

# Balancing the Scales: Reinforcement Learning for Fair Classification

Anonymous ACL submission

## Abstract

Fairness in classification tasks has traditionally focused on bias removal from neural representations, but recent trends favor algorithmic methods that embed fairness into the training process. These methods steer models towards fair performance, preventing potential elimination of valuable information that arises from representation manipulation. Reinforcement Learning (RL), with its capacity for learning through interaction and adjusting reward functions to encourage desired behaviors, emerges as a promising tool in this domain. In this paper, we explore the usage of RL to address bias in multi-class classification by scaling the reward function to mitigate bias. We employ the contextual multi-armed bandit framework and adapt three popular RL algorithms to suit our objectives, demonstrating a novel approach to mitigating bias<sup>1</sup>.

## 1 Introduction

In recent years, the issue of bias and fairness in Artificial Intelligence and Natural Language Processing has received significant attention (Mehrabi et al., 2021). In decision-making models such as classification algorithms, bias often stems directly from the training data leading to unfair outcomes between protected groups such as gender or race. To address this problem, previous work on fairness has often focused on achieving *representational fairness*, so that the information of the protected groups is lost (Ravfogel et al., 2020; Haghighatkah et al., 2022). However, recent work has shown that there is no meaningful correlation between representational fairness and *empirical fairness*, i.e. fairness on downstream tasks (Shen et al., 2022). To address empirical fairness directly, other work has explored the intersection of bias mitigation and class-imbalanced learning (Subramanian et al., 2021). Class-imbalanced learning

approaches aim to achieve fair performance by balancing the training data via sampling or reweighing the loss function.

At the same time, Reinforcement Learning (RL) has emerged as a promising alternative to traditional supervised learning methods for various NLP tasks, including syntactic parsing, conversational systems, and machine translation (Uc-Cetina et al., 2023). Unlike traditional supervised learning methods, RL is not bound to binary labels and is trained directly on the continuous value of each input, as illustrated in Figure 1. RL agents can learn from sparse reward signals, receiving the rewards for the action they choose, not necessarily the correct one. By exploring the environment and adapting their behavior based on the received rewards, RL agents find optimal actions under varying state values. In the context of classification, RL has been adapted to mitigate class imbalance through a scaling component of the reward function for binary classification (Lin et al., 2020). However, implementations for complex cases such as multi-class imbalanced classification remain largely unexplored.

In this work, we leverage RL to address fairness among protected groups in multi-class classification. First, we propose to frame the fair classification task as a Contextual Multi-Armed Bandit (CMAB) problem, see Figure 1 for an overview of our setup. To mitigate bias, we scale the reward function to counteract imbalances among protected groups within each class. We employ three different types of RL methods, each reflecting a key type of RL approach, and adapt them for our task. Additionally, we integrate the different scaling approaches into a supervised learning baseline to evaluate the effectiveness of our RL-based methods.

Experiments on two fair classification datasets demonstrate that our RL algorithms achieve competitive performance compared to existing baselines and that reward scaling is a powerful tool to mitigate bias in classification. We further inves-

<sup>1</sup>Our code is available at [https://anonymous.4open.science/r/RL\\_for\\_imbalanced\\_classification-755E](https://anonymous.4open.science/r/RL_for_imbalanced_classification-755E)

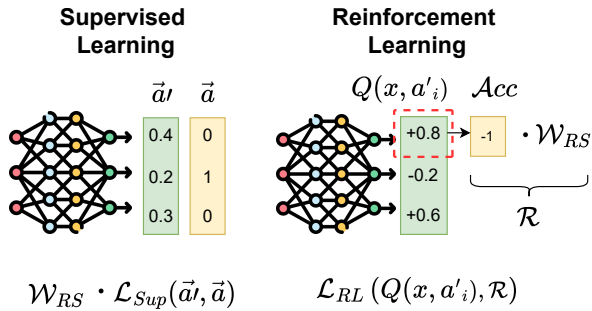


Figure 1: Overview of the classification setup with input vector  $x$ , and output class  $a$  for and Reinforcement Learning and Supervised Learning, highlighting the place of the reward scaling matrix  $\mathcal{W}_{RS}$

081 tigate how stable reward scaling is under various  
082 class and subclass imbalances as well as various  
083 degrees of *representational fairness*. Notably, the  
084 deep RL algorithms perform best on the multi-class  
085 dataset, while the classical CMAB algorithm excels  
086 on the binary dataset. Moreover, our scaled super-  
087 vised implementation improves existing implemen-  
088 tations and shows state-of-the-art performance on  
089 multi-class setting.

## 090 2 Related work

### 091 2.1 Bias Mitigation

092 Research on mitigating bias can be divided in those  
093 that tackle the training data (Wang et al., 2019),  
094 those that attempt to remove bias from representa-  
095 tions (Ravfogel et al., 2020; Haghighatkhah et al.,  
096 2022), and those that adjust the learning process  
097 (Elazar and Goldberg, 2018; Han et al., 2021).  
098 Within approaches that adjust the learning pro-  
099 cesses, we distinguish two main categories: those  
100 that add adversarial learners to ignore protected at-  
101 tributes (Wadsworth et al., 2018), and more closely  
102 to our work, approaches that adjust the loss func-  
103 tion to emphasize performance on minority classes.

104 Prior work that modified the training setup  
105 to increase fairness used methods such as down/upsampling (Wang et al., 2019) and reweighting the loss function (Höfler et al., 2005; Lahoti et al., 2020). Han et al. (2022a) evaluate both down-sampling and loss reweighting on two datasets for fair text classification. Both techniques are applied to align training with different definitions of fairness. Downsampling using the Equal Opportunity fairness metric demonstrated impressive results. In this paper, we take the first step to explore whether reward scaling in reinforcement learning can improve fairness in classification.

## 117 2.2 Reinforcement Learning for Classification

118 Literature on RL applications for classification pre-  
119 dominantly considers the following two theoretical  
120 frameworks: Markov Decision Process (MDP) and  
121 the Contextual Multi-Armed Bandit (CMAB).  
122

123 Early work by Wiering et al. (2011) casts clas-  
124 sification as a sequential decision-making task, by  
125 introducing a classification variant of the MDP. In  
126 their setup agents manipulate memory cells to en-  
127 code information by applying an action sequence  
128 on a single sample. They demonstrated competitive  
129 performance, but, remained limited to small tasks  
130 due to the computational complexity. Lin et al.  
131 (2020) extended this work, by introducing a vari-  
132 ant of the classification MDP and applying a Deep  
133 Q-learning Network (DQN) to binary classification  
134 of images and texts. They focused on mitigating  
135 bias arising from class imbalance by scaling the re-  
136 wards inversely proportional to the class frequency.  
137 However, in their setup the sequential component  
138 was taken over multiple data points, which assumes  
139 sequential dependency among data samples in the  
140 classification task.

141 The RL framework CMAB offers a promising  
142 alternative because it considers the input as a se-  
143 quence of independent states. We formalize our  
144 classification task as a CMAB problem, because  
145 this is consistent with the independence of data  
146 points in the commonly shuffled datasets. Dudík  
147 et al. (2014) use CMAB agents by modifying K-  
148 class classification as a K-armed bandit problem,  
149 where the agent receives a reward of 1 for correct  
150 and 0 for incorrect classification. Dimakopoulou  
151 et al. (2019) use this framework and modify differ-  
152 ent CMAB algorithms to balance exploration and  
153 exploitation and compare the original and modi-  
154 fied agents on 300 classification datasets. However,  
155 their analysis focused on datasets with either lim-  
156 ited classes, features, or observations. To the best  
157 of our knowledge, we are the first to extend reward  
158 scaling for fair multi-class classification or to apply  
159 reward scaling for classification with CMAB.

## 159 3 Methodology

160 In this section, we describe how we formalize our  
161 classification task as a CMAB. We introduce three  
162 RL methods and explain how we adapt them for  
163 fair classification.<sup>2</sup>

<sup>2</sup>Due to space limitations, we only summarize the key idea of algorithms and how we adapt them in the paper. Please refer to the Appendix and original papers for more details.

### 3.1 Contextual Multi-Armed Bandit

We formalize the multi-class classification task as a finite contextual multi-armed bandit (CMAB) problem. In each round  $t$ , an agent is presented with a context vector  $x_t \in \mathbb{R}^d$ . The agent chooses an action  $a_t \in A$  from a fixed set of arms, based on the policy  $a_t \sim \pi(x_t)$ . After the action is taken, the environment returns a reward:  $r_t \sim \mathcal{R}$ . In a multi-class classification framework, the action space is the set of all possible classes, while the context vector is a representation of the input, e.g. a contextual text embedding (see Section 4.1 for more information). Within a finite number of rounds, the agent aims to learn the optimal policy to maximize the total reward. In other words, given a set of testing data, we aim to learn the optimal policy to maximize the selection of correct classes.

