Balancing out Bias: Achieving Fairness Through Balanced Training

Anonymous ACL submission

Abstract

Bias in natural language processing manifests as disparities in error rates across author demographics, typically disadvantaging minority groups. Although dataset balancing has been shown to be effective in mitigating bias, existing approaches do not directly account for correlations between author demographics and linguistic variables. To achieve Equal Opportunity fairness, this paper introduces a simple but highly effective objective for countering bias using balanced training. We extend the method in the form of a gated model, which incorporates protected attributes as input, and show that it is effective at reducing bias in predictions through demographic input perturbation, outperforming all other bias mitigation techniques when combined with balanced training.

1 Introduction

001

003

007

800

014

017

019

021

029

034

040

Natural Language Processing (NLP) models have achieved extraordinary gains across a variety of tasks in recent years. However, naively-trained models often learn spurious correlations with demographics and socio-economic factors (Hendricks et al., 2018; Lu et al., 2018; Bolukbasi et al., 2016; Park et al., 2018), leading to disparities across author demographics in contexts including coreference resolution, sentiment analysis, and hate speech detection (Badjatiya et al., 2019; Zhao et al., 2018; Li et al., 2018a; Díaz et al., 2018).

Two popular approaches for mitigating such biases are: (1) balancing each demographic group in training, either explicitly via sampling (Zhao et al., 2018; Wang et al., 2019) or implicitly via balancing losses for each group (Höfler et al., 2005; Lahoti et al., 2020); and (2) removing demographic information from learned representations (Li et al., 2018a; Wang et al., 2019; Ravfogel et al., 2020; Han et al., 2021b).

While balancing methods have been shown to be successful, they have not been tested extensively

in NLP. In this paper, we focus on author bias and adapt three balanced training approaches for debiasing. In addition, we propose a new objective for balanced training, which can be used for proxy optimization of Equal Opportunity (Hardt et al., 2016). We first provide a theoretical justification for our approach, and then conduct experiments on two benchmark datasets which show that our proposed objective is highly effective in achieving Equal Opportunity fairness. 042

043

044

045

046

047

051

054

055

058

060

061

062

063

064

067

068

069

070

071

072

073

074

075

076

077

078

079

081

Even when the training data is balanced, ignoring demographic-specific features can lead to bias (Wang et al., 2019; Lahoti et al., 2020), due to differences in language use across demographics (Hovy, 2015). There is thus a fine line to be walked in terms of optimizing for linguistic variables associated with different demographic groups (potentially boosting overall model accuracy), and ensuring model fairness.

Inspired by work in domain adaptation on learning domain-specific representations that generalize across domains (Bousmalis et al., 2016; Li et al., 2018b), we propose a gated model, which incorporates author demographics as an *input* to generate group-specific representations but also generalizes across demographic groups. We show that when combined with instance reweighting during training, this technique leads to substantial bias reductions over leading debiasing techniques, typically with higher predictive accuracy. We also introduce a second means of bias reduction through tailoring gating coefficients of the trained model, which allows for fine-tuning of the accuracy-fairness tradeoff. Our experiments over two benchmark datasets for language debiasing show that our techniques are competitive with much more complex state-ofthe-art methods for debiasing in situations where the demographic attribute is not known at test time, and provide substantial gains over the state-of-theart when the protected attribute is observed. Codes will be released upon acceptance.

090

100

101

102

103

104

105

107

108

109

110

111

112

113

114 115

116

117

118 119

120

121

122

123

124

125

126

127

128

129

130

2 Balanced Training

Despite their simplicity and versatility, balanced training approaches have only received limited attention in prior work in NLP. In this section, we review three balanced training approaches, and discuss their objectives and applications. We further propose a novel objective for balanced training, which we will show to be a proxy optimization for the Equal Opportunity metric.

2.1 Problem Formulation

In this paper, we focus on bias mitigation for NLP classification tasks. Formally, we assume a dataset $\mathcal{D} = \{(x_i, y_i, g_i)\}_{i=1}^n$ where $x_i \in X$ is a *d*-dimensional input text representation vector, $y_i \in Y$ denotes the main task label (e.g. sentiment), and $g_i \in G$ represents the private attribute associated with x_i , e.g., author gender.

A standard model M is trained to predict Y given X, while debiasing methods generally aim to learn a model M' that is fair wrt G by considering $X \times G$ together.

2.2 Fairness Measurement

As in Barocas et al. (2019), the *separation* criterion acknowledges the correlation between G and Y, and is satisfied iff $G \perp \hat{Y}|Y$. A relaxation of the separation criterion known as *equality of opportunity* is widely used (Hardt et al., 2016; Ravfogel et al., 2020; Han et al., 2021a). *Equality of opportunity* measures the difference in true positive rate (TPR) across all groups, based on the notion that the positive outcome represents 'advantage', such as being accepted by a school or getting a loan. Essentially, the difference (gap) in TPR reflects whether different groups have equal opportunity.

2.3 Balanced Training Objectives

We now formally describe the objective functions of three established balanced training approaches, and discuss their applications.

Let \mathcal{X} be the task loss and n be the number of observed instances in the dataset \mathcal{D} . The overall empirical risk is written as $\mathcal{L} = \frac{1}{n} \sum_{i} \mathcal{X}(y_i, \hat{y}_i)$, which can be rewritten as the aggregation of subsets: $\mathcal{L} = \sum_{y} \sum_{g} \frac{n_{y,g}}{n} \mathcal{L}_{y,g}$, where $n_{y,g} := |\{i : y_i = y, g_i = g\}|$, the number of instances with target label y and demographic attribute g, and $\mathcal{L}_{y,g}$ is the empirical loss corresponding to the subset, $\mathcal{L}_{y,g} = \frac{1}{n_{y,g}} \sum_{i} \mathcal{X}(y_i, \hat{y}_i) \mathbb{1}(y_i = y, g_i = g)$. Furthermore, we use * as the notation for marginalization, for example, $n_{*,g} = \sum_{y} n_{y,g}$. Let p be the target objective, and \tilde{p} be the empirical probability based on the training dataset.

Given this notation, the three balanced training objectives are as follows:

Balanced Demographics Zhao et al. (2018) augment the dataset according to the demographic label distribution (making p(G) uniform) for in-textbias mitigation. Although their gender-swapping approach is not directly applicable to our tasks, we adapt the general objective function as $\mathcal{L}^G = \frac{1}{|G|} \sum_{y} \sum_{g} \frac{n_{y,g}}{n_{*,g}} \mathcal{L}_{y,g}$, where |G| is the number of distinct labels of G.

