g(x)=a', Y=y]  $

**Table 1: Three common measures of disparity for binary classification (extensions to multiclass are generally straightforward), expressed in terms of differences in expected values of predictions (where we take  $h : \mathcal{X} \rightarrow \{0, 1\}$ ). We denote by  $g(x)$  the group to which the example  $x$  belongs. In some work, disparities are taken to be *ratios* of expectations, rather than differences.**

Under this model, the active learning strategy is summarized in the example selection policy  $\pi$ . For example, margin-based active learning [40] can be realized by setting  $\pi(h, U)$  to assign a  $Q(x) = 1[x = x^*]$  where  $x^* = \operatorname{argmin}_{x \in U} |h(x)|$ , where  $h$  returns the margin. The goal in parity constraint actively learning is to construct a  $\pi$  with minimal expected loss subject to the constraint that  $\Delta(h) < \rho$ .

### 3.2 PANDA: Learning to Actively Learn under Parity Constraints

We develop a meta-learning approach, PANDA, to address the parity-constrained active learning problem: the selection policy  $\pi$  is trained to choose samples that, if labeled, are likely to optimize accuracy subject to a parity constraint. This learning happens on  $D$  itself: by simulating many possible ways additional data could be selected on the historic data, PANDA learns how to select additional examples, even if  $D$  itself was sampled in a biased manner.

To learn  $\pi$ , we construct a distribution of meta-training tasks,  $\mathfrak{M}$ ; samples  $(L, V) \sim \mathfrak{M}$  consist of a labeled dataset  $L$  (to simulate unlabeled data  $U$ ) and a validation set  $V$ . We form  $\mathfrak{M}$  by repeatedly reshuffling  $D$ , and produce a finite sample of meta-training tasks  $\mathfrak{D}$  i.i.d. from  $\mathfrak{M}$ . The meta-learning problem is then to optimize  $\pi$  on  $\mathfrak{D}$  to achieve high accuracy subject to a constraint on disparity. We begin by first writing the parity-constrained problem according to its characteristic function:

$$\hat{h} \in \operatorname{argmin}_{h \in \mathcal{H}} : \hat{\Delta}_V(h) < \rho \iff \hat{\mathbb{E}}_V \ell(\mathbf{y}, h(\mathbf{x})) \iff (1)$$

$$\hat{h} \in \operatorname{argmin}_{h \in \mathcal{H}} \hat{\mathbb{E}}_V \ell(\mathbf{y}, h(\mathbf{x})) + \chi_{\hat{\Delta}, \rho, V}(h) \quad (2)$$

where  $\chi_{\hat{\Delta}, \rho, V}(h) = 0$  if  $\hat{\Delta}_V(h) < \rho$  and  $+\infty$  otherwise; for brevity we write  $\chi(h)$  when  $\hat{\Delta}, \rho, V$  is clear from context. Under reasonable assumptions, both minimizers are finite.

Given this, the meta-learning optimization problem is:

$$\min_{\pi \in \Pi} \hat{\mathbb{E}}_{(L, V) \sim \mathfrak{D}} [\hat{\mathbb{E}}_V \ell(\mathbf{y}, h_L^\pi(\mathbf{x})) + \chi(h_L^\pi)] \quad (3)$$

$$\text{where } h_L^\pi = \text{ACTIVELEARNSIM}(\mathcal{A}, D, L, \pi) \quad (4)$$

Here,  $\text{ACTIVELEARNSIM}(\mathcal{A}, D, L, \pi)$  is the interactive algorithm in §3.1, where  $U$  is taken to be  $L$  (with labels hidden) and when a label is queried, it is retrieved from  $L$ .

When  $\mathcal{A}$  is, itself, an optimizer—as it is in most machine learning settings—then formulation Eq 4 is a constrained bilevel optimization problem. The outer optimization is over the sampling policy

$\pi$ , and the inner optimization is over the optimization over  $h$  in  $\text{ACTIVELEARNSIM}$ . We assume that  $\mathcal{A}$  can be written as a computational graph, in which case the outer objective can be optimized by unrolling the computational graph of  $\mathcal{A}$ . This introduces second-order gradient terms, but remains computationally feasible so long as the unrolled graph of  $\mathcal{A}$  is not too long: we ensure this by only running a few steps of SGD inside  $\mathcal{A}$  and using a simple hypothesis class for  $\mathcal{H}$ . In many cases,  $\mathcal{A}$  optimizes a model through stochastic gradient descent (SGD). The unrolled graph corresponds to the computations made by this SGD procedure. To ensure that the unrolled graph is not too long, we limit the number of steps of SGD that  $\mathcal{A}$  runs, and use a relatively simple hypothesis class for  $\mathcal{H}$ .

