Differentiable Optimal Adversaries for Learning Fair Representations

Anonymous Authors

Abstract

Fair representation learning is an important task in 1 many real-world domains, with the goal of finding a 2 performant model that obeys fairness requirements. 3 We present an adversarial representation learning 4 algorithm that learns an informative representation 5 6 while not exposing sensitive features. Our goal is 7 to train an embedding such that it has good performance on a target task while not exposing sensitive 8 information as measured by the performance of an 9 optimally trained adversary. Our approach directly 10 trains the embedding with these dual objectives in 11 mind by implicitly differentiating through the opti-12 mal adversary's training procedure. To this end, we 13 derive implicit gradients of the optimal logistic re-14 gression parameters with respect to the input train-15 ing embeddings, and use the fully-trained logistic 16 regression as an adversary. As a result, we are able 17 to train a model without alternating min max op-18 19 timization, leading to better training stability and 20 improved performance. Given the flexibility of our module for differentiable programming, we evalu-21 ate the impact of using implicit gradients in two ad-22 versarial fairness-centric formulations. We present 23 quantitative results on the trade-offs of target and 24 fairness tasks in several real-world domains. 25

26 Introduction

Deep learning models learn expressive data representations 27 which make them applicable in many settings such as health-28 care, criminal justice, or financial support. However, when 29 used in automatic processes, practitioners often want to en-30 sure that the model is performing fairly, with a variety of 31 approaches enforcing different forms of fairness [Mehrabi 32 et al., 2019]. One way to approach fairness is to ensure 33 the learned latent representation doesn't encode any sensi-34 tive information such as race or gender [Zemel et al., 2013]. 35 Several recent works learn fair representations through ad-36 versarial representation learning (ARL). In ARL approaches, 37 an embedding model is trained such that a classifier has 38 good performance on a target task, while also ensuring that 39 an optimally trained adversary has poor performance ex-40 tracting the sensitive information. Many of the ARL ap-41

proaches use a multi-agent approach, alternating between training the embedding and adversary [Xie *et al.*, 2017; 43 Roy and Boddeti, 2019]. However, these alternating ARL approaches disregard how changes in the embedding impact the corresponding new optimized adversary. As a result, they can suffer from training instability and suboptimality. 47

We propose an approach that directly trains the embedding 48 by treating the optimal adversary as a differentiable function 49 of the latent representation. We incorporate the adversarial 50 loss in the training, by considering adversary's model pa-51 rameters as an implicit differentiable function of the em-52 bedding. We derive gradients for the optimal logistic regres-53 sion solution with respect to the input embedding, thus en-54 abling backpropagation from the adversary loss and the ap-55 plication of the optimal adversary model to the embedding, 56 through the optimality conditions of the adversary, back to 57 the model parameters. 58

Our contributions are: 1) develop an end-to-end adversarial learning methodology that does not alternate between the target and sensitive attribute tasks, but instead optimizes both jointly; 2) derive how to incorporate optimal logistic regression as a differentiable layer in predictive models, which is interesting its own right; 3) show that our approaches often provide better tradeoffs between target and sensitive accuracy (as well as demographic parity) on diverse set of domains.

59

60

61

62

63

64

65

66

67

Problem Formulation

We consider that we are given data with features, target labels, and sensitive labels $\{(x^{(i)}, t^{(i)}, s^{(i)})\}_{i=1}^n$ with $x^{(i)} \in \mathbb{R}^{d_f}$ 69 being d_f – dimensional feature vectors, and target labels 70 $t^{(i)} \in \mathbb{R}^{d_t}$ and $s^{(i)} \in 2^{c_s}$ being one-hot sensitive labels 71 among c_s sensitive classes. 72

The goal is to find a classifier parameterized by embedding 73 parameters θ_e , and target classifier θ_t such that the feature 74 extractor with weights W, trained against our embedding θ_e , 75 has poor performance. We can consider that the sensitive ad-76 versary is a linear logistic function of the embedding as in 77 [Roy and Boddeti, 2019]. We consider the embedding func-78 tion $z(x^{(i)}; \theta_e) \in \mathbb{R}^{d_e}$ to return a representation of an exam-79 ple in the latent space of dimensionality d_e . 80