We extend the CMAB framework for fair classification by constructing a reward function that counters data imbalances. We assign a reward scale for each sensitive state  $(a, g)$ , comprising the desired class  $a$  (e.g. occupation) and protected attribute  $g$  (e.g. gender). The total reward for a given prediction is calculated as  $\mathcal{R}(a, a_{pred}, g) = \mathcal{Acc}(a, a_{pred}) \cdot \mathcal{W}(a, g)$ . It comprises an accuracy term  $\mathcal{Acc}$ , and a reward scale matrix  $\mathcal{W}$ . Unlike previous work (Dudík et al., 2014), which defines the accuracy term as  $\mathcal{Acc} \in \{0, 1\}$ , we define it as  $\mathcal{Acc} \in \{-1, 1\}$ . Which allows us to scale the reward for both correct (+1) and incorrect classifications (-1). We use the term *reward scale* to indicate that this approach adjusts the magnitude but not the sign of the reward. Section 3.3 presents various designs of the reward scale.

### 3.2 Reinforcement Learning Algorithms

We select three different RL algorithms and adapt them to learn optimal policies for fair classification in the formalized CMAB problem. These algorithms include one classical CMAB algorithm that addresses the linear relationship between the expected reward and the context, as well as two popular deep RL algorithms for MDP problems, Deep Q-Network (DQN) and Proximal Policy Optimization (PPO), which allow us to leverage non-linear approximations. The two deep RL algorithms are selected as they are representative of the two key types of deep RL approaches: value-based methods and policy gradient methods. By employing these three algorithms, we aim to investigate the application of diverse RL methods.

#### 3.2.1 LinUCB

The classical CMAB algorithm, disjoint Linear UCB (LinUCB) (Li et al., 2010) assumes a linear relationship between the context embedding  $x_t$  and the reward  $E[r_{t,a}|x_t] = x_t^\top \theta_a$ . A benefit of disjoint LinUCB over other CMAB algorithms is that each class has a unique learnable weight vector  $\theta_a$ , which makes it suitable for classification with many classes. In each round, the agent chooses the arm (i.e. class label) with the highest score  $\hat{\theta}_a^\top x_t + \alpha \sqrt{x_t^\top A_a^{-1} x_t}$ , based on the context vector  $x_t$ . This is a combination of the mean of the expected payoff,  $\hat{\theta}_a^\top x_t$ , and the standard deviation  $\sqrt{x_t^\top A_a^{-1} x_t}$ , weighted with parameter  $\alpha$  to control the level of exploration. The weight vector of each arm is defined as  $\hat{\theta}_{a_t} = A_{a_t}^{-1} b_{a_t}$ . Here the covariant matrix  $A_{a_t}$  is calculated with the history of context vectors chosen by that arm,  $A_a = \lambda I_d + \sum_{s=1}^{t-1} x_s x_s^\top$ . The vector  $b_a$  is the mean context vector of the arm weighted by the obtained rewards,  $b_{a_t} = \sum_{s=1}^t r_{s,a_t} x_{s,a_t}$ .

#### 3.2.2 DQN<sub>bandit</sub>

To adapt the MDP algorithms for a CMAB problem, our CMAB implementation is congruent with a one-step MDP, where each initial state is sampled from the existing set of context  $s_1 \in X$ , and each second state is a terminal state. In DQN (Mnih et al., 2015), the agent learns a Q-function, parameterized by  $\phi$ , to estimate the return for each state-action pair. According to the Bellman equation (Bellman, 1957), the optimal Q-value,  $Q^*$ , of two sequential states are linked by:

$$Q^*(s, a) = \mathbb{E}_\pi[r_t + \gamma \max_{a'} Q^*(s_{t+1}, a')] \quad (1)$$

In our case (the one-step MDP), each next state is the terminal state, after which there is no reward, thus we obtain,  $Q_\phi(s_{t+1}, \cdot) = 0$ , and  $G_t = r_t$ . The parameters of  $\phi$  are optimized using the mean-squared error between the current Q-value,  $Q_\phi(s_t, a_t)$ , and the updated value provided in Equation 1. The updated value is computed as  $r_t + \gamma \max_{a'} Q_\phi(s_{t+1}, a')$ , but since the next state is always the terminal state it reduces to  $r_t$ . We finalize the adaptation of DQN for the CMAB by casting the states as context vectors, obtaining the loss function:

$$L_{DQN}(\phi) = \mathbb{E}_{(x_t, a_t, r) \in B} [(r - Q_\phi(x_t, a_t))^2]$$

The network is updated by sampling a minibatch of tuples  $B$  from the replay buffer. The  $DQN_{bandit}$

enables exploration using an  $\epsilon$ -greedy policy for selecting actions.

### 3.2.3 PPO<sub>bandit</sub>

Different from DQN, in Proximal Policy Optimization (PPO) (Schulman et al., 2017), the policy  $\pi$  (parameterized by  $\theta$ ) is directly optimized under the objective of selecting the best action. The general objective in policy gradient methods is to maximize:  $\mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \cdot A_t \right]$ . The advantage  $A_t$  is computed as  $A_t = Q_\pi(s_t, a_t) - V_\phi(s_t)$ , where a critic network  $V_\phi$  is used to estimate the state value. PPO ensures the policy does not deviate too far during an update, by scaling the advantage with the probability ratio,  $r_t(\theta)$ . This ratio is clipped to create a conservative lower bound to control the policy’s change at each step. The actor’s objective function is thus defined as:

$$L_{actor}(\theta) = \mathbb{E}_t [\min(r_t(\theta)A_t, \text{clip}_\epsilon(r_t(\theta))A_t)]$$

To adapt PPO for CMAB, the sequential component is removed and the state  $s_t$  is replaced by the context vector  $x_t$ . For the actor loss, the advantage changes and is calculated as  $A_t = r_t - V_\phi(x_t)$ . The return again reduces to the reward, thereby simplifying the critic loss to:

$$L_{critic}(\phi) = \mathbb{E}_t [(V_\phi(x_t) - r_t)^2]$$

Lastly, the final loss of the PPO<sub>bandit</sub> agent contains a penalty that maximizes the policy’s entropy of the context vector to encourage exploration.

### 3.3 Reward Scales

Lin et al. (2020) implement reward scaling on an imbalanced binary classification dataset. Inspired by this, we propose different ways of reward scales for multi-class classification with imbalances of protected attributes. For context, we use the profession classification dataset, BiasBios, where reward scaling tackles the sub-class imbalance of the protected attribute gender. To illustrate the influence of various reward scales Figure 2 shows the scales of a balanced (Professor) and an imbalanced (Nurse) class for the protected groups with attribute gender.

For the first method, we extend the implementation of Lin et al. (2020) into the multi-classification setting and reduce the reward for the majority by scaling it with the imbalanced ratio  $\rho_{imb}^a = \frac{|D_{min}^a|}{|D_{maj}^a|}$ , which is the ratio between the number of samples

of the minority and majority class in class  $a$ .

$$\mathcal{W}_{\rho^+}(a, g) = \begin{cases} 1 & \text{if } g \text{ is minority in } a \\ \rho_{imb}^a & \text{if } g \text{ is majority in } a \end{cases}$$

Figure 2 demonstrates that  $\mathcal{W}_{\rho^+}(x)$  scales with a reverse of the bias within a class, however, compared to a balanced class, the reward scale of the majority is very low. Therefore, we propose a second design that keeps the scales of the majority group in the imbalanced class equal to the scales of the balanced class. Thus we set the majority value at 1 and only increase the minority value, based on the inversed imbalanced ratio.

$$\mathcal{W}_{\rho^-}(a, g) = \begin{cases} (\rho_{imb}^a)^{-1} & \text{if } g \text{ is minority in } a \\ 1 & \text{if } g \text{ is majority in } a \end{cases}$$