There remain two challenges to solve Eq 4. The first is to address the discontinuous nature of the characteristic function  $\chi$ ; we use forward-backward splitting to address this. The second is that the unrolling of  $\text{ACTIVELEARNSIM}$  yields a computational graph that has stochasticity (due to the sampling of unlabeled examples); we use the Gumbel reparameterization trick to address this.

**Forward-Backward Splitting** (FBS) is a class of optimization methods [35], which provide the simplicity of gradient descent methods while being able to enforce possibly non-differentiable constraints. In FBS, the objective is separated into a differentiable part  $f(x)$  and an arbitrary (not even necessarily smooth) part  $g(x)$ . The algorithm operates iteratively by first taking a gradient step just with respect to  $f$  to an intermediate value:  $x' = x - \eta \nabla f(x)$ . Next, it computes a proximal step that chooses the next iterate  $x$  to minimize  $\eta g(x) + \|x - x'\|^2/2$ . When  $f$  is convex, FBS converges rapidly; for non-convex problems (like Eq 4), theoretical convergence rates are unknown, but the algorithm works well in practice [24].

To apply FBS to our problem, we choose our policy class  $\Pi$  to be a differentiable neural network (see §3.3). We set  $f$  to be the expected loss term in Eq 4, and  $g$  to be the characteristic function on the disparity. The gradient step with respect to  $f$  can be computed by automatic differentiation of the unrolled computational graph. The proximal step requires projecting onto  $\chi$ ; for complex  $\Pi$  there is no closed-form solution; instead, we run a separate approximate projection step by running several steps of gradient descent on a smoothed version of  $\chi$ , which takes values 0 when the constraint is satisfied, and takes value  $\hat{\Delta}_V(h)$  otherwise. This remains non-continuous, but (sub)differentiable—empirically, this approximate projection always finds an iterate that satisfies the original constraint.

**Gumbel Reparameterization** is a generic technique to avoid back-propagating through stochastic sampling nodes in the computational graph [25, 28, 34, 36]. This trick allows us to sample from the categorical distribution  $Q$  by independently perturbing the log-probabilities  $Q_i$  with Gumbel noise and finding the largest element, thus enabling end-to-end differentiation through ACTIVELEARNSIM, so long as  $\mathcal{A}$  is differentiable.

The **Full Training Algorithm** for PANDA is summarized in **Algorithm 1**. Following the Forward-Backward Splitting template, PANDA operates in an iterative fashion. Over iterations, PANDA simulates a parity-constrained active learning setting for the current model parameters  $\theta^k$ . **Line 4** performs a simple forward gradient descent step to maximize the classifier performance. This step begins at iterate  $\theta^k$ , and then moves in the direction of the (negative) gradient of the performance loss, which is the direction of steepest descent. **Line 5** is the proximal (or backward) step, which enforces the parity constraint; this works even when the parity metric is non-differentiable. In both the gradient descent step and the proximal step, PANDA performs bilevel optimization. For example, the gradient step is a gradient with respect to the parameters of the selection policy, of the computational graph defined by ACTIVELEARNSIM. That function, itself, performs an optimization of the classifier  $h$ .

### 3.3 Network Structure of Selection Policy

The selection policy  $\pi$  takes as input the current classifier  $h$  and unlabeled dataset  $U$ , and produces a distribution  $Q$  over elements of  $U$ . We explore policies that are *agnostic* to changes in  $h$ , meaning that  $Q$  at round  $b$  is identical for all  $b$ . This introduces a limitation that  $\pi$  cannot directly adapt to changes in  $h$ ; however, since  $\pi$  is optimized end-to-end, we empirically found this to be a minor limitation. A significant advantage of this choice is that it means that ACTIVELEARNSIM can be unrolled much more easily: the simple Gumbel softmax can be replaced with Gumbel-top- $B$  [34, 47] and unrolled in a single step, rather than a sequence of  $B$ -many steps.