Conditionally Balance In a vision context, Wang et al. (2019) down-sample the majority demographic group within each class, so that on a per-class basis, it does not dominate the minority group (i.e. p(G|Y) is uniform for all Y), giving the objective function: $\mathcal{L}^{G|Y} = \frac{1}{|G|} \sum_{y} \frac{n_{y,*}}{n} \sum_{g} \mathcal{L}_{y,g}$.

Jointly Balance Lahoti et al. (2020) employ instance reweighting for structural data classification such that demographics and classes are jointly balanced, leading to the objective: $\mathcal{L}^{G,Y} = \frac{1}{|G| \times |Y|} \sum_{y} \sum_{g} \mathcal{L}_{y,g}$.

2.4 Achieving the objective

In this paper, we focus on two classic ways of achieving the target objective: instance reweighting, which manipulates the weight of each instance during training, and down-sampling, which preprocess the dataset before training.

Taking the objective $\mathcal{L} = \frac{1}{|G| \times |Y|} \sum_{y} \sum_{g} \mathcal{L}_{y,g}$ as an example, instance reweighting reweights each instance inversely proportional to the frequency of the combination of its main label and demographic label,

$$\frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i, g_i) \in \mathcal{D}} \tilde{p}^{-1} (G = g_i, Y = y_i) \mathcal{X}(y_i, \hat{y}_i),$$

where \mathcal{X} is the task loss, and \hat{y}_i denotes the model prediction given input text x_i .

The other approach, down-sampling, subsamples non-minority instances to derive a balanced training dataset, such that $\tilde{p}(g, y) = \frac{1}{|G| \times |Y|}, \forall g \in G, y \in Y$. Specifically, let $\mathcal{D}_{y,g}$ denote a subset of training instances s.t. $\mathcal{D}_{y,g} = \{(x_i, y_i, g_i) | y_i = y, g_i = g\}_{i=1}^n$. We sample without replacement to get a target subset $\mathcal{D}_{y,g}^*$ such that $|\mathcal{D}_{y,g}^*| = \min\{|\mathcal{D}_{y,g}|, \forall y \in Y, g \in G\}$. The

131

132

133

134

135

136

137

138

139

140

151 152 153

154

155

156

158

159

160

161

163

164

166

167

170

sampled subsets are merged to form the trainingset.

173 2.5 Towards equal opportunity

185

188

190

191

192

193

194

195

196

197

198

Without loss of generality, we illustrate with the 174 binary case of $y \in \{T, F\}$ and $g \in \{0, 1\}$. Re-175 call that the equal opportunity metric is satisfied 176 if a binary classification model has an equal positive prediction rate conditioned on the advantaged 178 class. Assuming the advantaged class is denoted as 179 y = T, i.e. the positive class, the equal opportunity 180 is measured by the TPR GAP between protected groups. Our proposed objective function for equal 183 opportunity is:

$$\mathcal{L} = \frac{n_{T,*}}{n} \frac{1}{2} \sum_{g \in \{0,1\}} \mathcal{L}_{T,g} + \sum_{g \in \{0,1\}} \frac{n_{F,g}}{n} \mathcal{L}_{F,g}$$
$$= \sum_{g \in \{0,1\}} \frac{n_{T,g}}{n} \frac{n_{T,g}}{2n_{T,g}} \mathcal{L}_{T,g} + \sum_{g \in \{0,1\}} \frac{n_{F,g}}{n} \mathcal{L}_{F,g}$$

Compared to the vanilla objective, the weights of instances with T target label are adjusted. Specifically, the reweighting term $\frac{n_{T,*}}{2n_{T,g}}$ is larger than 1 for the minority group, and less than 1 for the majority group.

From CE to TPR Cross-entropy is an estimate of the TPR at the mini-batch level when considering a subset of instances with the same target label. Recall that the CE loss for binary classification, of an instance is $-[y_i \cdot \log(\hat{p}(y_i)) + (1 - y_i) \cdot \log(1 - \hat{p}(y_i))]$, where $\hat{p}(y_i)$ is the predicted probability of y_i being True. Taking y = T for a certain demographic group g as an example,

$$\mathcal{L}_{T,g} = \frac{1}{n_{T,g}} \sum_{i} \mathcal{X}(y_i, \hat{y}_i) \mathbb{1}(y_i = T, g_i = g)$$
$$= -\frac{1}{n_{T,g}} \sum_{i} \hat{p}(y_i) \mathbb{1}(y_i = T, g_i = g).$$

Essentially, minimizing $\mathcal{L}_{T,g}$ is equivalent to maximizing the predicted probability of \hat{y} being True given target label y is True, within the demographic group g. That is, at the minibatch level, $-\mathcal{L}_{T,g}$ is an estimator of $p(\hat{y} = T|y = T, g = g)$, which is the TPR of group g. Given this, our proposed objective minimizes the TPR GAP by focusing on the TPR of the different demographic groups.

207Beyond binary labels&demographic attributes208Although we have introduced both target labels and209demographic attributes to be binary, our proposed



Figure 1: Gated model architecture. Given the input vector x, e.g. a text representation, the model has a shared encoder component and |G| encoder components, one for each demographic group.

objective generalizes trivially. The equal opportunity metric was originally designed for binary classification, under the assumption of a single advantaged class y = T. To satisfy the multi-class target label case, we adjust the equal opportunity to consider the one-vs-all setting, and measuring the TPR of each target class. Our proposed objective then becomes $\sum_{y} \sum_{g} \frac{n_{y,g}}{n} \frac{n_{y,*}}{|G| \times n_{y,g}} \mathcal{L}_{y,g}$. This recovers the formulation of Conditionally Balanced of section 2.3.

3 Demographic Factors Improve Fairness

Ignoring demographic-specific features can lead to bias even when the training data has been balanced (Wang et al., 2019; Lahoti et al., 2020). Our approach to dealing with this is, rather than removing demographic information from representations, to use a gated model that uses demographic labels as input.

As can be seen in Figure 1, the gated model consists of (1+|G|) encoders: one shared encoder, and a dedicated encoder for each demographic group in G.¹ Formally, let E denote the shared encoder, E_j denote the encoder for the j-th demographic group, C denote the classifier, and g_i be a 1-hot input such that $g_{i,j}$ is 1 if the instance (x_i, g_i, y_i) belongs to the j-th group, and 0 otherwise. The prediction for an instance is: $\hat{y}_i = C(h_i^s, h_i^g)$, where $h_i^s = E(x_i)$ and $h_i^g = \sum_{j=1}^{|G|} g_{i,j} E_j(x_i)$. The two inputs are concatenated and input to the classifier C.

Intuitively, the shared encoder learns a general representation, while each group-specific encoder captures group-specific representations during training.

3

¹Strictly speaking, it is possible to achieve a similar effect with |G| encoders by merging one group with the shared encoder, and using post-hoc correction to separate out the general from the group-specific representation (Kang et al., 2020).