The third implementation adopts the Equal Opportunity (EO) formalization used by Han et al. (2022a). Contrary to the previous two methods it ensures the average weights per class remain equal, providing an improved theoretical fairness among classes. The EO objective is achieved by aggregating the loss per sensitive state and then scaling it. However, our work scales per instance, thus we convert the EO objective to instance-specific weights (see Appendix B.3) and obtain:

$$\mathcal{W}_{EO}(a, g) = \frac{1}{2} \frac{1}{P(g|a)}$$

Lastly, we also employ the Inverse Probability Weighting (IPW) technique (Höfler et al., 2005). Full fairness across classes and protected groups is obtained by scaling with the joint probability, resulting in:

$$\mathcal{W}^{IPW}(a, g) = \frac{1}{P(a, g)}$$

### 3.4 Supervised Learning: Loss reweighting

Parallel to reward scaling in RL is (instance) reweighting in supervised learning (Han et al., 2022a; Lahoti et al., 2020), here loss reweighting for clarity. Loss reweighting has been a popular technique for imbalanced datasets, where the loss of each data sample is scaled to mitigate the class imbalance, traditionally using the IPW (Höfler et al., 2005). The weighted cross-entropy loss using the true probability  $p$ , predicted probability  $q$ :

$$L^{CE} = - \sum_{x, g} \sum_a \mathcal{W}(a, g) p(a|x) \log q(a|x)$$

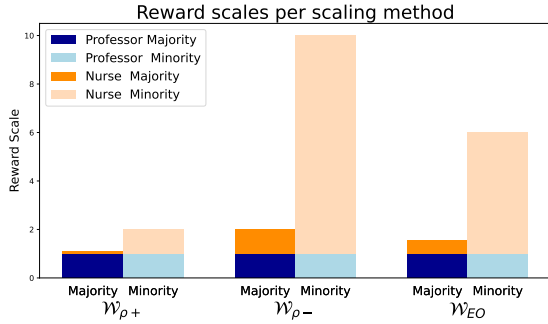


Figure 2: Reward scales for the professions *Professor* and *Nurse* using the different scaling functions. Professor (50/50) and Nurse (90/10)

We implement supervised learning with loss reweighing to serve as a strong baseline and highlight the connection between loss reweighing and reward scaling.

## 4 Experiments

### 4.1 Dataset

The BiasBios (De-Arteaga et al., 2019) consists of 393,423 biographies labeled with one of 28 professions, and a binary gender label. Following De-Arteaga et al. (2019), the data is randomly split according to 65% training, 25% testing, and 10% for validating. The dataset contains two imbalances: varying frequencies of the professions and a difference in gender percentage for each class.

Following Ravfogel et al. (2020) and Han et al. (2022a) we also evaluate on the Emoji (Elazar and Goldberg, 2018) sentiment analysis task of Twitter data (Blodgett et al., 2016). The task involves binary sentiment classification evaluation with race as the protected attribute, approximated through the provided labels Standard American English (SAE) and African American English (AAE). As per Han et al. (2021), the dataset is composed of Happy (40% AAE, 10% SAE), and Sad: (10% AAE, 40% SAE). We use the same train, dev, and test splits of 100k/8k/8k instances, respectively.

**Context Vectors** Each textual data sample is embedded into a context vector via a pretrained encoder, enabling the algorithms for classification. Following Ravfogel et al. (2020) we use the same fixed pretrained encoder for each dataset. For the BiasBios dataset, each biography is encoded using the [CLS] output of the uncased BERT-base model (Devlin et al., 2019). For the Emoji dataset, we use the DeepMoji encoder (Felbo et al., 2017), which

has been demonstrated to capture a diverse range of moods and demographic information.

### 4.2 Metrics

Following prior work, we evaluate performance using accuracy and fairness using the True Positive Ratio (TPR) gap (De-Arteaga et al., 2019; Ravfogel et al., 2020). The TPR gap of a class  $a \in A$  is calculated as:  $TPR_{gap}^a = TPR_g^a - TPR_{\sim g}^a$ , where  $g$  and  $\sim g$  represent the two options for the sensitive states. The global TPR metric, GAP, is then calculated as the root mean square of the individual metrics:

$$GAP = \sqrt{\frac{1}{|A|} \sum_{a \in A} (TPR_{gap}^a)^2} \quad (2)$$

To quantify performance and fairness as a single metric we use the *Distance To the Optimum* (DTO) introduced in Han et al. (2022a). DTO combines the metrics (accuracy, 1-GAP) as dimensions of evaluation space and computes the Euclidean distance between the achieved and Utopian point. The smaller the distance to the Utopian point (lower DTO), the better. We report the DTO with the Utopian accuracy and GAP as the best values across all evaluated models.

While accuracy measures the overall performance and GAP the disparity among protected groups within a class, these metrics do not capture imbalance performance across classes. Therefore we also evaluate our algorithms using the macro-averaged F1 metric to detect if minority classes are ignored. All metrics are scaled by 100 for ease of reading and all metrics are represented in the tables as the mean  $\pm$  std over 5 random seeds, except DTO which is taken over the mean.

### 4.3 Hyperparameters and model selection

Each algorithm uses the same classifier architecture, except LinUCB, which has a custom set of learnable parameters. The classifier has one hidden layer MLP. All models are trained for 10 epochs, except LinUCB, which achieved optimal performance within 2 epochs. All models are evaluated on the validation set after 50k iterations to account for different convergence speeds of models. The best model throughout training and across hyperparameters is selected using DTO. We apply hyperparameter optimization on both datasets for each of the algorithms, for details see Appendix A.3.

	PPO <sub>bandit</sub>		Sup	
	Accuracy $\uparrow$	GAP $\downarrow$	Accuracy $\uparrow$	GAP $\downarrow$
$\mathcal{W}_{\rho+}$	74.6 $\pm$ 0.7	9.9 $\pm$ 0.8	79.3 $\pm$ 0.1	7.9 $\pm$ 0.3
$\mathcal{W}_{\rho-}$	78.8 $\pm$ 0.1	<b>8.4 <math>\pm</math> 0.6</b>	79.8 $\pm$ 0.3	6.9 $\pm$ 0.2
$\mathcal{W}_{EO}$	<b>79.2 <math>\pm</math> 0.2</b>	8.5 $\pm$ 0.2	<b>80.1 <math>\pm</math> 0.2</b>	7.1 $\pm$ 0.5
$\mathcal{W}_{IPW}$	45.8 $\pm$ 6.9	10.5 $\pm$ 0.9	72.1 $\pm$ 0.7	<b>6.1 <math>\pm</math> 0.3</b>

Table 1: Results with different reward scales for Supervised Learning (Sup) and PPO on BiasBios

#### 4.4 Comparison Models

Besides the supervised implementation in Section 3.4, abbreviated to **Sup**, we also compare our models against various existing debiasing methods. **INLP** (Ravfogel et al., 2020) debiases embeddings by iteratively training classifiers to predict the protected attribute, it then removes this information from the embedding using a projection of the classifier’s nullspace. **MP** (Haghighatkhah et al., 2022) simplifies the INLP setup by using a single Mean Projection (MP) between the representation of each class’s protected groups. **DAdv**, (Han et al., 2021) removes sensitive information from the embeddings by applying adversarial training using diverse adversaries. Lastly, **BTEO** (Han et al., 2022a) subsamples the dataset to establish equal opportunity. We implement these existing methods with the same training settings for fair comparison. Notably, we highlight how Supervised  $\mathcal{W}^{EO}$  is theoretically equal to instance reweighing in Han et al. (2022a), but our implementation achieves significantly higher performance.

### 5 Results and Analysis

We train each of the RL algorithms with and without reward scaling. As a strong baseline, we also train a supervised learning model in the standard fashion and use loss reweighting with the same reward scale function.

#### 5.1 Different Reward Functions

We first evaluate the effect of the different reward scales discussed in Section 3.3 by providing the results for the implementations of PPO and Sup, see Table 1, for full table see Appendix D.4.

The results presented in the table demonstrate that the imbalance ratio  $\rho$  yields substantial gains in fairness and accuracy when applied to increase the reward for the minority class ( $\mathcal{W}_{\rho-}$ ) as opposed to decreasing the reward for the majority class ( $\mathcal{W}_{\rho+}$ ). Especially the accuracy of PPO is sensitive to this,

suggesting that PPO might not work very well for states with low reward scales.

As hypothesized, scaling with the joint probability of class and protected attribute as in  $\mathcal{W}_{IPW}$  proves to be too unstable. It results in the worst accuracy for both algorithms, with a minor improvement in fairness for supervised learning. The overall difference between  $\mathcal{W}_{EO}$  and  $\mathcal{W}_{\rho-}$  is minimal as expected from their similar reward scales depicted in Figure 2. We use EO in our experiments because of its better theoretical foundation.