Because  $\pi$  must effectively make all selections in a single step, it is important that  $\pi$  consider each  $\mathbf{x}$  in the context of all other items in  $U$ , and not consider each  $\mathbf{x}$  individually. We implement this using a Transformer architecture [45], in which a self-attention mechanism essentially combines information across different  $\mathbf{x}$ s in  $U$ . Specifically, we treat the examples in  $U$  as an unordered sequence as input to the Transformer encoder<sup>1</sup>. The Transformer architecture employs several layers of self-attention across  $U$  with independent feed-forward networks for each position. The final layer of the Transformer can be interpreted as a contextual representation for each  $\mathbf{x} \in U$ , where the context is “the rest of  $U$ .” We use a final linear softmax layer to map these contextual representations to the probability distribution  $Q$ .

Because this model architecture is flexible, it is also data-hungry, and training all of its parameters based just on a small set of  $B$  examples is unlikely to be sufficient. This is where the initial dataset  $D^0$  comes in: we pretrain the parameters of the Transformer on  $D^0$

<sup>1</sup>Recall that although Transformers are typically used for *ordered* problems like language modeling [14] and translation [15], this is not how they “naturally” work: ordered inputs to Transformers require additional “position” tags, which we omit.

and use the  $B$ -many actively selected samples to fine-tune the final layer, thus keeping the required sample complexity small.

## 4 EXPERIMENTS

We conduct experiments in the standard active learning manner: pretend that a labeled dataset is actually unlabeled, and use its labels to answer queries. Experimentally, given a complete dataset, we first split it 50/50 into meta-training and meta-testing sections. We use meta-training to pretrain the Transformer model (see §3.3), and also to train PANDA. All algorithms use the same Transformer representation. The meta-testing section is split again 50/50 into the “unlabeled” set and the test set.

Picking a good dataset for parity-constrained active learning is challenging: it needs to contain a protected attribute, be sufficiently large that an active sample from unlabeled portion is representative (i.e., as the size of the sample approaches the size of the unlabeled data, all algorithms will appear to perform identically), and be sufficiently rich that learning does not happen “too quickly.”

We considered three standard datasets: COMPAS [5], Adult Income [17], and Law School [50]. Law School has only two features and we found only a few examples are needed to learn; and COMPAS we found to be too small, and is also problematic to use in conditions separated from its true use case [7]. COMPAS consists of just under  $8k$  samples, so after two splits, each set contains only  $2k$  samples. We anticipated that this would be too small for three reasons. First, with a budget  $B = 400$ , many algorithms will end up sampling very similar sets, resulting in difficulties telling approaches apart. Second, we found that after pre-training, 15–20 completely random samples suffice to learn a classifier that is as good as one trained on all the remaining data. Nonetheless, we performed experiments on COMPAS for all baselines and found that while PANDA can fit the meta-training data well, and this generalizes well with respect to *loss*, it has poor generalization with respect to *disparity*. We also ran Fairlearn (described below) on this dataset randomly sampled subsets of the training data, and found that, while it eventually is able to achieve a target disparity level of 0.04 once  $B = 400$ , at any point with  $B < 300$  the test-time disparity is significantly larger. We therefore drop COMPAS from consideration; it seems ill-suited to a warm-start active learning paradigm. This left only the Adult Income dataset for experiments. This dataset consists of 48,842 examples and 251 features (with one-hot encodings of categorical variables) and the binary prediction task is whether someone makes more than 50k per year, with binary gender as the group attribute (the dataset does not contain information about gender beyond male/female).

### 4.1 Baseline Active Learning Approaches

Our experiments aim to determine how PANDA compares to alternative active learning strategies, including those that explicitly take disparity into account as well as those that do not. Among those that do not consider disparity, we compare to:

**Random Sampling** – select examples to label randomly.

**Margin Sampling** – uncertainty-based active learning that selects the example closest to the current decision boundary [40].

**Algorithm 1** Parity-constrained Active Learning via PANDA

**Input:** pre-existing datasets  $D$ , budget  $B$ , loss function  $\ell$ , disparity metric  $\Delta$ , threshold  $\rho$ , meta-learning learning rate schedule  $\langle \eta^k \rangle_{k \geq 1}$ , and inner learning rate  $\eta'$