We consider the 3-player game proposed in [Roy and Boddeti, 2019], where the adversary minimizes a loss $V_a(\theta_e, W)$, and the target classifier and embedding minimize their own



Figure 1: Fair representation learning model computation diagram.

loss, linearly weighting a penalty from the performance of the adversary $V_p(\theta_e, W)$ and the predictive performance on the target data $V_t(\theta_e, \theta_t)$. The adversarial penalty coefficient α is a tradeoff parameter that determines the weight on the adversarial penalty V_p . This setting is represented as the bilevel optimization problem:

$$\min_{\theta_e,\theta_t,W^*} \quad V_t(\theta_e,\theta_t) + \alpha V_p(\theta_e,W^*)$$
(1a)

s.t.
$$W^* = \arg\min_{W} V_a(\theta_e, W)$$
 (1b)

Here Equation 1a represents the overall loss, a linear combination of the target classification performance and the sensitive penalty. Similarly, Equation 1b ensures the adversarial weights W^* optimize the adversary's objective V_a .

Considering that our setting consists of supervised learn-94 ing tasks, we consider the target and adversary clas-95 sifiers output predictions for targets $\hat{t}(z(x;\theta_e);\theta_t)$ and 96 sensitive labels $\hat{s}(z(x;\theta_e);W)$ respectively. We de-97 fine the target and adversary objective functions us-98 ing standard supervised losses, with target classifier loss 99 $V_t(\theta_e, \theta_t) = L_t(t, \hat{t}(z(x; \theta_e); \theta_t)),$ and adversary classifier 100 loss $V_a(\theta_e, W) = L_a(s, \hat{s}(z(x; \theta_e); W))$. We now define the 101 target and adversary loss functions as well as the adversarial 102 penalty to fully specify our problem. 103

104 Target loss function: V_t

This loss function represents the performance of the classifier on the target class. It is a supervised loss $V_t(\theta_e, \theta_t) = L_t(t, \hat{t}(z(x; \theta_e); \theta_t))$ with L_t being a differentiable supervised loss function such as cross-entropy loss.

109 Adversary loss function: V_a

We consider the adversary to be solving a logistic regression 110 problem, so our loss function on the adversary's weights W111 is considered to be the logistic loss with L2 penalty. Given 112 the one-hot encoded sensitive targets s, and softmax pre-113 dictions $\hat{s}(z(x;\theta_e);W)) = \sigma(W^T z^{(i)}(x;\theta_e))$, the softmax 114 regression loss is $V_a(\theta_e, W) = L_a(s, \hat{s}(z(x; \theta_e); W)) =$ 115 $-\sum_{i=1}^{n} s^{(i)T} \sigma(W^T z^{(i)}(x; \theta_e)) + \|W\|_2^2$. Although the func-116 tions here are known to be differentiable, our approach will 117 take gradients of the optimal weights W^* with respect to the 118 input embeddings $z(x; \theta_e)$ to perform backpropagation. 119

Adversarial penalty: V_p

Lastly, given our flexible formulation, we can consider both formulations of adversarial representation learning (ARL) presented in [Roy and Boddeti, 2019], one penalizing the embedding based on the entropy of the optimal adversary (referred to as MaxEnt-ARL), and another based on adversary's classification performance (referred to as ML-ARL).

120

145

Optimizing the entropy considers that we want to max-127 imize the entropy of the sensitive classifier's predictions. 128 For simplicity, we can consider minimizing the cross-129 entropy between the uniform distribution and the predictions 130 $\hat{s}(x;\theta_e,W^*)$. Thus we can formulate entropy maximiza-131 tion as minimizing $V_p(\theta_e, W^*) = L_p(s, \hat{s}(z(x; \theta_e); W^*)) =$ 132 $CE(1/c_s, \hat{s}(z(x; \theta_e); W^*))$, with CE(p, q) being the cross 133 entropy between p and q i.e. $CE(p,q) = -\sum_{i=1}^{c_s} p \log_2 q$. 134 Note that in this setting, the adversary penalty disregards the 135 sensitive labels, but the sensitive labels will still be used in 136 the training of the adversary. 137

To encode ML-ARL in our formulation, we can consider the adversary penalty V_p to be the negative of the classification performance of the worst-case adversary. In this case we would have $V_p(\theta_e, W^*) = L_p(s, \hat{s}(z(x; \theta_e); W^*)) =$ $-CE(s, \hat{s}(z(x; \theta_e); W^*))$, or the negative of the cross entropy between the sensitive labels and the adversary's predictions of the sensitive labels.