#### 5.2 Main Results

The main results of our experiments are summarized in Table 2.

On the BiasBios dataset, our DQN and PPO implementations achieve strong results, with PPO outperforming DQN in fairness, as PPO’s lower GAP shows. Our supervised implementation, Sup  $\mathcal{W}^{EO}$ , surpasses all other baselines. LinUCB performs significantly worse on this task, obtaining one of the worst DTO scores. In contrast, on the Emoji dataset, LinUCB achieves one of the best performance-fairness trade-offs, as indicated by the low DTO score, and is only surpassed by BTEO, which has a slightly higher accuracy. On the other hand, DQN obtains a lower accuracy and PPO lower fairness compared to other metrics.

These results suggest that the classical CMAB algorithm LinUCB excels on binary datasets, while the deep RL implementations, DQN and PPO, perform better in multi-class settings. Notably, although DQN and PPO obtain competitive results on the BiasBios, their F1 score is considerably lower than the baselines. Further analysis of per class metrics (see Appendix D.2) reveals that while the F1 for most classes was on par with the supervised setup, both deep RL algorithms failed to recall two of the very sparse classes.

Compared to baseline methods such as BTEO and DAdv, our DQN and PPO implementations demonstrate competitive performance on the BiasBios dataset. Table 2 also shows that, contrary to previous work (Han et al., 2022a),<sup>3</sup> loss-scaling for supervised learning (Sup  $\mathcal{W}^{EO}$ ) achieves superior overall performance to downsampling (BTEO). Downsampling only seems to outperform scaling when enough data is present, as demonstrated by its lower GAP for the Emoji dataset.

<sup>3</sup>The EO scaled supervised implementation of Han et al. (2022a) achieves an Accuracy of 75.7 and GAP of 13.9

Algorithm	BiasBios (28 Classes)					Emoji (2 Classes)			
	Accuracy $\uparrow$	GAP $\downarrow$	DTO $\downarrow$	F1 $\uparrow$	Time $\downarrow$	Accuracy $\uparrow$	GAP $\downarrow$	DTO $\downarrow$	Time $\downarrow$
Sup	81.0 $\pm$ 0.1	16.4 $\pm$ 0.5	9.3	73.8 $\pm$ 0.3	1.0	72.3 $\pm$ 0.1	38.1 $\pm$ 0.6	28.3	1.0
INLP	80.2 $\pm$ 0.6	9.7 $\pm$ 0.4	2.8	71.7 $\pm$ 1.4	50.1	63.5 $\pm$ 3.6	24.1 $\pm$ 5.4	18.6	3.6
MP	<b>81.1 <math>\pm</math> 0.1</b>	13.9 $\pm$ 0.6	6.8	<b>74.0 <math>\pm</math> 0.2</b>	2.6	71.8 $\pm$ 0.3	17.1 $\pm$ 1.0	8.1	2.3
BTEO	79.2 $\pm$ 0.3	8.4 $\pm$ 0.6	2.3	68.1 $\pm$ 0.4	1.7	75.4 $\pm$ 0.1	10.4 $\pm$ 1.0	<b>0.4</b>	0.8
DAdv	80.8 $\pm$ 0.2	8.5 $\pm$ 0.6	1.4	72.9 $\pm$ 0.4	4.8	<b>75.6 <math>\pm</math> 0.3</b>	11.6 $\pm$ 1.7	1.6	5.7
Sup <sup>EO</sup>	80.1 $\pm$ 0.2	<b>7.1 <math>\pm</math> 0.5</b>	<b>1.0</b>	71.7 $\pm$ 0.5	1.0	75.5 $\pm$ 0.1	11.4 $\pm$ 1.1	1.4	1.0
LinUCB <sup>EO</sup>	74.6 $\pm$ 0.2	12.2 $\pm$ 0.5	8.3	59.8 $\pm$ 1.1	31.9	75.3 $\pm$ 0.2	10.4 $\pm$ 0.7	0.5	2.8
DQN <sup>EO</sup> <sub>bandit</sub>	79.2 $\pm$ 0.1	10.1 $\pm$ 0.4	3.6	66.4 $\pm$ 0.2	57.4	70.8 $\pm$ 0.8	<b>10.0 <math>\pm</math> 1.0</b>	4.8	30.2
PPO <sup>EO</sup> <sub>bandit</sub>	79.2 $\pm$ 0.2	8.5 $\pm$ 0.2	2.4	66.0 $\pm$ 0.8	2.9	75.4 $\pm$ 0.1	14.4 $\pm$ 0.6	4.4	3.0

Table 2: Results on the BiasBios and Emojis classification datasets for our own models (in grey) and the baselines. Metrics are provided as mean  $\pm$  std over 5 random seeds, except DTO which is computed over the mean Accuracy, and GAP, and Time which is the relative time compared to the supervised baseline (first row).

Algo + $\mathcal{W}^{EO}$	Accuracy $\uparrow$	GAP $\downarrow$	F1 $\uparrow$
Sup	<b>81.0</b> (- 0.9)	16.4 (- 9.3)	<b>73.8</b> (- 2.1)
LinUCB	78.4 (- 3.8)	15.5 (- 3.3)	67.3 (- 7.5)
DQN <sub>bandit</sub>	80.1 (- 0.9)	<b>13.7</b> (- 3.6)	66.5 (- 0.1)
PPO <sub>bandit</sub>	79.7 (- 0.5)	14.4 (- 5.9)	67.5 (- 1.5)

Table 3: Results on the Bias dataset without reward scaling, presented as mean and difference from the case without EO, where red (worse), blue (better).

### 5.3 Reward Scaling Impact

We investigate the influence of reward scaling on our models by training them with and without scaling. Table 3 presents the results on BiasBios as the mean performance without scaling and the change in metrics when EO scaling is applied.

Without reward scaling the three RL algorithms achieve similar accuracy to the supervised approach but at the cost of a lower F1 score. As mentioned above, the RL algorithms fail on two very sparse classes, which explains the drop in GAP and F1. Failing to classify any instances of a class correctly results in a TPR gap of 0 for that class, since the result is "fair" among both genders.

The EO reward scale significantly reduces the GAP of all implementations, at the cost of a slight decrease in Accuracy and F1 for most models. However, on LinUCB the scaling causes a large performance reduction with only a small GAP reduction, suggesting that scaling hinders the performance more than it improves the fairness.

To inspect the weak effect of reward scaling on LinUCB, we analyze the TPR gap per profession against the gender imbalance for both cases in Figure 3. Without scaling, LinUCB’s performance follows a predictable positive correlation with gender imbalance. For instance, in the Nurse class

( $\sim 90\%$  women), the model performs better for the majority group (women), resulting in a positive TPR gap. However, after reward scaling this correlation is inverted, causing the model to perform worse for the majority group and better for the minority group. In case of the Nurse class, the model obtains a negative TPR gap. This suggests LinUCB is oversensitive to scaling on the BiasBios dataset, causing it to overcompensate and penalize the majority group.

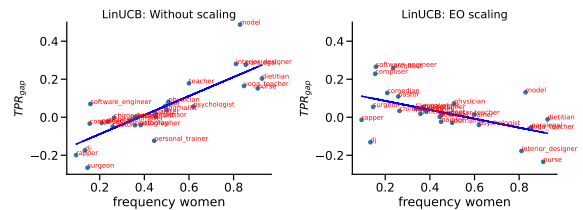


Figure 3: TPR gap plotted against the gender distribution per profession for LinUCB. Left without reward scaling and right with EO reward scaling

### 5.4 Subclass Imbalance Sensitivity

For a more in-depth analysis of each model’s sensitivity to subclass imbalance, we train them on the Emoji dataset under various stereotyping ratios. A stereotyping ratio represents the proportion of the AAE and SAE samples in each class. For example, a stereotyping ratio of 0.2 means the data is distributed as Happy (20% AAE, 80% SAE), Sad (80% AAE, 20% SAE).