**Output:** Selection policy  $\pi$

```

1: Initialize selection policy  $\pi(\cdot; \theta^0)$  parameterized by  $\theta^0$ 
2: for iteration  $k = 1 \dots$  convergence do
3:   Split  $D$  into pool  $L$  and validation set  $V$ 
4:    $\hat{\theta}^{k+1} = \theta^k - \eta^k \nabla_{\theta} \hat{\mathbb{E}}_V \ell(y, \text{ACTIVELEARNSIM}(\mathcal{A}, D, L, \pi(\cdot; \theta^k)(\mathbf{x})))$ 
5:    $\theta^{k+1} = \operatorname{argmin}_{\theta} \eta^k \chi_{\Delta, \rho, V}(\text{ACTIVELEARNSIM}(\mathcal{A}, D, L, \pi(\cdot; \theta))) + 1/2 \|\theta - \hat{\theta}^{k+1}\|^2$ 
6: end for
7: return  $\pi(\cdot; \theta^{\text{final}})$ 

8: function ACTIVELEARNSIM( $\mathcal{A}, D, L, \pi$ )
9:   Let  $\langle \mathbf{x}_i, y_i \rangle_{i=1}^{|L|}$  be an indexing of  $L$ 
10:  for  $b = 1 \dots B$  do
11:    set  $\tilde{Q}_i = \pi(h^{b-1}, \mathbf{x}_i) + \text{GUMBEL}(0)$  for all  $i$  and pick  $j = \operatorname{argmax}_i \tilde{Q}_i$ 
12:    take (a/several) gradient step(s) of the form:  $h^b = h^{b-1} - \eta' \nabla_h \ell(y_j, h(\mathbf{x}_j))$ 
13:  end for
14:  return  $h^B$ 
15: end function

```

**Entropy Sampling** – uncertainty-based active learning that selects the example with highest entropy of predicted label [41, 42].

We also consider alternate approaches that take groups and/or disparity into account explicitly.

**Group Aware Random Sampling** – select examples to label uniformly at random, restricted to the group on which worse performance is achieved by  $h^0$ . Closely related to active learning in domain adaptation [39, 43, 49].

**Fair Active Learning** – the “fair” active learning approach described in §2 that optimizes an interpolation between Entropy Sampling and expected disparity [3].

**Fairlearn** – select examples to label uniformly at random, and then run fairlearn to train a classifier to optimize accuracy subject to a parity constraint [2].

**Fairlearn+AL** – combine uncertainty-based active learning with a fairlearn, select examples closest to the decision boundary to label, and then run fairlearn to train a classifier to optimize accuracy subject to a parity constraint.

## 4.2 Implementation Details and Hyperparameter Tuning

We use the Transformer Model [45] implemented in PyTorch [38]. We use the standard transformer encoder with successive encoder layers. Each layer contains a self-attention layer, followed by a fully connected feed-forward layer. We use the feed-forward layer for decoding, where we sample  $B$  items from the predicted probability distribution in a single decoding step. To ensure a fair-comparison among all approaches, we use the same Transformer architecture as a feature extractor for all approaches. This ensures that PANDA has no additional advantage by observing more training data.

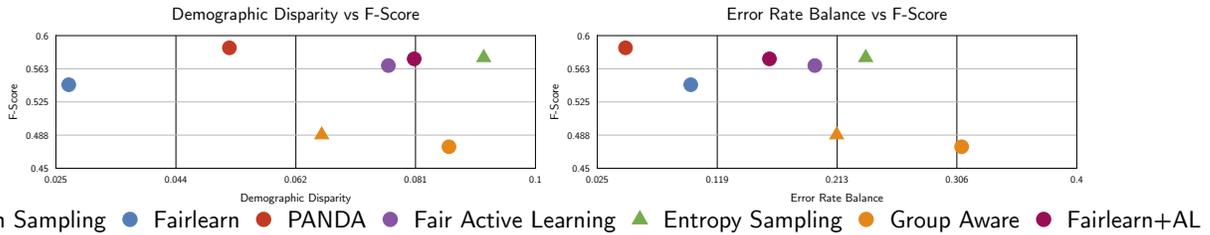
The model is optimized with Adam [31]. We optimize all hyperparameters with the Bayes search algorithm implemented in comet.ml using an adaptive Parzen-Rosenblatt estimator. We search for the best parameters for learning rate ( $10^{-2}$  to  $10^{-7}$ ), number of layers in the transformer encoder (1, 3, 5), embedding dimensions for the

encoder hidden-layer (16, 32, 64), as well as the initial value for the Gumbel-Softmax temperature parameter (1 to  $10^{-6}$ ) which is then learned adaptively as meta-training progresses. The sampled examples are used to train a linear classifier, again we optimize the hyper-parameters for the learning rate and batch size using Bayes search. For active learning model selection, we sweep over parameters using the random sampling active learning method. We found that hyper-parameters for random sampling work well for other alternative approaches too. We scale all the features to have a mean zero and unit standard deviation.