Evaluating the objective: Equation 1a

Given this problem formulation, we can clearly evaluate the objective function we are trying to minimize given embedding and target parameters θ_e, θ_t .

Examining the pipeline in algorithm 1 and visualized in 149 Figure 1, we can now begin to see that what is easily dif-150 ferentiable in parameters θ_e, θ_t . Clearly step 1, 2, and 3 are 151 known differentiable functions of the weights so standard li-152 braries will handle backpropagation. Furthermore, step 5 is 153 clearly a differentiable function of both the embedding and 154 the optimal logistic layer so a standard autograd library will 155 chain together gradients from softmax and product rule for 156 differentiating $W^{*T}\mathbf{z}$. In step 6, the adversarial penalty loss 157 is a standard cross entropy loss on the predictions. Lastly, 158 the returned loss is a simple linear combination. Therefore, 159 the only component that does not yet have readily-available 160 gradient computation is step 4. 161

Algorithm 1: Compute objective function

1 Embed $\mathbf{z} \leftarrow z(x; \theta_e)$

- ² Predict targets $\hat{\mathbf{t}} \leftarrow \hat{t}(\mathbf{z}; \theta_t)$
- 3 Compute $V_t \leftarrow L_t(t, \hat{\mathbf{t}})$
- 4 Optimize Logistic Regression

$$W^* \leftarrow \arg\min - \sum_{i=1}^{n} s^{(i)T} \sigma(W^T \mathbf{z}) + ||W||_2^2$$

- ${}^{S} \frac{1}{W} \qquad \sum_{i=1}^{N}$
- 5 Predict sensitive $\hat{\mathbf{s}} \leftarrow \sigma(W^{*T}\mathbf{z})$
- 6 Compute $V_p \leftarrow L_p(s, \hat{\mathbf{s}})$
- 7 Return $V_t + \alpha V_p$

Our approach derives gradients of the optimal solution to 162 163 the logistic regression problem W^* with respect to the input feature embeddings z so that we can backpropagate from the 164 loss function, through the logistic regression training, to the 165 original embedding for training. 166

Differentiating Through Adversary 167

Optimization 168

Given that the rest of the pipeline is specified for both for-169 ward and backward passes, here we investigate gradients for 170 the remaining step: step 4 in algorithm 1. We derive gradients 171 of the optimal logistic regression parameters $W^* \in \mathbb{R}^{d_e imes c_s}$ 172 with respect to the input features z. Here the logistic regres-173 sion makes predictions for c_s classes from d_e features. Given 174 that the objective of the logistic regression is convex in it's 175 weights [Boyd and Vandenberghe, 2004], we know that the 176 optimal solution is defined as the solution where the gradient 177 of the objective function is 0. Thus we know that W^* must 178 satisfy the constraints 179

$$0 = \nabla_W \left|_{W^*} \left(-\sum_{i=1}^n s_i^T \sigma \left(W^T \mathbf{z} \right) + \|W\|_2^2 \right).$$

We write the gradients of the logistic regression objective 180 with respect to the model parameters evaluated at optimality: 181

$$\nabla_{W} \bigg|_{W^{*}} \left(-\sum_{i=1}^{n} s_{i}^{T} \sigma \left(W^{T} \mathbf{z} \right) + \|W\|_{2}^{2} \right)$$
$$= \sum_{i=1}^{n} \left(\sigma \left(W^{*T} \mathbf{z} \right) - s_{i}^{T} \right) \mathbf{z} + 2W^{*}$$