Figure 4 reveals a strong inverse relationship between LinUCB’s fairness and the stereotyping ratio. Although the stereotypical ratios are symmetric at the value of 0.5 the fairness of LinUCB is asymmetric at this value. Thus there is a residual representation bias in the data that is not addressed

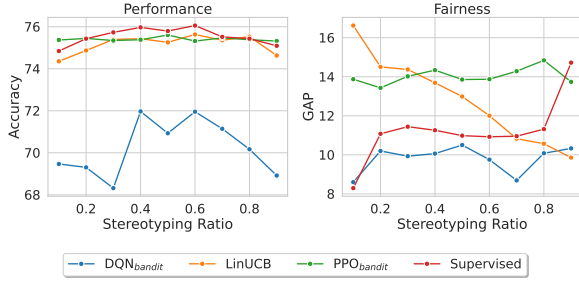


Figure 4: Performance (Accuracy) and Fairness (GAP) on the Emoji dataset using different stereotyping ratios. All models use the scaling of  $\mathcal{W}^{EO}$ .

by the reward scaling. In contrast, supervised learning maintains a mostly stable fairness, except for the most extreme ratios. Interestingly, LinUCB reveals a reverse pattern in best and worst fairness. The relatively low accuracy of DQN and poor performance on fairness of PPO are consistent across ratios. However, PPO does have the most constant fairness and performance across stereotyping ratios, indicating good training stability.

### 5.5 Signal strength vs. Scaling

We now examine how the strength of the protected information affects the efficacy of reward scaling. We focus on two scenarios that modify the gender signal in the representations: 1) adding explicit gender information, thus increasing the gender signal strength 2) debiasing the embeddings using MP, which reduces it. Table 4 presents the mean and relative difference compared to the results without the specified modification.<sup>4</sup>

Providing the model with gender information increases the overall accuracy. However, the impact on fairness, as indicated by the GAP score varies among algorithms. The GAP score increases for the two algorithms with the lowest GAPs (Sup, PPO) and decreases for the two with the highest GAP (LinUCB, DQN). Only the algorithms that perform less well on fairness benefited from access to protected attribute.

Removing the bias with MP reduces the test accuracy for nearly all algorithms, indicating some useful information is removed. Again, the modification increased relatively low GAP scores, and decreased relative high scores. As such, changing to a representation with relatively low bias helps LinUCB and DQN, whereas Sup and PPO that already achieved better fairness mainly see their overall performance hindered.

<sup>4</sup>Full tables in Appendix D.5

Algo+ $\mathcal{W}^{EO}$	Explicit Gender Info		MP-Debiased	
	Accuracy $\uparrow$	GAP $\downarrow$	Accuracy $\uparrow$	GAP $\downarrow$
Sup	<b>80.2 (+0.1)</b>	<b>7.2 (+0.1)</b>	<b>80.0 (-0.1)</b>	<b>7.4 (+0.3)</b>
LinUCB	74.5 (+0.1)	11.7 (-0.5)	74.3 (-0.3)	11.5 (-0.7)
DQN <sub>bandit</sub>	79.2 (+0.0)	10.0 (-0.1)	79.0 (-0.2)	8.6 (-1.5)
PPO <sub>bandit</sub>	79.3 (+0.1)	8.7 (+0.2)	79.2 (+0.0)	9.7 (+1.2)

Table 4: Results on the BiasBios with added gender info (left) and MP-debiased (right), presented as mean, and difference without change: **red** (worse), **blue** (better).

Notably, the differences in Table 4 are relatively small and hardly ever surpass the standard deviation provided in Table 2. This suggests that while the strength of the protected information influences performance and fairness, the impact might be less pronounced than the choice algorithmic design. Moreover, all methods reduced the GAP score compared to only applying MP (which yields a GAP of 13.9, see Table 2).

## 6 Discussion and Conclusion

This paper introduces a novel approach to fair classification using the Contextual Multi-Armed Bandit (CMAB) framework and explores various Reinforcement Learning (RL) algorithms. Our findings demonstrate the potential of different RL algorithms for this task and the efficacy of reward scaling in mitigating imbalances of protected groups. The results show the MDP-derived deep RL methods perform best on the multi-class dataset, while the classical bandit algorithm, LinUCB, excels on the binary dataset. Moreover, our scaled supervised learning implementation achieved new state-of-the-art results on the complex BiasBios dataset.

Our experiments also revealed two limitations 1) RL algorithms can ignore some very sparse classes, despite performing well under most class imbalances. 2) Reward scaling for LinUCB can impair majority group performance beyond that of the minority group. However, the effect of reward scaling remains robust across varying strengths of the protected information, highlighting its potential as a powerful tool for achieving fair outcomes.

Despite these challenges, we believe that the proposed framework presents a promising approach to leverage RL algorithms for fair classification, opening up new research avenues. We encourage future work to extend upon our framework by exploiting different RL characteristics, such as model updates for MDP algorithms based on non-differentiable fairness metrics.



## 610 Limitations

611 Important limitations of this work can be divided  
612 into two sections: 1) Limitations of the dataset  
613 and data requirements of our models 2) Limita-  
614 tions specific to our algorithms and experiments,  
615 independent of the data.

616 **Data limitation** Firstly, all datasets considered  
617 in this study used English text, which restricts the  
618 analysis and might miss other types of biases re-  
619 lated to different linguistic and cultural contexts.  
620 Secondly, the protected groups evaluated in this  
621 study simplified to binary labels, which excludes  
622 people who do not fall into this category such as  
623 non-binary individuals and the multidimensional  
624 nature of ethnicity.

625 Our reward scaling approach also requires these  
626 labels for classification. Although our setup could  
627 easily be extended to cases with more labels, it  
628 would be interesting to see fair classification with  
629 protected attributes as continuous values. But due  
630 to lack of good benchmarks restricts the evaluation  
631 of such cases.

632 **Algorithmic limitation** Firstly, our paper used  
633 two deep RL MDP algorithms and one linear clas-  
634 sical CMAB agent. We recognize that while linear  
635 agents have a significant focus in the CMAB lit-  
636 erature, the fast field of CMAB agents includes  
637 options with non-linear algorithms that could also  
638 be applied to this task. The choice of LinUCB does  
639 not represent the state-of-the-art, but rather a clas-  
640 sical high-performance implementation.

641 Second, the various hyperparameters limit the ex-  
642 tent of general statements about each algorithm.  
643 We have documented our hyperparameter search  
644 and training methods in the appendix, to ensure the  
645 interpretability of our experiments, but our results  
646 only demonstrate the capabilities of our best im-  
647 plementation. Moreover, the use of DTO to select  
648 the best model throughout training fails to account  
649 for potential trade-offs between fairness and accu-  
650 racy at different points in training. For example,  
651 on the Emoji dataset, PPO underperformed in Fair-  
652 ness and DQN in accuracy. However, it is possible  
653 that at another pointing training with a higher DTO  
654 score, the trade-off between fairness and accuracy  
655 was reversed.

## 656 Ethics Statement

657 The application of the paper was to improve fair-  
658 ness among protected groups in classification.

659 However, no algorithm is able to obtain perfect  
660 fairness and remove the bias perfectly. Therefore  
661 applications of the mentioned debiasing methods  
662 should always strongly take the mentioned limita-  
663 tions into account. Moreover, the current experi-  
664 ments are limited to specific datasets and real world  
665 use cases may be different. Careful evaluation and  
666 testing system behavior in the intended setting with  
667 input from experts who can judge the consequences  
668 of remaining bias is essential.

## 669 References

- 670 Carlos Aguirre, Kuleen Sasse, Isabel Cachola, and Mark  
671 Dredze. 2023. Selecting shots for demographic fair-  
672 ness in few-shot learning with large language models.  
673 *arXiv preprint arXiv:2311.08472*. 673
- 674 Richard Bellman. 1957. *Dynamic Programming*, 1  
675 edition. Princeton University Press, Princeton, NJ,  
676 USA. 676
- 677 Su Lin Blodgett, Lisa Green, and Brendan O’Connor.  
678 2016. Demographic Dialectal Variation in Social  
679 Media: A Case Study of African-American English.  
680 In *Proceedings of the 2016 Conference on Empiri-  
681 cal Methods in Natural Language Processing*, pages  
682 1119–1130. 682
- 683 Maria De-Arteaga, Alexey Romanov, Hanna Wal-  
684 lach, Jennifer Chayes, Christian Borgs, Alexandra  
685 Chouldechova, Sahin Geyik, Krishnamurthy Suresh,  
686 and Adam Tauman Kalai. 2019. Bias in bios: A case  
687 study of semantic representation bias in a high-stakes  
688 setting. In *proceedings of the Conference on Fairness,  
689 Accountability, and Transparency*, pages 120–128. 689
- 690 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and  
691 Kristina Toutanova. 2019. Bert: Pre-training of deep  
692 bidirectional transformers for language understand-  
693 ing. In *Proceedings of NAACL-HLT*, pages 4171–  
694 4186. 694
- 695 Maria Dimakopoulou, Zhengyuan Zhou, Susan Athey,  
696 and Guido Imbens. 2019. Balanced linear contextual  
697 bandits. In *Proceedings of the AAAI Conference on  
698 Artificial Intelligence*, volume 33, pages 3445–3453. 698
- 699 Miroslav Dudík, Dumitru Erhan, John Langford, and  
700 Lihong Li. 2014. Doubly robust policy evaluation  
701 and optimization. *Statistical Science*, 29(4):485–511. 701
- 702 Yanai Elazar and Yoav Goldberg. 2018. Adversarial  
703 removal of demographic attributes from text data. In  
704 *Proceedings of the 2018 Conference on Empirical  
705 Methods in Natural Language Processing*, pages 11–  
706 21. 706
- 707 Bjarke Felbo, Alan Mislove, Anders Søgaard, Iyad Rah-  
708 wan, and Sune Lehmann. 2017. Using millions of  
709 emoji occurrences to learn any-domain representa-  
710 tions for detecting sentiment, emotion and sarcasm. 710