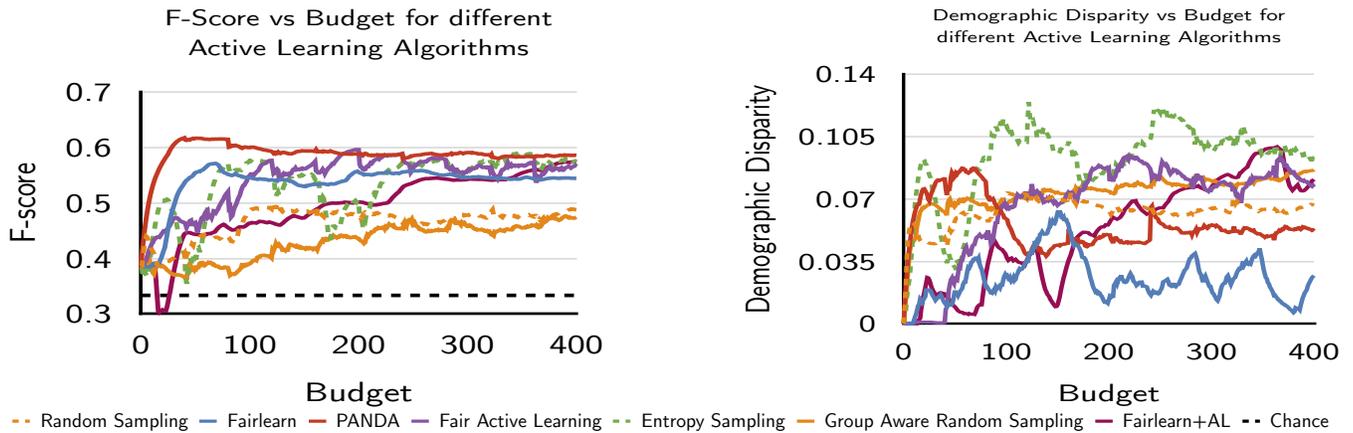
## 4.3 Evaluation Metrics and Results

We evaluate the performance of the learned classifiers using the overall F-score on the evaluation set  $V$ , and report violations for parity-constraints in terms of demographic disparity and error rate balance (Table 1), as these account for different ends of the constrained spectrum of parity metrics. In order to set trade-off parameters (the convex combination  $\alpha$  for FAL and the constraints for fairlearn and PANDA), we first run FAL with several different trade-off parameters to find a value for  $\alpha$  large enough that disparity matters but small enough that a non-zero F-score is possible. All results are with  $\alpha = 0.6$ . We then observed the final disparity for FAL of 0.08 and set a constraint for PANDA and FAL of half of that: 0.04. This choice was made to ensure that FAL has an overall advantage.

The main results are shown in Figure 1, where we consider performance for a fixed budget. Here, we first observe (unsurprisingly) that the baselines that do not take parity into account (Random Sampling and Entropy Sampling) do quite poorly (we do not plot margin-based sampling as it was dominated by Entropy sampling in all experiments). For example, while entropy sampling gets a very high F-score, it has quite poor disparity. Somewhat surprisingly, group-aware random sampling does worse in terms of disparity than even plain random sampling. FAL is able to achieve higher accuracy than random sampling, but, again, it’s disparity is no better despite the fact that it explicitly optimizes for the trade-off.



**Figure 1:** (Left) A scatterplot of demographic disparity versus F-score for a fixed budget  $B = 400$ , for PANDA and baseline approaches. (Right) A similar scatterplot for error rate balance versus F-score. In both cases, the upper-left is optimal behavior. Overall, we see that fairlearn and PANDA are the most competitive algorithms, with flipped behavior with respect to disparity on the two metrics. Dotted curves are algorithms unaware of parity/groups; solid lines are algorithms that are.



**Figure 2:** Learning curves for all algorithms, (Left) budget versus F-score and (Right) budget versus demographic disparity. The constraint value for fairlearn and PANDA is 0.04. We see that PANDA and fairlearn are able to (approximately) achieve the target parity, with PANDA achieving a higher F-score than FAL (which has higher disparity). The black dotted line shows the F-score for a random classifier.