Here we can see that the trained parameters W^* are an implicitly defined function of the embedding z, namely those which ensure the gradients are 0. Thus, to find gradients of the optimal parameters W^* with respect to a single embedding $\mathbf{z}^{\mathbf{i}_0}$ of example i_0 , we can relate changes in W^* to changes in z^{i_0} as those satisfying a set of equations. Specifically, we have that for each sensitive class $k \in [c_s]$,

$$\sum_{i=1}^{n} \left[\sum_{c=1}^{c_s} (\delta_{c,k} - \hat{s}_c^{(i)}) s_k^{(i)} \left(dW_c^{*T} \mathbf{z}^{(i)} + W_c^{*T} d\mathbf{z}^{(i)} \right) \right] \mathbf{z}^{(i)} + (\hat{s}_k^{(i_0)} - s_k^{(i_0)}) d\mathbf{z}^{(i_0)} = 0,$$

where δ is the Kronecker delta. 182

Experiments

We train all methods with early stopping based on the valida-184 tion loss of the encoder. We selected model hyperparameters 185 and architectures for the embedding model and target classi-186 fier from [Roy and Boddeti, 2019]. 187

Methods

MLP is a fairness-unaware neural network classifier to mini-189 mize a target loss without regard for the sensitive classifier. 190

CE-ARL [Xie et al., 2017], Ent-ARL [Roy and Boddeti, 191 **2019**] are standard alternating approaches. CE-ARL imposes 192 an adversarial penalty on the embedding of the negative cross entropy loss. EntARL uses the prediction entropy as the ad-194 versaryial loss. 195

CE-OptARL, Ent-OptARL are the corresponding vari-196 ants of our method which penalize our embedding using the 197 negative of the adversary's cross-entropy and the adversary's 198 output entropy respectively. This method follows the same 199 mathematical program as [Xie et al., 2017] but fully opti-200 mizes the adversary model instead of iteratively training the 201 embedding and the adversary. 202

Datasets

COMPAS [Angwin *et al.*, 2016] has defendant data where 204 we aim to predict whether the person will recidivate within 205 2 years, being sensitive to race. Heritage Health data con-206 tains features about 60,000 patients from insurance claims 207 and physician records. As in [Madras et al., 2018; Song et al., 208 2018], we consider the target task of predicting whether the 209 Charlson Index is nonzero being sensitive to age group (9 age 210 groups total). Adult is a UCI dataset [Frank and Asuncion, 211 2010] of 40,000 adults where the task is to predict whether 212 the income is above \$50,000, while being sensitive to gender. 213 German is another UCI dataset [Frank and Asuncion, 2010] 214 of 1,000 people where the task is to predict low or high credit 215 score while being sensitive to gender. 216

Evaluation

Sensitive Accuracy evaluates the sensitive information in an 218 embedding. We train a logistic regression classifier to predict 219 the sensitive features from the embeddings of the training set 220 and evaluate the test accuracy of that fully-trained model. 221

Demographic Parity Difference. The demographic parity 222 difference Δ_{DP} [Dwork *et al.*, 2011] measures the difference 223 in selection rates between sensitive groups and is defined for 224 targets predictions t and sensitive labels s as 225

$$\Delta_{DP} = |P(\hat{t} = 1|s = 1) - P(\hat{t} = 1|s = 0)|.$$

Results. Table 1 reports best target accuracy achieved by 226 each method at different cutoffs of sensitive accuracy and 227 demographic parity (Δ_{DP}). Results spanning this tradeoff 228 are collected by varying the adversarial penalty coefficient α 229 between 0.1 and 1000 by factors of 10, for all methods but 230 MLP. Each method and parameter setting is run with 5 ran-231 dom seeds. We observe that our approaches, CE-OptARL 232 and Ent-OptARL, outperform their respective standard ARL 233 counterparts. The OptARL approaches provide better target 234 accuracy at the given sensitive accuracy cutoffs, demonstrat-235 ing that differentiating through the adversary's optimization 236