Our setting differs from other debiasing methods in that we assume the demographic attribute is available at training and prediction time, while techniques such as adversarial training (Li et al., 2018a) and INLP (Ravfogel et al., 2020) only require the attribute for training. This richer input allows for more accurate predictions, courtesy of the demographic-specific encoder, but limits applicability at test time. As suggested by Hovy and Yang (2021), addressing demographic factors is essential for NLP to get closer to the goal of human-like language understanding, and increase fairness. For better applicability, we also relax this requirement by replacing demographic factors with non-informative prior in section 4.8.

4 Experimental Results

4.1 Evaluation Metrics

245

246

247

248

257

258

261

262

265

266

267

272

273

277

278

281

290

291

Following Ravfogel et al. (2020), we use overall accuracy as the performance metric, and the separation criterion to measure fairness in the form of TPR GAP and TNR GAP: the true positive rate and true negative rate differences between demographic groups. For both GAP metrics, smaller is better, and a perfectly fair model will achieve 0. For multi-class classification tasks, we follow Ravfogel et al. (2020) in reporting the quadratic mean (RMS) of TPR GAP over all classes. In a binary classification setup, TPR and TNR are equivalent to the TPR of the positive and negative classes, respectively, so we employ the RMS TPR GAP in this case also.

Throughout this paper, we report accuracy and GAP results as mean values \pm standard deviation over the test set, averaged across five independent runs with different random seeds.

For ease of comparison between approaches, we introduce 'distance to the optimum' (DTO), a single metric to incorporate accuracy and GAP into a single figure of merit, which is calculated by: (1) converting GAP to 1 - GAP (denoted as fairness; higher is better); (2) normalizing each of accuracy and fairness, by dividing by the best result for the given dataset (i.e., highest accuracy and fairness); and (3) calculating the Euclidean distance to the point (1, 1), which represents the hypothetical system which achieves highest accuracy and fairness for the dataset. Lower is better for this statistic, with minimum 0.

In addition to performance and fairness, we are also interested in the efficiency of the different debiasing approaches and report each method's average training time.² We present normalized training times relative to the standard method, i.e., the average training time divided by that of the standard model.

294

295

296

297

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

4.2 Dataset

Following Ravfogel et al. (2020), we conduct experiments over two NLP classification tasks — sentiment analysis and biography classification — using the same dataset splits as prior work.

4.2.1 Moji

This sentiment analysis dataset was collected by Blodgett et al. (2016), and contains tweets that are either African American English (AAE)-like or Standard American English (SAE)-like. Each tweet is annotated with a binary 'race' label (based on language use: either AAE or SAE), and a binary sentiment score determined by (redacted) emoji contained in it. We use the train, dev, and test splits from Han et al. (2021b) of 100k/8k/8k instances, respectively.

4.2.2 BIOS

The second task is biography classification (De-Arteaga et al., 2019; Ravfogel et al., 2020), where biographies were scraped from the web, and annotated for binary gender and 28 classes of profession. Since the data is not directly available, in order to construct the dataset, we use the scraping scripts of Ravfogel et al. (2020), leading to a dataset with 396k biographies.³ Following Ravfogel et al. (2020), we randomly split the dataset into train (65%), dev (10%), and test (25%).

4.3 Models

We first implement a "STANDARD" model on each dataset, without explicit debiasing. On the MOJI dataset, we follow Ravfogel et al. (2020); Han et al. (2021b) in using DeepMoji (Felbo et al., 2017) as the encoder to get 2304d representations of input texts. The DeepMoji model contains 22.4 million parameters and was pretrained over 1246 million tweets to predict one of 64 common emojis. Ravfogel et al. (2020) and Subramanian et al. (2021) used uncased BERT-base (Devlin et al., 2019) as their

 $^{^{2}}$ Testing on Titan X and RTX 3090, all models have roughly near-identical inference time.

³There are slight discrepancies in the dataset composition due to data attrition: the original dataset (De-Arteaga et al., 2019) had 399k instances, while 393k were collected by Ravfogel et al. (2020).

Model	Accuracy↑	$\mathbf{GAP}\downarrow$	$\textbf{DTO}\downarrow$
STANDARD	82.3 ± 0.0	16.0 ± 0.5	0.093
BD (Zhao et al., 2018)	82.3 ± 0.0	15.6 ± 0.2	0.089
JB (Lahoti et al., 2020)	74.7 ± 0.3	7.4 ± 0.3	0.092
EO	79.4 ± 0.1	9.7 ± 0.6	0.043

Table 1: Results for balanced training methods (BT) on the BIOS test set. **EO**: our proposed objective in section 2.5. **BD** and **JB** refer to baselines, balanced demographics and jointly balance in section 2.3.

STANDARD model for the BIOS dataset, taking the 'CLS' token as the source of a fixed text representation, without further fine-tuning. However, we found that taking the average of all contextualized token embeddings led to an accuracy improvement of 1.4% and GAP fairness improvement of 2.4%. Given this, we use 768d 'AVG' representations extracted from the pretrained uncased BERT-base model.

337

338

339

341

343

344

351

354

357

367

371

373

For INLP (Ravfogel et al., 2020), we take the fixed STANDARD model for the given dataset, and iteratively train a linear classifier and perform nullspace projection over the learned representation. For the other baseline models — ADV and DADV— we jointly train the adversarial discriminators and classifier. In order to ensure a fair comparison, we follow Han et al. (2021a) in using a model consisting of the same fixed-parameter encoder as ours followed by a trainable 3-layer MLP.

4.4 Balanced Training Approaches

Since the MOJI dataset has been artificially balanced for main task labels and demographic labels, balanced training corresponding to p(g) makes no difference, and moreover, the results for p(g|y) and p(g, y) will be identical. Given this, we focus on the BIOS dataset for comparing different balanced training objectives.⁴

Table 1 shows the results of balanced training using the different objectives. Compared to the STANDARD model, balanced training with different objectives are all able to reduce bias, and the objective proposed by Lahoti et al. (2020) achieves the best TPR GAP. However, in terms of accuracy– fairness trade-off, our proposed approach outperforms all other models, which is not surprising as our proposed objective is designed to achieve better equal opportunity fairness. Based on these results, hereafter, we only report balanced training with our proposed objective.

4.5 Main Results

We report results over the sentiment analysis and biography classification tasks in Table 2. The baseline models are: **STANDARD**, which is a naively-trained MLP classifier; **INLP** (Ravfogel et al., 2020), which removes demographic information from text representations through iterative nullspace projection; **ADV** (Li et al., 2018a; Wang et al., 2019), which performs protected information removal through adversarial training with a single discriminator; and **DADV** (Han et al., 2021b), which also uses adversarial training but with multiple adversaries subject to an orthogonality constraint, and represents the current state-of-the-art (SOTA).