Finally, fairlearn and PANDA dominate in terms of the trade-off, with PANDA achieving higher accuracy, better error rate balance, but worth demographic disparity.

We also consider the dynamic nature of these algorithms as they collect more data. In Figure 2, we plot budget versus f-score and disparity for a fixed parity constraint of 0.04. Unsurprisingly, we see that entropy sampling outperforms random sampling, though they perform essentially the same for disparity. We also see a clear trade-off in FAL between F-score (goes up as budget increases) and disparity (goes up).

Here, we see that both fairlearn and PANDA are able to keep the disparity low (after an initial peak for PANDA). There is a generalization gap between PANDA’s training disparity (which always exactly satisfies the 0.04 constraint) and its validation disparity, which is somewhat higher, as anticipated by concentration bounds on disparity like those of Agarwal et al. [2]. The initial peak in disparity (where it does not satisfy the constraint) for PANDA is not surprising: it is trained end-to-end to pick a good sample of 400 points; it is not optimized for smaller budgets. Similarly, in terms of F-score, PANDA achieves a very high initial F-score, essentially a zero-shot learning type effect. However, as it lowers the disparity, the F-score also

drops slightly. In all cases, the “Fairlearn+AL” baseline achieves better F-score in comparison to the “Fairlearn” baseline, however, this comes at the expense of both the demographic disparity and error rate balance.

At test time, PANDA is the fastest of all the active learning algorithms we compare to from §4.1 with matching runtime performance to random sampling. This is because at test time we only need a single forward pass through the selection policy to select the  $B$  samples to label. Entropy sampling requires computing then entropy in every time step. Fairlearn is much slower as the learning reduction refits a mixture of experts model with different weights. Fair Active Learning is the slowest approach as it needs to compute the “expected fairness” that requires learning a new classifier for every data point in the pool of unlabeled data. For meta-training, learning the policy for PANDA converged after few hours.

## 5 DISCUSSION, LIMITATIONS AND CONCLUSION

We presented PANDA, a meta-learning approach for learning to active learn under parity constraints, motivated by the desire to build an algorithm to mitigate unfairness in machine learning by

collecting more data. We have seen that empirically PANDA is effective experimentally, even in a setting in which it essentially has to choose all 400 points to label at once, rather than one at a time. An obvious direction of future work is to incorporate features of the underlying classifier into the selection policy; the major challenge here is the computational expense of unrolling the corresponding computational graph. One major advantage of PANDA over all other alternatives is that in principle it does not need access to group information at test time. So long as it can be trained with group information available (for measuring disparities on the meta-test data), there is nothing in the algorithm that requires this information at test time. The only other setting in which this is possible is FAL with demographic disparity (precisely because demographic disparity does not need access to *labels*). Exploring this experimentally is a potential next step. Finally, there is the broader question of: how does one know what is the right intervention to mitigate disparities? Should we constrain our classifier? Should we collect more data? More features? Change the architecture? These are all important questions that are only beginning to be explored [9, 22].

Standard concentration bounds on disparity like those of [2] hold for PANDA: if a given disparity is achieved at training time, then PANDA guarantees a bound on the generalization error for the validation disparity. What one would like to show is that after sampling using PANDA, we can guarantee improvement in fairness parity; this theoretical analysis remains a future work item.

## 6 BROADER IMPACT

The motivation of this work is precisely to have positive broader impacts, by giving machine learning practitioners who care about fairness in machine learning another tool in their toolbox to build models with fewer disparities. Our primary target stakeholder population is such machine learning practitioners and data scientists. Secondarily, as that primary stakeholder population builds and deploys algorithms, those who are impacted by those algorithms through direct or indirect use will, we hope, suffer fewer disparities as a result.

There are several risks. The first is a false sense of security. For instance, we do not prove formally that this approach is guaranteed to work in all cases, and our empirical results are limited to a small number of tests over a single dataset. On the positive side, Agarwal et al. [2] prove a generalization bound for disparity that applies to our algorithm (as well as any other algorithm); thus, so long as practitioners properly test the resulting disparities of their models, they can consult these generalization bounds to get estimates of worst case behavior.

A second risk is around, if deployed, how the new data is collected. We have seen news stories recently around unethical practices for data collection. Any additional labeling that is performed as a result of running this or similar algorithms should be done consistent with standard ethical guidelines for data collection.