711	In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 1615–1625.	765
712		766
713		767
714	Pantea Haghighatkhah, Antske Fokkens, Pia Sommerauer, Bettina Speckmann, and Kevin Verbeek. 2022. Better hit the nail on the head than beat around the bush: Removing protected attributes with a single projection. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 8395–8416.	768
715		769
716		770
717		771
718		772
719		773
720		774
721	Xudong Han, Timothy Baldwin, and Trevor Cohn. 2021. Diverse adversaries for mitigating bias in training. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 2760–2765.	775
722		776
723		777
724		778
725		779
726	Xudong Han, Timothy Baldwin, and Trevor Cohn. 2022a. Balancing out bias: Achieving fairness through balanced training. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 11335–11350.	780
727		781
728		782
729		783
730		784
731	Xudong Han, Aili Shen, Yitong Li, Lea Frermann, Timothy Baldwin, and Trevor Cohn. 2022b. Fairlib: A unified framework for assessing and improving fairness. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 60–71.	785
732		786
733		787
734		788
735		789
736		790
737	Michael Höfler, Hildegard Pfister, Roselind Lieb, and Hans-Ulrich Wittchen. 2005. The use of weights to account for non-response and drop-out. <i>Social psychiatry and psychiatric epidemiology</i> , 40:291–299.	791
738		792
739		793
740		794
741	Preethi Lahoti, Alex Beutel, Jilin Chen, Kang Lee, Flavien Prost, Nithum Thain, Xuezhi Wang, and Ed Chi. 2020. Fairness without demographics through adversarially reweighted learning. <i>Advances in neural information processing systems</i> , 33:728–740.	795
742		796
743		797
744		798
745		799
746		800
747	Lihong Li, Wei Chu, John Langford, and Robert E Schapire. 2010. A contextual-bandit approach to personalized news article recommendation. In <i>Proceedings of the 19th international conference on World wide web</i> , pages 661–670.	801
748		802
749		803
750		804
751		805
752	Enlu Lin, Qiong Chen, and Xiaoming Qi. 2020. Deep reinforcement learning for imbalanced classification. <i>Applied Intelligence</i> , 50:2488–2502.	806
753		807
754		808
755	Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A survey on bias and fairness in machine learning. <i>ACM computing surveys (CSUR)</i> , 54(6):1–35.	809
756		810
757		811
758		812
759	Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. 2015. Human-level control through deep reinforcement learning. <i>nature</i> , 518(7540):529–533.	813
760		814
761		815
762		816
763		
764		
	Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. 2020. Null it out: Guarding protected attributes by iterative nullspace projection. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7237–7256.	
	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> .	
	Aili Shen, Xudong Han, Trevor Cohn, Timothy Baldwin, and Lea Frermann. 2022. Does representational fairness imply empirical fairness? In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2022</i> , pages 81–95.	
	Jack Sherman and Winifred J Morrison. 1950. Adjustment of an inverse matrix corresponding to a change in one element of a given matrix. <i>The Annals of Mathematical Statistics</i> , 21(1):124–127.	
	Shivashankar Subramanian, Afshin Rahimi, Timothy Baldwin, Trevor Cohn, and Lea Frermann. 2021. Fairness-aware class imbalanced learning. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 2045–2051.	
	Victor Uc-Cetina, Nicolas Navarro-Guerrero, Anabel Martin-Gonzalez, Cornelius Weber, and Stefan Wermter. 2023. Survey on reinforcement learning for language processing. <i>Artificial Intelligence Review</i> , 56(2):1543–1575.	
	Christina Wadsworth, Francesca Vera, and Chris Piech. 2018. Achieving fairness through adversarial learning: an application to recidivism prediction. In <i>Proceedings of the FAT/ML Workshop</i> .	
	Tianlu Wang, Jieyu Zhao, Mark Yatskar, Kai-Wei Chang, and Vicente Ordonez. 2019. Balanced datasets are not enough: Estimating and mitigating gender bias in deep image representations. In <i>Proceedings of the IEEE/CVF international conference on computer vision</i> , pages 5310–5319.	
	Marco A Wiering, Hado Van Hasselt, Auke-Dirk Pietersma, and Lambert Schomaker. 2011. Reinforcement learning algorithms for solving classification problems. In <i>2011 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning (ADPRL)</i> , pages 91–96. IEEE.	
	<b>A Reproducibility</b>	
	<b>A.1 Data Analysis</b>	
	Because the BiasBios dataset needs to be scraped online, we provide the full composition of the BiasBios dataset split up in profession and gender in Table 5.	

Profession	Female	Male
Professor	53290	64820
Physician	19579	18986
Attorney	12494	20113
Photographer	8689	15635
Journalist	9873	10077
Nurse	17236	1735
Psychologist	11385	6910
Teacher	9768	6428
Dentist	5153	9326
Surgeon	1972	11301
Architect	2398	7715
Painter	3543	4193
Model	6214	1288
Poet	3441	3570
Filmmaker	2310	4699
Software Engineer	1089	5817
Accountant	2081	3571
Composer	918	4682
Dietitian	3689	289
Comedian	592	2207
Chiropractor	690	1908
Pastor	609	1923
Paralegal	1499	268
Yoga Teacher	1406	257
Dj	211	1274
Interior Designer	1183	280
Personal Trainer	654	778
Rapper	136	1271
rapper	136	1271

Table 5: Class and gender composition of the BiasBios dataset

## A.2 Model Selection

Selecting the best model throughout training or across hyperparameters is strongly dependent on the selection metric. To balance fairness and performance we use the proposed method of Han et al. (2022a), and select using DTO. The full equation of DTO is provided below, where the obtained metrics are determined by the point  $(Acc, (1 - GAP))$ , and the utopian metrics are  $(Acc^{utop}, (1 - GAP^{utop}))$ .

$$DTO = \sqrt{(Acc^{utop} - Acc)^2 + ((1 - GAP^{utop}) - (1 - GAP))^2}$$

The best training timestep according to DTO is determined with utopian values (1,1), and the best hyperparameters setting utopian values as the best metric values during training (i.e. the highest performance and fairness each individually obtained,

which do not necessarily belonging to the same algorithm).

The reported DTO values in table X and Y are obtained using the best performance and accuracy method as: [performance, fairness] BiasBios 28C = [0.811, 0.929], BiasBios 8C = [0.868, 0.978], Moji=[0.756, 0.900]

## A.3 Hyperparameters

The architecture of the neural network for each algorithm is fixed and consists of 2 layers MLP. For the critic in PPO the architecture is the same except for the output size which is 1. Hyperparameter optimization is applied for each of the parameters of the algorithms using grid search. Table 7 shows the ranges and the best values.

	Type	Dimensions
Layer 1	Linear	$n\_features \times 128$
Layer 2	Linear	$128 \times n\_actions$
Activation	ReLU	
Optimizer	Adam	

Table 6: Neural Network Architecture

**Related work implementations** Following previous work (Ravfogel et al., 2020; Han et al., 2022b), we use INLP and MP in a post hoc manner to the features extracted from the last hidden layer of the supervised model and train a logistic classifier for the final classification. For our MP debiasing experiments in section 5.5 we use MP to debias the context vectors before training, instead of poshoc on the hidden layer of the trained network.