On the MOJI dataset, compared to the STAN-DARD model, BT simultaneously increases main task accuracy and mitigates bias, leading to results competitive with ADV and better than INLP. Although BT does not outperform the SOTA DADV, it leads to performance–fairness trade-offs that are competitive with the other debiasing methods.

On the BIOS dataset, BT again leads to performance–fairness trade-offs that outperform the baseline methods. However, different to the MOJI dataset, BT does not further improve accuracy, increasing fairness by 5.3% absolute at the cost of 2.9% accuracy.

In terms of training time, existing debiasing methods (esp. DADV on MOJI) incur a substantial overhead, while balanced training is much more frugal: around 1.3 times faster (because of the reduction in training data volume).

In addition to evaluating BT, we also combine GATE with BT, which achieves a better performance–fairness balance, as shown in Table 2. This is consistent with our argument that, rather than removing demographic information, properly used demographic factors can further reduce biases. Indeed, the BT +GATE consistently outperforms the current SOTA model DADV on both datasets.

Combining balanced training with benchmark methods The baseline methods INLP and DADV as presented above were used in a manner consistent with their original formulation, i.e., without balanced training. An important question is whether balanced training might also benefit these methods. It is trivial to combine

423

374

⁴As BIOS is a multi-class classification task and our proposed approach generalizes to p(g|y) in this case, there is no need to include Wang et al. (2019) in our comparison.

		Мојі			BIOS				
Method	Model	Accuracy↑	$\mathbf{GAP}\downarrow$	$\textbf{DTO}\downarrow$	Time↓	Accuracy↑	$\mathbf{GAP}\downarrow$	$\mathbf{DTO}\downarrow$	Time↓
Baselines	STANDARD	71.6 ± 0.1	31.0 ± 0.3	0.261	1.0	82.3 ± 0.0	16.0 ± 0.5	0.110	1.0
	INLP	68.5 ± 1.1	33.8 ± 3.9	0.300	14.0	70.5 ± 0.5	6.7 ± 0.9	0.145	6.3
	Adv	74.3 ± 0.4	22.2 ± 3.7	0.163	36.1	81.1 ± 0.1	12.7 ± 0.3	0.077	1.3
	DADV	74.5 ± 0.3	18.5 ± 2.0	0.123	109.4	81.1 ± 0.1	12.6 ± 0.3	0.076	2.4
Ours	BT	74.0 ± 0.2	21.5 ± 0.4	0.155	0.8	79.4 ± 0.1	9.7 ± 0.6	0.057	0.7
	+Gate $*$	74.9 ± 0.2	13.8 ± 0.3	0.072	0.8	79.4 ± 0.1	9.2 ± 0.2	0.053	0.7
Combination	+DADV	72.2 ± 0.2	14.3 ± 0.2	0.085	90.4	79.3 ± 0.1	9.9 ± 0.2	0.059	2.7
	+INLP	72.3 ± 1.9	15.7 ± 3.1	0.099	8.9	73.6 ± 0.6	5.6 ± 0.7	0.107	3.8

Table 2: Results over sentiment analysis (MOJI) and biography classification (BIOS) tasks. DTO are measured by the normalized Euclidean distance between each model and the ideal model, and lower is better. **Bold** = best trade-off within category. Training time is reported relative to STANDARD, which takes 35 secs and 16 mins for MOJI and BIOS, respectively. *: requires demographic attribute at test time.

Model	Size	Accuracy↑	$\mathbf{GAP}\downarrow$	$\mathbf{DTO}\downarrow$
STANDARD	257k	82.3 ± 0.0	16.0 ± 0.5	0.093
$\mathbf{RW} + p(g)$	257k	82.3 ± 0.0	15.6 ± 0.2	0.089
$\mathbf{RW} + p(g y)$	257k	75.7 ± 0.2	13.9 ± 0.4	0.107
$\mathbf{RW} + p(g, y)$	257k	74.7 ± 0.3	7.4 ± 0.3	0.092
DS + p(g)	237k	82.1 ± 0.1	15.9 ± 0.3	0.092
DS + p(g y)	37k	79.4 ± 0.1	9.7 ± 0.6	0.043
DS + $p(g, y)$	5k	66.1 ± 0.1	10.9 ± 0.4	0.200

Table 3: Results for balanced training methods on the BIOS test set. "RW" = balancing through instance reweighting; "DS" = balancing through dataset down-sampling; and "Size" = the number of instances in the training dataset.

downsampling with INLP and DADV, as the method simply prunes the training dataset, but does not impact the training objective. To combine instance reweighting with DADV, we modify the training objective such that the cross-entropy term is scaled by \tilde{p}^{-1} , while leaving the adversarial term unmodified, i.e., solve for $\min_M \max_A \sum_{(x_i,y_i,g_i)\in D} \tilde{p}^{-1} \mathcal{X}(y_i, \hat{y}_i) - \lambda_{adv} \mathcal{X}(g, \hat{g})$. For INLP, we simply train a BT model, and then iteratively perform INLP linear model training and nullspace projection over the learned representations.

494

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

Results are presented in the final section of Table 2 ("Combination"), and indicate that the combined methods appreciably outperform both the standalone demographic removal methods and balanced training approaches, without extra training time cost. That is, demographic information removal and balanced training appear to be complementary.

4.6 Reweighting vs. Down-sampling

Table 3 shows the results of the naively-trained MLP model ("STANDARD") and six balancedtraining methods, all based on the same MLP model architecture as STANDARD. Corpus downsampling ("DS") removes instances from majority groups and thus leads to less training data and overall lower accuracy than instance reweighting ("RW"). 444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

When using p(g) as the objective, both RW and DS perform similarly to the STANDARD model, as the overall gender distribution is quite balanced, which can also be seen in the size of the training data for DS + p(g). Both RW + p(g,y) and RW + p(g|y) reduce bias and performance, but RW + p(g,y) outperforms RW + p(g|y) in terms of the performance–fairness trade-off, in that RW + p(g,y) achieves similar performance but substantially better fairness (6.6% absolute improvement in GAP). However, p(g,y) is not as effective as p(g|y) when combined with DS, due to the big drop in the volume of training data.

4.7 Gated Model

If the training dataset is imbalanced and contains spurious correlations between task labels and demographic attributes, a naively trained model will learn and possibly amplify dataset biases. The gated model, with its explicit conditioning and group-specific encoding, will be particularly vulnerable to bias.

Table 4 shows that, on both datasets, the GATE model increases the accuracy but amplifies bias (e.g., GAP of 65 on MOJI): as it uses demographic information directly to make predictions, it is highly vulnerable to bias in the training dataset.