Overall, while there are real concerns about how this technology might be deployed, our hope is that the positive impacts outweigh the negatives, specifically because standard best-use practices should mitigate most of the risks.

## ACKNOWLEDGMENTS

The authors would like to thank Tom Goldstein for his valuable comments and suggestions regarding the forward-backward splitting method for optimization. The authors would also like to thank the anonymous referees as well as members of the Computational Linguistics and Information Processing (CLIP) lab for their valuable comments and helpful suggestions.

## REFERENCES

- [1] Behnoush Abdollahi and Olfa Nasraoui. 2018. Transparency in fair machine learning: The case of explainable recommender systems. In *Human and Machine Learning*. Springer, Springer, 21–35.
- [2] Alekh Agarwal, Alina Beygelzimer, Miroslav Dudik, John Langford, and Hanna Wallach. 2018. A Reductions Approach to Fair Classification. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, USA, 60–69. <https://proceedings.mlr.press/v80/agarwal18a.html>
- [3] Hadis Anahideh, Abolfazl Asudeh, and Saravanan Thirumuruganathan. 2022. Fair active learning. *Expert Systems with Applications* 199 (2022), 116981. <https://doi.org/10.1016/j.eswa.2022.116981>
- [4] Dana Angluin. 1988. Queries and Concept Learning. *Mach. Learn.* 2, 4 (April 1988), 319–342. <https://doi.org/10.1023/A:1022821128753>
- [5] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2016. Machine bias. ProPublica.
- [6] Philip Bachman, Alessandro Sordani, and Adam Trischler. 2017. Learning Algorithms for Active Learning. In *Proceedings of the 34th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 70)*, Doina Precup and Yee Whye Teh (Eds.). PMLR, PMLR, 301–310. <https://proceedings.mlr.press/v70/bachman17a.html>
- [7] Michelle Bao, Angela Zhou, Samantha Zottola, Brian Brubach, Sarah Desmarais, Aaron Horowitz, Kristian Lum, and Suresh Venkatasubramanian. 2021. It's COMPASlicated: The Messy Relationship between RAI Datasets and Algorithmic Fairness Benchmarks.
- [8] Ángel Alexander Cabrera, Will Epperson, Fred Hohman, Minsuk Kahng, Jamie Morgenstern, and Duen Horng Chau. 2019. FairVis: Visual analytics for discovering intersectional bias in machine learning.
- [9] Irene Y. Chen, Fredrik D. Johansson, and David Sontag. 2018. Why is My Classifier Discriminatory?. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems (Montréal, Canada) (NIPS'18)*. Curran Associates Inc., Red Hook, NY, USA, 3543–3554.
- [10] Alexandra Chouldechova. 2017. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data* 5, 2 (2017), 153–163.
- [11] David Cohn, Les Atlas, and Richard Ladner. 1994. Improving generalization with active learning. *Machine learning* 15, 2 (1994), 201–221.
- [12] Patrick L Combettes and Valérie R Wajs. 2005. Signal recovery by proximal forward-backward splitting. *Multiscale Modeling & Simulation* 4, 4 (2005), 1168–1200.
- [13] Kate Crawford and Ryan Calo. 2016. There is a blind spot in AI research. *Nature* 538, 7625 (2016), 311–313.
- [14] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. 2019. Transformer-xl: Attentive language models beyond a fixed-length context.
- [15] Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Lukasz Kaiser. 2018. Universal transformers.
- [16] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding.
- [17] Dheeru Dua and Casey Graff. 2017. UCI Machine Learning Repository. <http://archive.ics.uci.edu/ml>
- [18] Meng Fang, Yuan Li, and Trevor Cohn. 2017. Learning how to active learn: A deep reinforcement learning approach.
- [19] Valerii Vadimovich Fedorov. 2013. Theory of optimal experiments.
- [20] Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. 2015. Certifying and Removing Disparate Impact. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (Sydney, NSW, Australia) (KDD '15)*. Association for Computing Machinery, New York, NY, USA, 259–268. <https://doi.org/10.1145/2783258.2783311>
- [21] Batya Friedman and Helen Nissenbaum. 1996. Bias in computer systems. *ACM Transactions on Information Systems (TOIS)* 14, 3 (1996), 330–347.
- [22] Sainyam Galhotra, Yuriy Brun, and Alexandra Meliou. 2017. Fairness testing: Testing software for discrimination.
- [23] Vivian Giang. 2018. The potential hidden bias in automated hiring systems.
- [24] Tom Goldstein, Christoph Studer, and Richard Baraniuk. 2014. A field guide to forward-backward splitting with a FASTA implementation.