183

193

188

203

217

COMPAS	sens acc < 0.98	sens acc < 0.99	sens acc < 1.00	$\Delta_{DP} < 0.10$	$\Delta_{DP} < 0.15$	$\Delta_{DP} < 0.20$
MLP	-	-	0.6961	-	0.6945	0.6961
CE-ARL	-	0.5429	0.6848	0.6005	0.6572	0.6848
CE-OptARL (ours)	0.701	0.701	0.701	0.6969	0.701	0.701
Ent-ARL	-	0.6921	0.6921	0.6669	0.6872	0.6921
Ent-OptARL (ours)	0.701	0.701	0.701	0.7002	0.7002	0.7002
Health	sens acc < 0.30	sens acc < 0.32	sens acc < 0.34	$\Delta_{DP} < 0.40$	$\Delta_{DP} < 0.60$	$\Delta_{DP} < 0.80$
MLP	-	0.8177	0.8192	-	0.8192	0.8192
CE-ARL	-	0.8176	0.8176	0.708	0.8176	0.8176
CE-OptARL (ours)	0.8165	0.8178	0.8178	-	0.8178	0.8178
Ent-ARL	0.7492	0.8184	0.8194	0.7066	0.8194	0.8194
Ent-OptARL (ours)	0.8203	0.8203	0.8203	0.6883	0.8203	0.8203
Adult	sens acc < 0.68	sens acc < 0.69	sens acc < 0.70	$\Delta_{DP} < 0.10$	$\Delta_{DP} < 0.15$	$\Delta_{DP} < 0.20$
MLP	0.8216	0.8242	0.8242	-	0.8242	0.8242
CE-ARL	0.8163	0.8163	0.8163	0.814	0.8163	0.8163
CE-OptARL (ours)	0.8248	0.8248	0.8248	0.8167	0.8248	0.8248
Ent-ARL	0.8186	0.821	0.821	0.8153	0.821	0.821
Ent-OptARL (ours)	0.8192	0.827	0.827	0.8013	0.827	0.827
German	sens acc < 0.90	sens acc < 0.95	sens acc < 1.00	$\Delta_{DP} < 0.02$	$\Delta_{DP} < 0.03$	$\Delta_{DP} < 0.04$
MLP	0.6933	0.73	0.73	0.6933	0.6933	0.6933
CE-ARL	0.6967	0.71	0.71	0.6967	0.6967	0.6967
CE-OptARL (ours)	0.72	0.72	0.72	0.7	0.72	0.72
Ent-ARL	0.7067	0.7067	0.7067	0.69	0.7	0.7067
Ent-OptARL (ours)	0.7333	0.7333	0.7333	0.7267	0.7267	0.7333

Table 1: Target accuracy at fairness cutoffs: We present test results for maximum target accuracy at given cutoffs on the accuracy of a fullytrained adversary (sens acc), as well as on the demographic parity (Δ_{DP}). These cutoffs are selected for each dataset to span the distribution in the results. Metrics are obtained by varying the adversarial penalty coefficient α between 0.1 and 1000 by factors of 10.

procedure is able to improve the desired effect of adversar-ial representation learning. In addition, we observe that our

methods provide better target accuracy at most Δ_{DP} cutoffs, with the exception of the Adult and Health datasets only at

the lowest Δ_{DP} threshold.

242 Related Work

In [Zemel et al., 2013], the authors optimize clusters of indi-243 viduals to generate discrete and fair representations. [Calmon 244 et al., 2017] optimize a random data transformation preserv-245 ing utility for downstream tasks but obfuscating sensitive at-246 tributes. Approaches with alternating training such as [Roy 247 and Boddeti, 2019; ?] iteratively train an embedding along 248 with an adversary by optimizing the models with respective 249 parameters and objectives. These approaches generally for-250 mulate the objective of the embedding using an optimal ad-251 versary; however, the optimiziation procedures don't differ-252 entiate through the adversary's optimization procedure, and 253 instead treat the adversary's parameters as constants during 254 backpropagation to the embedding model. Previous work has 255 considered a similar differentiable optimization approach for 256 meta-learning, proposing a differentiable svm optimization 257 algorithm [Lee et al., 2017], closed-form ridge-regression 258 formulation, or iterative logistic regression solver [Bertinetto 259 et al., 2019] as a last-layer fine tuning methodology. In our 260 work, we consider adversarial representation learning, and di-261 262 rectly differentiate through the optimality condition of logistic regression rather than the unrolled solver iterates. 263