		Мојі				
Model	Accuracy↑	$\mathbf{GAP}\downarrow$	DTO \downarrow	Accuracy↑	$\mathbf{GAP}\downarrow$	DTO \downarrow
BT	74.0 ± 0.2	21.5 ± 0.4	0.155	79.4 ± 0.1	9.7 ± 0.6	0.057
+GATE $*$	74.9 ± 0.2	13.8 ± 0.3	0.072	79.4 ± 0.1	9.2 ± 0.2	0.053
GATE *	64.8 ± 0.1	65.2 ± 0.9	0.640	82.4 ± 0.1	19.2 ± 0.3	0.144
GATE $_{0.5}^{\text{soft}}$	72.7 ± 0.2	30.2 ± 0.3	0.250	80.8 ± 0.1	11.6 ± 0.3	0.066
GATE Soft *	74.8 ± 0.2	20.3 ± 0.3	0.142	81.1 ± 0.1	19.8 ± 0.4	0.151
Gate $_{RMS}^{soft}$ *	73.5 ± 0.2	7.1 ± 0.3	0.019	80.5 ± 0.1	11.1 ± 0.3	0.063

Table 4: Results over MOJI and BIOS. *: demographic attributes available at test time.

Intuitively, the only objective of GATE training is standard cross-entropy loss, which has been shown to lead to bias amplification under imbalanced training without regularisation. The gate components explicitly rely on demographic information, and thus become a strong indicator of main label predictions due to spurious correlations between the main task label and demographic labels in the training set. Balanced training approaches act as regularizers in preventing the model from learning and amplifying spurious correlations in training.

4.8 Soft Averaging

Although the gated model naturally requires the demographic attribute at test time, we also evaluate a condition where this is not available.

GATE ^{soft} Instead, we take a Bayesian approach by evaluating $p(y|x) = \sum_g p(g)p(y|g,x)$, where we can control the prior explicitly. For example, under a uniform demographic attribute prior, we simply average the predictions p(y|x,g) and $p(y|x,\neg g)$. This Bayesian approach can be approximated by soft averaging, whereby the activation of all demographic-specific encoders are uniformly averaged inside the model, i.e., $g_{i,j} = \frac{1}{|G|}$, rather than selecting only one in the standard gated model (i.e., $g_{i,:}$ is 1-hot).⁵

GATE soft Acc and GATE soft RMS When the protected attribute is observed at test time the soft averaging method may still prove useful, which we use as a means for fine-tuning the balance between accuracy and bias. Specifically, we consider non-uniform encoder averaging conditioned on the gold protected attribute, g^* . Let α and β denote to what extend the 1-hot labels are manipulated according to the value of g^* as 0 and 1 respectively, the soft labels are $\left[\alpha \quad 1-\alpha\right]$ and $\left[1-\beta \quad \beta\right]$. I.e., the

two specific encoders are weighted by either α and $1 - \alpha$, or $1 - \beta$ and β , respectively, according to the value of g^* . Values of $\alpha, \beta < 0.5$ mean the protected label is (softly) preserved, while values > 0.5 mean the label is flipped. In cases where the model is biased towards or against a demographic group, it may be advantageous to use these two additional parameters to correct for this bias, by disproportionately using the other group's encoder.

Results We next look to the Bayesian "soft averaging" approach to gating, and mitigating bias at inference time. Note that this does not involve retraining the model, as the soft averaging happens at test time. We first evaluate the effectiveness of using a uniform averaging in gated model predictions, where $\alpha = \beta = 0.5$. We label this method "GATE soft", and present results in Table 4. The method results in a much better performance–fairness trade-off than the standard GATE, and for the BIOS dataset, its results are competitive with the best debiasing methods.

We next take the demographic labels into consideration and search for the best gating coefficients for each group. Figure 2 shows accuracy and GAP results from tuning the coefficients on development data for the basic GATE model. The results show that $\alpha = \beta = 0.5$ is a reasonable default setting, however small gains may be possible for non-uniform parameter settings.

To demonstrate the power of adjusting these parameters, we take the trained GATE model, and then optimize α and β over the development set, and report the corresponding results on the test set. We select the parameter values that achieve either: (1) the highest development accuracy; or (2) the lowest development GAP, provided accuracy is above a threshold.⁶ The results are reported in Table 2, un-

⁵Results for Bayesian averaging vs. soft in-network averaging were near identical, hence we report only the latter.

⁶The $[\alpha, \beta]$ values are [0.64,0.66] and [0.00,0.08] for accuracy over MOJI and BIOS, respectively, and [0.64, 0.99] and [0.38, 0.72] for GAP optimised models. We also experi-

616

617

618

619

620

621

622

623

624

625

576

577



Figure 2: Accuracy and GAP of α and β settings for MOJI and BIOS. The axes refer to the propensity to change the gold group in gating the encoder components, and the bottom left point $\alpha = \beta = 0$ is the GATE model using true demographic inputs. Lighter shading denotes better performance.

der GATE ^{soft} and GATE ^{soft}_{RMS}, respectively. On the MOJI dataset, our results show that GATE with soft averaging can consistently outperform the STAN-DARD and GATE models without balanced training. In terms of GAP, the model is substantially better than all other models, while remaining competitive in terms of accuracy. The BIOS dataset is noisier, meaning there are bigger discrepancies between the development and test datasets. As a result, the accuracy-optimized model under performs below the standard GATE model in terms of both accuracy and fairness. However, we achieve a good performance–fairness trade-off when optimizing for GAP, at a level comparable to the much more complex INLP and DADV models.

5 Related Work

553

555

557

559

561

562

563 564

570

571

573

574

575

Fairness Much work on algorithmic fairness has focused on group fairness, i.e. disparities in error rates across groups defined by protected attributes, such as gender, age, or race. Many criteria have been proposed for group fairness, such as statistical parity (Dwork et al., 2012) and equal opportunity (Hardt et al., 2016). Broadly speaking, fairness can be classified into three categories: independence, separation, and sufficiency (Barocas et al., 2019), with the most recent work addressing separation criteria, i.e. potential correlations between main task labels and protected attributes.

Mitigating bias Many approaches for bias mitigation haven been proposed in recent work, including removing protected information form hidden representations (Li et al., 2018a; Wang et al., 2019; Ravfogel et al., 2020; Han et al., 2021b), preprocessing data to remove bias (Zhao et al., 2018; Vanmassenhove et al., 2018; Saunders and Byrne, 2020), modifying the training algorithm (Badjatiya et al., 2019), and post-hoc correction (Hardt et al., 2016).

In the context of NLP, the best results have been achieved through protected information removal. Iterative nullspace projection (**INLP**: Ravfogel et al. (2020)) takes hidden representations and projects them onto the nullspace of the weights of a linear classifier for each protected attribute. The classifier training and projection are carried out over multiple iterations to more comprehensively remove protected information.