- [25] Emil Julius Gumbel. 1948. Statistical theory of extreme values and some practical applications: a series of lectures.
- [26] Moritz Hardt, Eric Price, and Nati Srebro. 2016. Equality of opportunity in supervised learning. , 3315–3323 pages.
- [27] Kenneth Holstein, Jennifer Wortman Vaughan, Hal Daumé III, Miro Dudik, and Hanna Wallach. 2019. Improving fairness in machine learning systems: What do industry practitioners need? , 16 pages.
- [28] Eric Jang, Shixiang Gu, and Ben Poole. 2016. Categorical reparameterization with gumbel-softmax.
- [29] Jun Jiang and Horace Ho-Shing Ip. 2008. Active learning for the prediction of phosphorylation sites. , 3158–3165 pages.
- [30] Nathan Kallus and Angela Zhou. 2018. Residual unfairness in fair machine learning from prejudiced data.
- [31] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization.
- [32] Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. 2016. Inherent trade-offs in the fair determination of risk scores.
- [33] Ksenia Konyushkova, Raphael Sznitman, and Pascal Fua. 2017. Learning active learning from data. , 4225–4235 pages.
- [34] Wouter Kool, Herke Van Hoof, and Max Welling. 2019. Stochastic beams and where to find them: The gumbel-top-k trick for sampling sequences without replacement.
- [35] Pierre-Louis Lions and Bertrand Mercier. 1979. Splitting algorithms for the sum of two nonlinear operators. *SIAM J. Numer. Anal.* 16, 6 (1979), 964–979.
- [36] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. 2016. The concrete distribution: A continuous relaxation of discrete random variables.
- [37] Arvind Narayanan. 2018. Translation tutorial: 21 fairness definitions and their politics.
- [38] 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. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. , 8024–8035 pages. <http://papers.nips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
- [39] Piyush Rai, Avishek Saha, Hal Daumé III, and Suresh Venkatasubramanian. 2010. Domain adaptation meets active learning. , 27–32 pages.
- [40] Dan Roth and Kevin Small. 2006. Margin-based active learning for structured output spaces. , 413–424 pages.
- [41] Burr Settles. 2009. *Active learning literature survey*. Technical Report. University of Wisconsin-Madison Department of Computer Sciences.
- [42] Claude E Shannon. 1948. A note on the concept of entropy. *Bell System Tech. J* 27, 3 (1948), 379–423.
- [43] Xiaoxiao Shi, Wei Fan, and Jiangtao Ren. 2008. Actively transfer domain knowledge. , 342–357 pages.
- [44] Megha Srivastava, Hoda Heidari, and Andreas Krause. 2019. Mathematical notions vs. human perception of fairness: A descriptive approach to fairness for machine learning. , 2459–2468 pages.
- [45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. , 5998–6008 pages.
- [46] Michael Veale and Reuben Binns. 2017. Fairer machine learning in the real world: Mitigating discrimination without collecting sensitive data. *Big Data & Society* 4, 2 (2017), 2053951717743530.
- [47] Tim Vieira. 2014. Gumbel-max trick and weighted reservoir sampling. (2014). <https://timvieira.github.io/blog/post/2014/08/01/gumbel-max-trick-and-weighted-reservoir-sampling/>.
- [48] Sara Wachter-Boettcher. 2017. AI recruiting tools do not eliminate bias.
- [49] Xuezhi Wang, Tzu-Kuo Huang, and Jeff Schneider. 2014. Active Transfer Learning under Model Shift. In *Proceedings of the 31st International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 32)*, Eric P. Xing and Tony Jebara (Eds.), PMLR, Beijing, China, 1305–1313. <https://proceedings.mlr.press/v32/wangi14.html>
- [50] L. Wightman. 1998. LSAC National Longitudinal Bar Passage Study. (1998). LSAC.
- [51] Benjamin Wilson, Judy Hoffman, and Jamie Morgenstern. 2019. Predictive inequity in object detection.
- [52] Tal Zarsky. 2016. The trouble with algorithmic decisions: An analytic road map to examine efficiency and fairness in automated and opaque decision making. *Science, Technology, & Human Values* 41, 1 (2016), 118–132.