## B Algorithms

### B.1 Single-Step Markov Decision Process

To formalize how the policy-gradient methods such as PPO relate to the Contextual Multi-Armed Bandit framework, we define below the single-step Markov Decision Process. An MDP is defined by the tuple  $(S, A, P, R, \gamma)$ , and our single-step variant contains only two states  $S = \{s_1, s_2\}$ . The initial state is sampled each time from the environment and for our classification setup is part of the set of context embeddings,  $s_1 \in \{x_j\}$ . To ensure data samples are treated independently the second state is always the terminal state  $s_2 = s_{terminal}$ . The action space is equal to the number of classes:

Algorithm	Parameter	Min	Max	Best	
				BiosBias	Emoji
PPO	lr (actor)	$3.0 \times 10^{-4}$	$1.0 \times 10^{-6}$	$1.0 \times 10^{-4}$	$3.0 \times 10^{-5}$
	lr (critic)	$1.0 \times 10^{-3}$	$1.0 \times 10^{-5}$	$1.0 \times 10^{-3}$	$1.0 \times 10^{-4}$
	Batch size	64	512	512	512
	Entropy $c_2$	0.01	0.1	0.2	0.1
	$\epsilon$ -clip	0.05	0.3	0.1	0.3
Supervised	lr	$1.0 \times 10^{-3}$	$1.0 \times 10^{-6}$	$3.0 \times 10^{-4}$	$1.0 \times 10^{-3}$
	Batch size	64	512	128	512
DQN	lr	$3.0 \times 10^{-4}$	$1.0 \times 10^{-6}$	$3.0 \times 10^{-6}$	$3.0 \times 10^{-4}$
	Batch size	32	256	256	32
	Eps_end	0.001	0.1	0.1	0.01
	Eps decay	0.5	1.0	0.5	0.5
LinUCB	$\alpha$	0.1	3.0	1.5	2.5

Table 7: Hyperparameter ranges and best values for different algorithms. For PPO the "Entropy  $c_2$ " refers to the coefficient of the entropy in the loss.

$A = C = \{c_1, c_2, \dots, c_{28}\}$ . The reward function  $R$  is equal to that of the CMAB and is defined in section 3.1. Lastly, each trajectory is defined as  $\tau = \{s_1, a_1, s_{terminal}\}$  and both the transition probability,  $P$ , and the discount factor  $\gamma$  are irrelevant since each action results in the terminal state.

## B.2 LinUCB

The full algorithm of LinUCB from Li et al. (2010), used in the paper is shown in Algorithm

### Algorithm 1 LinUCB Algorithm

**Require:** Context features  $x_{t,a}$  for context at time  $t$  and arm  $a \in \mathcal{A}$ , exploration parameter  $\alpha$ .  
Initialize  $A_a$  and  $b_a$  for each arm  $a \in \mathcal{A}$   
**for** each sample  $t$  **do**  
  **for** each arm  $a$  **do**  
     $\hat{\theta}_{at} = A_{at}^{-1} b_{at}$   
     $p_{t,a} = \hat{\theta}_{at}^\top x_{t,a} + \alpha \sqrt{x_{t,a}^\top A_{at}^{-1} x_{t,a}}$   
  **end for**  
  Choose arm  $a_t = \arg \max_{a \in \mathcal{A}} (p_{t,a})$ , and observe real-valued payoff  $r_t$   
  Update  $A_{a_t} \leftarrow A_{a_t} + x_{t,a_t} x_{t,a_t}^\top$   
  Update  $b_{a_t} \leftarrow b_{a_t} + r_t x_{t,a_t}$   
**end for**

## B.3 Equal Opportunity Weights

Where Han et al. (2022a) used EO for supervised learning, their implementation achieved this objective by grouping the loss per class and then averaging

over them. In this section, we see how we can use this to obtain the weights for each data sample based on the class  $a$  and protected attribute  $g$ . For two protected groups  $g_1$  and  $g_2$  in class  $a$ , let  $C_1$  and  $C_2$  be the number of samples for  $g_1$  and  $g_2$ , and  $\mathcal{W}_1$  and  $\mathcal{W}_2$ , be the weights. To get a statement of the weights with EO for each sensitive state,  $(a, g)$ , we need two axioms.

**Axiom 1.** The weight scale ratio between the two protected groups of a class should be inversely proportional to their probability in the dataset:

$$\mathcal{W}_1 \cdot C_1 = \mathcal{W}_2 \cdot C_2$$

**Axiom 2:** To ensure fairness across classes, the average weight per profession should be a fixed value  $B$  so that:

$$\frac{1}{C_1 + C_2} (\mathcal{W}_1 \cdot C_1 + \mathcal{W}_2 \cdot C_2) = B$$

Combining these two axioms we obtain the formulation:

$$\mathcal{W}_2 = \frac{B (C_1 + C_2)}{2 C_2}$$

$$\mathcal{W}_2 = \frac{B}{2} \frac{1}{P(C_2)}$$

For the multi-class classification task the average reward scale,  $B$ , should be 1, and the probability is conditional on the class  $a$ , obtaining the final  $W_{EO}$  equation:

$$W_{EO}(g, y) = \frac{1}{2} \frac{1}{P(g|a)}$$

## C Ablation Experiments

Here we add our experiments that did not make the main paper.

### C.1 Analysis: Model and Data Efficiency

An important aspect for evaluation is related to the data and computational of each algorithm. For ease of comparison, all algorithms except LinUCB were trained for 10 epochs. However, DQN and PPO each reuse the seen data in a different way to deal with the data sparsity of standard RL settings. DQN is updated using a replay-buffer from which it samples a minibatch of  $N$  triplets  $(s, a, r)$  for each iteration. In contrast, PPO collects  $N$  samples during the observation phase after which it updates the model with this batch  $K_{epoch}$  number of times. Lastly, LinUCB achieves optimal results after 1 epoch but is constrained by the computations of its weight matrices, which require the inverse of a square matrix with dimension  $n_{features}$ . For computational efficiency, we use the Sherman–Morrison formula which updates the previous computed inverse with a rank one update (Sherman and Morrison, 1950)

The time complexities in Table 2, demonstrate that PPO is closest to supervised learning and that DQN takes significantly more time since it needs to sample from the buffer at each iteration. Notably, LinUCB is strongly dependent on the number of classes, reducing its relative efficiency from 32 to 3 times that of Supervised Learning. The bottleneck here is that it needs to compute an upper confidence bound for each class.

Another important feature is the sensitivity to hyperparameters. PPO and DQN are sensitive to several hyperparameters that determine the level of its exploration, such as DQN’s mini-batch size or exploration parameter, or PPO’s entropy and clipping coefficients. LinUCB is easiest to implement in this regard and does not require any neural network hyperparameters, but only one exploration parameter  $\alpha$ , see section B.2.

## D Full result for experiments

To distinguish the sensitivity of gender imbalance and data-sparsity we also run experiments with a subset of the data, following Aguirre et al. (2023), and select only the professions that have at least 1000 samples for both genders in the test set, resulting in 8 professions.

## D.1 BiasBios: training performance over time

In Reinforcement Learning literature it is common to provide the performance of an algorithm throughout training for evaluation. Therefore we provide the evaluation accuracy of our four algorithms in Figure 5

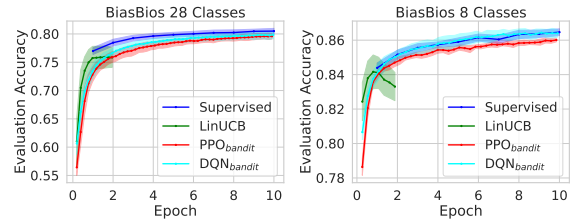


Figure 5: Evaluation accuracy of the different algorithms the full 28 classes and the 8 class subset of the Bias in Bios dataset

## D.2 BiasBios: Recall per profession

As a further analysis of the lacking F1 score of the RL algorithms compared to the supervised implementation, we provide the Recall scores as a percentage of the class. Since class 21, Professor appears significantly more often than the most common class after it, we leave it out for clarity.

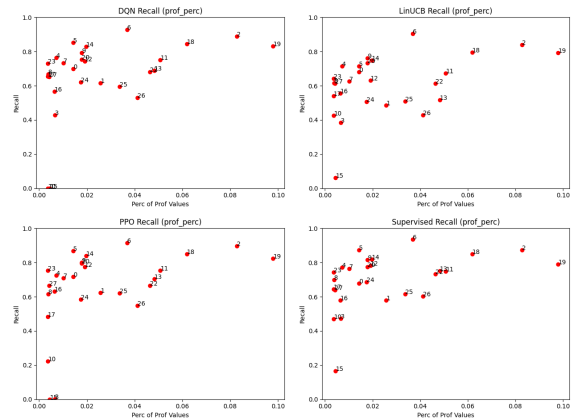


Figure 6: Recall of each class on the BiasBios dataset for the four algorithm implementations

## D.3 Full tables: BiasBios (28C and 8C)

Some of our results in section 5.3 are presented as the mean only. The full results of our algorithms as the mean and std over the five seeds is provided in the tables here. Table 8 shows the performance of our algorithms with and without reward scaling on the BiasBios dataset with the 28 and 8 classes.