Another popular approach is adversarial training, which jointly optimizes the removal of sensitive information and main task performance, through the incorporation of adversarial discriminator(s) to identify protected attributes from the hidden representations (Li et al., 2018a; Elazar and Goldberg, 2018; Wang et al., 2019). Differentiated adversarial learning (**DADV**: Han et al. (2021b)) uses an ensemble of adversaries for each protected attribute, subject to an orthogonality constraint.

6 Conclusions and Future Work

This paper proposed the adoption of balanced training approaches to mitigate bias, and demonstrated their effectiveness relative to existing methods, as well as their ability to further enhance existing methods. We also proposed a gated model based on demographic attributes as an input, and showed that while the simple version was highly biased, with a simple Bayesian extension at inference time, the method was highly effective at mitigating bias.

For future work, it is important to consider settings where there are multiple protected attributes, such as author age, gender, and ethnicity. A simple extension would be to treat G as being *intersectional classes*, defined as the Cartesian product of the multiple demographic groups. E.g., k binary groups would result in 2^k intersectional classes.

mented with adjusting the gating coefficients for the GATE + RW model, in which case there was no benefit to accuracy or GAP from using non-zero α or β .

References

626

631

632

641

643

646

647

649

651

652

653

654

655

657

658

671 672

673

674

675

676

677

- Pinkesh Badjatiya, Manish Gupta, and Vasudeva Varma. 2019. Stereotypical bias removal for hate speech detection task using knowledge-based generalizations. In *The World Wide Web Conference*, pages 49–59.
- Solon Barocas, Moritz Hardt, and Arvind Narayanan. 2019. *Fairness and Machine Learning*. http:// www.fairmlbook.org.
- Su Lin Blodgett, Lisa Green, and Brendan O'Connor.
 2016. Demographic dialectal variation in social media: A case study of African-American English. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 1119–1130.
 - Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In Advances in Neural Information Processing Systems, pages 4349–4357.
- Konstantinos Bousmalis, George Trigeorgis, Nathan Silberman, Dilip Krishnan, and Dumitru Erhan. 2016. Domain separation networks. In Advances in Neural Information Processing Systems, volume 29.
- Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. 2019. Class-balanced loss based on effective number of samples. In *Conference on Computer Vision and Pattern Recognition*.
- Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, and Adam Tauman Kalai. 2019. Bias in bios: A case study of semantic representation bias in a high-stakes setting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 120–128.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186.
- Mark Díaz, Isaac Johnson, Amanda Lazar, Anne Marie Piper, and Darren Gergle. 2018. Addressing agerelated bias in sentiment analysis. In *Proceedings* of the 2018 CHI Conference on Human Factors in Computing Systems, pages 1–14.
- Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, pages 214–226.

Yanai Elazar and Yoav Goldberg. 2018. Adversarial removal of demographic attributes from text data. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 11–21, Brussels, Belgium. Association for Computational Linguistics. 681

682

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

702

703

704

705

706

707

709

710

711

712

713

714

715

716

717

719

720

721

722

723

724

725

726

727

729

730

732

733

734

- Bjarke Felbo, Alan Mislove, Anders Søgaard, Iyad Rahwan, and Sune Lehmann. 2017. Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Xudong Han, Timothy Baldwin, and Trevor Cohn. 2021a. Decoupling adversarial training for fair NLP. In Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021, pages 471–477.
- Xudong Han, Timothy Baldwin, and Trevor Cohn. 2021b. Diverse adversaries for mitigating bias in training. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2760–2765.
- Moritz Hardt, Eric Price, and Nati Srebro. 2016. Equality of opportunity in supervised learning. *Advances in Neural Information Processing Systems*, 29:3315– 3323.
- Lisa Anne Hendricks, Kaylee Burns, Kate Saenko, Trevor Darrell, and Anna Rohrbach. 2018. Women also snowboard: Overcoming bias in captioning models. In *Proceedings of the European Conference on Computer Vision*, pages 793–811.
- Michael Höfler, Hildegard Pfister, Roselind Lieb, and Hans-Ulrich Wittchen. 2005. The use of weights to account for non-response and drop-out. *Social Psychiatry and Psychiatric Epidemiology*, 40(4):291– 299.
- Dirk Hovy. 2015. Demographic factors improve classification performance. In Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 752–762.
- Dirk Hovy and Diyi Yang. 2021. The importance of modeling social factors of language: Theory and practice. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 588–602, Online. Association for Computational Linguistics.
- Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis Kalantidis. 2020. Decoupling representation and classifier for long-tailed recognition. In *Eighth International Conference on Learning Representations (ICLR)*.

- 792 793 794 795
- 796 797 798 799

801

Diederick P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*.

736

741

742

743

745

746 747

748

749 750

751

753

754

756

757

758

759

761

771

772

773

774

775

779

784 785

786

787

789

790

- 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. In Advances in Neural Information Processing Systems, volume 33, pages 728–740.
- Yitong Li, Timothy Baldwin, and Trevor Cohn. 2018a. Towards robust and privacy-preserving text representations. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics* (Volume 2: Short Papers), pages 25–30.
- Yitong Li, Timothy Baldwin, and Trevor Cohn. 2018b. What's in a domain? learning domain-robust text representations using adversarial training. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers), pages 474–479, New Orleans, Louisiana. Association for Computational Linguistics.
- Kaiji Lu, Piotr Mardziel, Fangjing Wu, Preetam Amancharla, and Anupam Datta. 2018. Gender bias in neural natural language processing. *arXiv preprint arXiv:1807.11714*.
- Ji Ho Park, Jamin Shin, and Pascale Fung. 2018. Reducing gender bias in abusive language detection. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2799–2804.
- Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. 2020. Null it out: Guarding protected attributes by iterative nullspace projection. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7237–7256.
- Danielle Saunders and Bill Byrne. 2020. Reducing gender bias in neural machine translation as a domain adaptation problem. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7724–7736.
- Shivashankar Subramanian, Xudong Han, Timothy Baldwin, Trevor Cohn, and Lea Frermann. 2021. Evaluating debiasing techniques for intersectional biases. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Eva Vanmassenhove, Christian Hardmeier, and Andy Way. 2018. Getting gender right in neural machine translation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3003–3008.
- 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 *Proceedings of the IEEE International Conference on Computer Vision*, pages 5310–5319.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender bias in coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 15–20.

805

807

810

811

812

813

814

815

816

A Dataset distribution

А.1 Мојі

This training dataset has been artificially balanced according to demographic and task labels, but artificially skewed in terms of race–sentiment combinations, as follows: AAE–happy = 40%, SAE–happy = 10%, AAE–sad = 10%, and SAE–sad = 40%.