Discussion

We improve adversarial representation learning approaches 265 by implicitly defining the fully-trained adversary as a differ-266 entiable function of the embedding, allowing us to directly 267 train the representation with gradient information from the 268 adversary's optimality conditions. In particular, we provide 269 a novel methodology for computing gradients of the optimal 270 logistic regression adversary with respect to the input embed-271 dings. This approach can be viewed in several lights. One in-272 terpretation is that we fully backpropagate the global loss (the 273 penalty of the adversary and the target performance) through 274 the adversary optimization to the embedding model's param-275 eters. Another facet is that we train the embedding with ex-276 plicit information about how the fully-trained adversary will 277 change due to changes in the embedding. Lastly, we can view 278 the overall optimization procedure as optimizing the embed-279 ding for the loss it observes at equilibrium in the 3-player 280 game formulation suggested in [Roy and Boddeti, 2019]. 281

264

Since our contribution enables logistic regression fitting as a differentiable layer in any end-to-end learning, we hope in future work to evaluate other relevant settings. 291

292 **References**

- [Angwin *et al.*, 2016] Julia Angwin, Jeff Larson, Surya
 Mattu, and Lauren Kirchner. Machine bias. *ProPublica*,
 May, 23:2016, 2016.
- 296 [Bertinetto et al., 2019] Luca Bertinetto, Joao F. Henriques,
- Philip Torr, and Andrea Vedaldi. Meta-learning with dif-ferentiable closed-form solvers. In *International Confer-*
- *ence on Learning Representations*, 2019.
- Boyd and Vandenberghe, 2004] Stephen P Boyd and Lieven
 Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [Calmon *et al.*, 2017] Flavio Calmon, Dennis Wei, Bhanuki ran Vinzamuri, Karthikeyan Natesan Ramamurthy, and
- 305 Kush R Varshney. Optimized pre-processing for discrimi-
- nation prevention. In I. Guyon, U. V. Luxburg, S. Bengio,
- H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Sys*-
- *tems 30*, pages 3992–4001. Curran Associates, Inc., 2017.
- [Dwork *et al.*, 2011] Cynthia Dwork, Moritz Hardt, Toniann
 Pitassi, Omer Reingold, and Richard Zemel. Fairness
 Through Awareness, apr 2011.
- [Frank and Asuncion, 2010] Andrew Frank and Arthur
 Asuncion. Uci machine learning repository [http://archive.
 ics. uci. edu/ml]. irvine, ca: University of california. *School of information and computer science*, 213(11),
 2010.
- [Lee *et al.*, 2017] Hsin-Ying Lee, Jia-Bin Huang, Maneesh
 Singh, and Ming-Hsuan Yang. Unsupervised representation learning by sorting sequences. In *Proceedings of the IEEE International Conference on Computer Vision*, pages
 667–676, 2017.
- Imadras *et al.*, 2018] David Madras, Elliot Creager, Toniann
 Pitassi, and Richard Zemel. Learning adversarially fair
 and transferable representations. In *International Confer- ence on Machine Learning*, pages 3384–3393, 2018.
- [Mehrabi *et al.*, 2019] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A
 survey on bias and fairness in machine learning. *arXiv* preprint arXiv:1908.09635, 2019.
- [Roy and Boddeti, 2019] Proteek Chandan Roy and
 Vishnu Naresh Boddeti. Mitigating information leakage
 in image representations: A maximum entropy approach.
 Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019.
- [Song *et al.*, 2018] Jiaming Song, , Aditya Grover, Shengjia
 Zhao, and Stefano Ermon. Learning controllable fair representations. *arXiv preprint arXiv:1812.04218*, 2018.
- [Xie *et al.*, 2017] Qizhe Xie, Zihang Dai, Yulun Du, Ed-
- uard Hovy, and Graham Neubig. Controllable invariance
 through adversarial feature learning. In *Advances in Neu-*
- ral Information Processing Systems, pages 585–596, 2017.
- [Zemel *et al.*, 2013] Rich Zemel, Yu Wu, Kevin Swersky,
 Toni Pitassi, and Cynthia Dwork. Learning fair representations. In Sanjoy Dasgupta and David McAllester,

editors, *Proceedings of the 30th International Conference* 346 on Machine Learning, volume 28 of *Proceedings of Machine Learning Research*, pages 325–333, Atlanta, Georgia, USA, 17–19 Jun 2013. PMLR. 349