Algorithm	28 Classes				8 Classes			
	Accuracy $\uparrow$	GAP $\downarrow$	DTO $\downarrow$	F1 $\uparrow$	Accuracy $\uparrow$	GAP $\downarrow$	DTO $\downarrow$	F1 $\uparrow$
Sup	<b>81.0 <math>\pm</math> 0.1</b>	16.4 $\pm$ 0.5	10.0	<b>73.8 <math>\pm</math> 0.3</b>	<b>86.8 <math>\pm</math> 0.1</b>	8.3 $\pm$ 0.7	6.2	<b>82.7 <math>\pm</math> 0.1</b>
LinUCB	78.4 $\pm$ 0.1	15.5 $\pm$ 0.3	9.6	67.3 $\pm$ 0.4	85.3 $\pm$ 0.2	<b>7.6 <math>\pm</math> 0.3</b>	5.8	80.6 $\pm$ 0.2
DQN <sub>bandit</sub>	80.1 $\pm$ 0.2	<b>13.7 <math>\pm</math> 0.3</b>	<b>7.2</b>	66.5 $\pm$ 1.3	86.5 $\pm$ 0.2	<b>7.6 <math>\pm</math> 0.3</b>	<b>5.5</b>	82.2 $\pm$ 0.2
PPO <sub>bandit</sub>	79.7 $\pm$ 0.5	14.4 $\pm$ 0.7	8.0	67.5 $\pm$ 2.0	86.0 $\pm$ 0.2	8.7 $\pm$ 0.4	6.7	81.6 $\pm$ 0.2
Sup <sup>EO</sup>	<b>80.1 <math>\pm</math> 0.2</b>	<b>7.1 <math>\pm</math> 0.5</b>	<b>1.1</b>	<b>71.7 <math>\pm</math> 0.5</b>	<b>86.3 <math>\pm</math> 0.2</b>	2.4 $\pm$ 0.1	<b>0.6</b>	<b>82.0 <math>\pm</math> 0.2</b>
LinUCB <sup>EO</sup>	74.6 $\pm$ 0.2	12.2 $\pm$ 0.5	9.6	59.8 $\pm$ 1.1	83.4 $\pm$ 0.2	7.6 $\pm$ 0.3	6.8	77.6 $\pm$ 0.3
DQN <sub>bandit</sub> <sup>EO</sup>	79.2 $\pm$ 0.1	10.1 $\pm$ 0.4	3.9	66.4 $\pm$ 0.2	86.2 $\pm$ 0.1	<b>2.2 <math>\pm</math> 0.2</b>	0.7	81.6 $\pm$ 0.2
PPO <sub>bandit</sub> <sup>EO</sup>	79.2 $\pm$ 0.2	8.5 $\pm$ 0.2	2.7	66.0 $\pm$ 0.8	85.8 $\pm$ 0.1	2.8 $\pm$ 0.6	1.3	81.4 $\pm$ 0.2

Table 8: Results on the BiasBios dataset for the full dataset (28 classes) and a subset of the most common professions (8 classes). The first rows use a constant reward scale, and the last four (in grey) use the EO reward scale

#### D.4 Full tables: four reward scaling methods

The results from reward scaling using the four described scales and our four algorithms are shown in Table 9.

Algo		Accuracy $\uparrow$	GAP $\downarrow$	F1
SUP	$\mathcal{W}_{\rho+}$	79.3 $\pm$ 0.1	7.9 $\pm$ 0.3	69.3 $\pm$ 0.3
	$\mathcal{W}_{\rho-}$	79.8 $\pm$ 0.3	6.9 $\pm$ 0.2	<b>71.8 <math>\pm</math> 0.6</b>
	$\mathcal{W}_{EO}$	<b>80.1 <math>\pm</math> 0.2</b>	7.1 $\pm$ 0.5	71.7 $\pm$ 0.5
	$\mathcal{W}_{IPW}$	72.1 $\pm$ 0.7	<b>6.1 <math>\pm</math> 0.3</b>	64.8 $\pm$ 0.8
PPO	$\mathcal{W}_{\rho+}$	74.6 $\pm$ 0.7	9.9 $\pm$ 0.8	49.7 $\pm$ 2.2
	$\mathcal{W}_{\rho-}$	78.8 $\pm$ 0.1	<b>8.4 <math>\pm</math> 0.6</b>	64.7 $\pm$ 0.8
	$\mathcal{W}_{EO}$	<b>79.2 <math>\pm</math> 0.2</b>	8.5 $\pm$ 0.2	<b>66.0 <math>\pm</math> 0.8</b>
	$\mathcal{W}_{IPW}$	45.8 $\pm$ 6.9	10.5 $\pm$ 0.9	45.3 $\pm$ 5.8
DQN	$\mathcal{W}_{\rho+}$	76.2 $\pm$ 1.1	10.4 $\pm$ 0.7	57.2 $\pm$ 4.8
	$\mathcal{W}_{\rho-}$	<b>79.3 <math>\pm</math> 0.1</b>	11.1 $\pm$ 0.6	65.8 $\pm$ 1.4
	$\mathcal{W}_{EO}$	79.2 $\pm$ 0.1	<b>10.1 <math>\pm</math> 0.4</b>	<b>66.4 <math>\pm</math> 0.2</b>
	$\mathcal{W}_{IPW}$	74.6 $\pm$ 0.3	12.8 $\pm$ 0.2	56.6 $\pm$ 0.3
LinUCB	$\mathcal{W}_{\rho+}$	72.8 $\pm$ 0.1	12.0 $\pm$ 0.5	54.6 $\pm$ 0.9
	$\mathcal{W}_{\rho-}$	74.1 $\pm$ 0.4	11.6 $\pm$ 0.5	59.3 $\pm$ 1.7
	$\mathcal{W}_{EO}$	<b>74.6 <math>\pm</math> 0.2</b>	12.2 $\pm$ 0.5	<b>59.8 <math>\pm</math> 1.1</b>
	$\mathcal{W}_{IPW}$	37.3 $\pm$ 2.5	<b>10.3 <math>\pm</math> 0.7</b>	35.4 $\pm$ 1.0

Table 9: Results with different reward scaling on BiasBios for various algorithms

#### D.5 Full results: Explicit gender information and Ensemble techniques

This section includes the full results of Section 5.5, after adding the gender information explicitly and after removing it with MP. The results are presented as mean and standard deviation over 5 seeds in Table 10 and Table 11

Algo + g	Accuracy $\uparrow$	GAP $\downarrow$	F1 $\uparrow$
Sup <sup>EO</sup>	<b>80.2 <math>\pm</math> 0.2</b>	<b>7.2 <math>\pm</math> 0.5</b>	<b>71.9 <math>\pm</math> 0.7</b>
LinUCB <sup>EO</sup>	74.5 $\pm$ 0.2	11.7 $\pm$ 0.5	59.6 $\pm$ 0.8
DQN <sub>bandit</sub> <sup>EO</sup>	79.2 $\pm$ 0.2	10.0 $\pm$ 0.5	66.1 $\pm$ 0.4
PPO <sub>bandit</sub> <sup>EO</sup>	79.3 $\pm$ 0.1	8.7 $\pm$ 0.3	66.1 $\pm$ 0.6

Table 10: Results on the BiasBios dataset with explicit gender information added to the context.

Algo + MP	Accuracy $\uparrow$	GAP $\downarrow$	F1 $\uparrow$
Sup <sup>EO</sup>	80.0 $\pm$ 0.2	7.4 $\pm$ 0.4	71.9 $\pm$ 0.3
LinUCB <sup>EO</sup>	74.3 $\pm$ 0.4	11.5 $\pm$ 0.1	59.4 $\pm$ 1.0
DQN <sub>bandit</sub> <sup>EO</sup>	79.0 $\pm$ 0.2	8.6 $\pm$ 0.3	65.8 $\pm$ 0.6
PPO <sub>bandit</sub> <sup>EO</sup>	79.2 $\pm$ 0.2	9.7 $\pm$ 0.6	66.8 $\pm$ 1.4

Table 11: Performance on the BiasBios dataset, using MP debiased embeddings