A.2 BIOS



Figure 3: Bios dataset statistics.

Figure 3 shows the statistic of the BIOS dataset. Each row corresponds to a profession, including the total number of instances and number of female instances. Besides, each profession is also annotated with the percentage of female instances.

B Reproducibility

B.1 Hyperparameter Tuning

All approaches proposed in this paper share the same hyperparameters as the standard model. Hyperparameters are tuned using grid-search, in order to maximise accuracy for the standard model, and to minimise the fairness GAP for debiasing methods, subject to the accuracy exceeding a given threshold. The accuracy threshold is chosen to ensure the selected model achieves comparable performance to baseline methods, defined as up to 2% less than best baseline accuracy. Taking RW as an example, the best baseline accuracy on the BIOS development dataset is 75.7% and accordingly the (development) accuracy threshold is set to 73.7%; among models in the hyperparameter search space that exceed this threshold, we take the model with minimum GAP. We report test results for the selected models.

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

In terms of the baseline models, both DADV and INLP have additional hyperparameters: for DADV these are the weight of the adversarial loss, which controls the performance–fairness trade-off; the number of sub-adversaries; and the weight of the difference loss, to better remove demographic information; while INLP also has a trade-off hyperparameter, the number of null-space projection iterations, and other hyperparameters related to linear attackers and classifiers.

The trade-off hyperparameter makes such models more flexible in performing model selection. However, it also requires manual selection for better trade-offs, and different strategies have been introduced. For example, INLP manually selects the model at a iteration where the accuracy is minimally damaged while the fairness improves greatly. Similar manual selection for better trade-offs is also required for ADV and DADV, but the strategies proposed in the original papers are slightly different to one another, and are also task-specific.

In order to reproduce previous methods, we follow the original paper in setting the accuracy threshold, and then tuning hyperparameters for the best fairness.

For the ADV and DADV models, following the work of Han et al. (2021b), we tune extra hyperparameters separately, such as the trade-off hyperparameter, while using the same shared hyperparameters to the selected base models. Similarly, the number of iterations for the INLP model is tuned once other hyperparameters have been fixed.

B.2 Training Details

We conduct all our experiments on a Windows server with a 16-core CPU (AMD Ryzen Threadripper PRO 3955WX), two NVIDIA GeForce RTX 3090s with NVLink, and 256GB RAM.

В.2.1 Мојі

For all baseline models, we follow the method of Han et al. (2021b). Specifically, we train the

		Best assignment						
Hyperparameter	Search space	STANDARD	ADV	DADV	DS	RW	DADV + DS	DADV + RW
number of epochs	-				100			
patience	-	10						
encoder	-	DeepMoji (Felbo et al., 2017)						
embedding size	-				2304			
hidden size	-				300			
number of hidden layers	choice-integer[1, 3]	2						
batch size	loguniform-integer[64, 2048]	1024	1024	1024	512	1024	512	1024
output dropout	uniform-float[0, 0.5]	0.4	0.4	0.4	0.5	0.5	0.2	0.1
optimizer	-			Adam (Kir	ngma an	d Ba, 20)15)	
learning rate	$loguniform$ -float $[10^{-6}, 10^{-1}]$	3×10^{-5}	$3 imes 10^{-5}$	$3 imes 10^{-5}$	10^{-5}	10^{-4}	3×10^{-5}	$3 imes 10^{-4}$
learning rate scheduler	-			redu	ce on pl	ateau		
LRS patience	-				2 epoch	s		
LRS reduction factor	-				0.5			
ADV loss weight	$loguniform$ -float $[10^{-4}, 10^2]$	-	$10^{-0.1}$	$10^{-0.1}$	-	-	$10^{0.2}$	$10^{0.0}$
ADV hidden size	loguniform-integer[64, 1024]	-	256	256	-	-	256	256
number of adversaries	choice-integer[1, 8]	-	1	3	-	-	3	3
DADV loss weight	$loguniform$ -float $[10^{-5}, 10^5]$	-	-	$10^{3.7}$	-	-	10^{2}	$10^{2.6}$

Table 5: Search space and best assignments on the MOJI dataset

		Best assignment						
Hyperparameter	Search space	STANDARD	ADV	DADV	DS	RW	DADV + DS	DADV + RW
number of epochs	-					100		
patience	-	10						
encoder	-	uncased BERT-base (Devlin et al., 2019)						
embedding size	-					768		
embedding type	choice{'CLS', 'AVG'}				•4	WG'		
hidden size	-	300						
number of hidden layers	choice-integer[1, 3]	2						
batch size	loguniform-integer[64, 2048]	512	128	128	128	256	256	512
output dropout	uniform-float[0, 0.5]	0.5	0.3	0.2	0.3	0.5	0.2	0.4
optimizer	-			Adam	ı (Kingr	na and Ba, 2	015)	
learning rate	$loguniform$ -float $[10^{-6}, 10^{-1}]$	$3 imes 10^{-3}$	10^{-3}	10^{-3}	10^{-3}	$3 imes 10^{-5}$	$3 imes 10^{-3}$	$3 imes 10^{-4}$
learning rate scheduler	-				reduce	on plateau		
LRS patience	-				2 e	pochs		
LRS reduction factor	-					0.5		
ADV loss weight	$loguniform$ -float $[10^{-8}, 10^2]$	-	$10^{-2.3}$	$10^{-2.3}$	-	-	$10^{-2.8}$	10^{-5}
ADV hidden size	loguniform-integer[64, 1024]	-	256	256	-	-	256	256
number of adversaries	choice-integer[1, 8]	-	1	3	-	-	3	3
DADV loss weight	$loguniform$ -float $[10^{-5}, 10^5]$	-	-	10^{2}	-	-	10^{3}	$10^{3.3}$

Table 6: Search space and best assignments on the BIOS dataset

875 STANDARD model for 100 epochs with the Adam 876 optimizer (Kingma and Ba, 2015), learning rate 877 of 3×10^{-5} , and batch size of 1024. For ADV, 878 the main model is jointly trained together with ad-879 versaries which are implemented as 3-layer MLP, 880 and the weight of adversarial loss is 0.8. For each 881 iteration (epoch) of the main model, an adversary is trained for 60 epochs, keeping the checkpoint model that performs best on the dev set. Three sub-adversaries are employed by the DADV, with the difference losss weight of $10^{3.7}$. For INLP, logistic regression models are used for both identifying null-space to the demographic information at each iteration, and making the final predictions

		Мојі					BIOS		
Method	Model	Accuracy↑	$\mathbf{GAP}\downarrow$	$\textbf{DTO}\downarrow$	Time↓	Accuracy ↑	$\mathbf{GAP}\downarrow$	$\textbf{DTO}\downarrow$	Time↓
D1:	STANDARD	71.6 ± 0.1	31.0 ± 0.3	0.261	1.0	82.3 ± 0.0	16.0 ± 0.5	0.110	1.0
	INLP	68.5 ± 1.1	33.8 ± 3.9	0.300	14.0	70.5 ± 0.5	6.7 ± 0.9	0.145	6.3
Dasennes	Adv	74.3 ± 0.4	22.2 ± 3.7	0.163	36.1	81.1 ± 0.1	12.7 ± 0.3	0.077	1.3
	DADV	74.5 ± 0.3	18.5 ± 2.0	0.123	109.4	81.1 ± 0.1	12.6 ± 0.3	0.076	2.4
Dalamaa	DS	71.9 ± 0.1	23.2 ± 0.2	0.178	0.5	79.4 ± 0.1	9.7 ± 0.6	0.057	0.3
Balance	RW	74.0 ± 0.2	21.5 ± 0.4	0.155	1.0	74.7 ± 0.3	7.4 ± 0.3	0.095	1.0
	GATE	64.8 ± 0.1	65.2 ± 0.9	0.640	1.0	82.4 ± 0.1	19.2 ± 0.3	0.144	1.0
Gate	GATE + DS	72.5 ± 0.0	16.3 ± 0.7	0.104	0.6	79.4 ± 0.1	9.2 ± 0.2	0.053	0.3
	GATE + RW	74.9 ± 0.2	13.8 ± 0.3	0.072	1.1	74.9 ± 0.2	7.1 ± 0.2	0.092	1.0
	GATE $_{0.5}^{\text{soft}}$	72.7 ± 0.2	30.2 ± 0.3	0.250	1.0	80.8 ± 0.1	11.6 ± 0.3	0.066	1.0
Bayesian	GATE $_{Acc}^{soft}$	74.8 ± 0.2	20.3 ± 0.3	0.142	1.0	81.1 ± 0.1	19.8 ± 0.4	0.151	1.0
	$GATE ^{soft}_{RMS}$	73.5 ± 0.2	7.1 ± 0.3	0.019	1.0	80.5 ± 0.1	11.1 ± 0.3	0.063	1.0
Combination	DADV + DS	72.2 ± 0.2	14.3 ± 0.2	0.085	72.1	79.3 ± 0.1	9.9 ± 0.2	0.059	2.3
	INLP + DS	72.1 ± 1.6	18.4 ± 3.1	0.127	6.3	73.2 ± 0.6	5.9 ± 0.8	0.112	1.3
	DADV + RW	74.6 ± 0.1	18.9 ± 0.3	0.127	108.2	74.1 ± 0.2	7.2 ± 0.4	0.102	3.0
	INLP + RW	72.3 ± 1.9	15.7 ± 3.1	0.099	13.9	73.6 ± 0.6	5.6 ± 0.7	0.107	6.3

Table 7: Results over the sentiment analysis (MOJI) and biography classification (BIOS) tasks. Trade-offs are measured by the normalized Euclidean distance between each model and the ideal model, and lower is better. **Bold** = best trade-off within category. Training time is reported relative to STANDARD, which takes 35 secs and 16 mins for MOJI and BIOS, respectively.

given debiased hidden representations. Since the number of iterations in INLP is highly affected by the random seed at each run, we re-select it at each iteration.

As for our models, the DS model is trained with the learning rate of 10^{-5} and batch size of 512; the RW is trained with the learning rate of 10^{-4} and batch size of 1024; and the GATE is trained with the the set of hyperparameters to the base model.

B.2.2 BIOS

890

891

892

893

899

900

901

902

903

904

905

906

907

908

910

911

912

913

914

Models are trained with similar hyperparameters as models on the MOJI dataset. We thus only report main differences for each of them: the STANDARD model is trained with the batch size of 512 and learning rate of 3×10^3 ; DS models are trained with the batch size of 128 and learning rate of 10^{-3} , and RW models are trained with the batch of 256 and learning rate of 3×10^{-5} .

We train the ADV model with the adversarial loss weight of $10^{-2.3}$, learning rare for adversarial training of 10^{-1} , learning rate of 10^{-3} , and batch size of 128. The DADV is trained with same setting as the ADV, excepting the difference loss weight of 10^2 . For details of the assignment of other hyperparameters and hyperparameter searching space, refer to Supplementary Materials.

C The calculation of trade-off

We calculate DTO based on all results shown in Table 7. Taking the DAdv model on the Moji dataset for example, the trade-off is calculated as follows: 915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

- 1. Find the best accuracy and fairness (1-GAP) separately; i.e., 74.9 (GATE + RW) and 92.9 (GATE $_{\text{RMS}}^{\text{soft}}$), resp.
- 2. Normalize the accuracy and fairness metric of DADV, resulting in $0.995 = \frac{74.5}{74.9}$ and $0.877 = \frac{81.5}{92.9}$.
- 3. Calculate the Euclidean distance between (1, 1) and (0.995, 0.877), giving 0.123.

D Training time estimation

Given that the training time is affected by factors such as batch size, hidden size, and learning rate, to perform a fair comparison between different models, we estimate the training time of a model based on hyperparameter tuning results, over a shared search space of base hyperparameters (i.e., the hyperparameters related to the standard model), with any other approach-specific hyperparameters fixed.

E Balancing toward anti-stereotyping

As shown in Table 2, even with DS or RW balancing, the model still shows biases in its predictions.

We conduct preliminary experiments on MOJI with 939 RW and DS, while controlling for stereotyping 940 skew in training using values for 0.8 to 0.2. In 941 standard rebalancing we use as target 0.5, which describes a balanced situation. A larger skew > 0.5943 will amplifying stereotyping, and < 0.5 describes 944 a different type of stereotyping operating in the op-945 posite direction. Balancing towards a 0.4 training 946 skew leads to the best test results, with an accuracy of 71.7% and GAP of 11.8% for DS, and accuracy 948 of 74.5% and GAP of 11.3% for RW. Comparing to the corresponding values in Table 2 (rows Bal-950 ance DS and RW, for MOJI), both results show a 951 substantial reduction in GAP. 952

954

955

956

957

958

961

962

963

964

965

966

967 968

969

This idea is related to existing reweighting approaches in long-tail learning. For example, Cui et al. (2019) infer the effective number of samples which group each instance with its neighbours within a small region instead of using all data points, and reweight the loss of each class inversely proportional to the effective number of samples. We leave this further exploration of this line of research to future work.

We also experiment with GATE +RW and GATE +DS with a 0.4 training skew, however, the gated model does not show the same behaviour, as it just amplifies the training biases. This implies that, for the gated model, balanced training can help remove spurious correlations between protected attributes and main task labels, which is similar in nature to the effects of adversarial training.