

Differentiable Optimal Adversaries for Learning Fair Representations

Aaron Ferber¹, Umang Gupta², Greg Ver Steeg², Bistra Dilkina¹,

¹ University of Southern California

² USC Information Sciences Institute

{aferber, umanggup, dilkina}@usc.edu, {gregv}@isi.edu

Abstract

Fair representation learning is an important task in many real-world domains, with the goal of finding a performant model that obeys fairness requirements. We present an adversarial representation learning algorithm that learns an informative representation while not exposing sensitive features. Our goal is to train an embedding such that it has good performance on a target task while not exposing sensitive information as measured by the performance of an optimally trained adversary. Our approach directly trains the embedding with these dual objectives in mind by implicitly differentiating through the optimal adversary’s training procedure. To this end, we derive implicit gradients of the optimal logistic regression parameters with respect to the input training embeddings, and use the fully-trained logistic regression as an adversary. As a result, we are able to train a model without alternating min max optimization, leading to better training stability and improved performance. Given the flexibility of our module for differentiable programming, we evaluate the impact of using implicit gradients in two adversarial fairness-centric formulations. We present quantitative results on the trade-offs of target and fairness tasks in several real-world domains.

Introduction

Deep learning models learn expressive data representations which make them applicable in many settings such as health-care, criminal justice, or financial support. However, when used in automatic processes, practitioners often want to ensure that the model is performing fairly, with a variety of approaches enforcing different forms of fairness [Mehrabi *et al.*, 2019]. One way to approach fairness is to ensure the learned latent representation doesn’t encode any sensitive information such as race or gender [Zemel *et al.*, 2013]. Several recent works learn fair representations through adversarial representation learning (ARL). In ARL approaches, an embedding model is trained such that a classifier has good performance on a target task, while also ensuring that an optimally trained adversary has poor performance extracting the sensitive information. Many of the ARL ap-

proaches use a multi-agent approach, alternating between training the embedding and adversary [Xie *et al.*, 2017; 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.

We propose an approach that directly trains the embedding by treating the optimal adversary as a differentiable function of the latent representation. We incorporate the adversarial loss in the training, by considering **adversary’s model parameters as an implicit differentiable function of the embedding**. We derive gradients for the optimal logistic regression solution with respect to the input embedding, thus enabling backpropagation from the adversary loss and the application of the optimal adversary model to the embedding, through the optimality conditions of the adversary, back to the model parameters.

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 in 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.

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 $x^{(i)} \in \mathbb{R}^{d_f}$ being d_f – dimensional feature vectors, and target labels $t^{(i)} \in \mathbb{R}^{d_t}$ and $s^{(i)} \in 2^{c_s}$ being one-hot sensitive labels among c_s sensitive classes.

The goal is to find a classifier parameterized by embedding parameters θ_e , and target classifier θ_t such that the feature extractor with weights W , trained against our embedding θ_e , has poor performance. We can consider that the sensitive adversary is a linear logistic function of the embedding as in [Roy and Boddeti, 2019]. We consider the embedding function $z(x^{(i)}; \theta_e) \in \mathbb{R}^{d_e}$ to return a representation of an example in the latent space of dimensionality d_e .

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

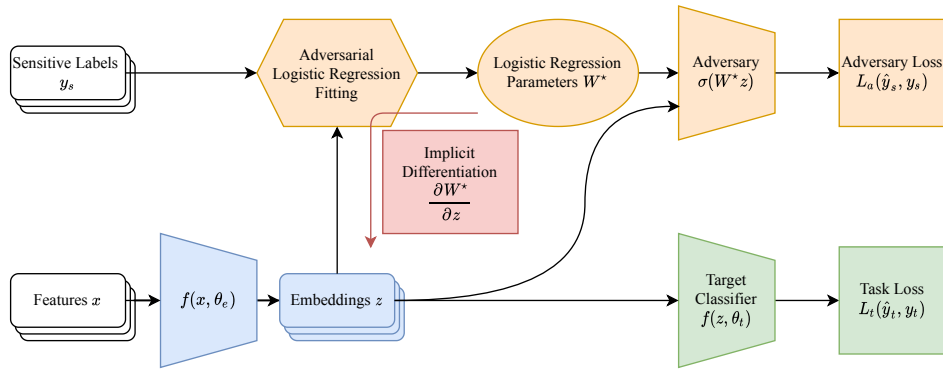


Figure 1: Fair representation learning model computation diagram.

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

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

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

90 Here Equation 1a represents the overall loss, a linear combina-
 91 tion of the target classification performance and the sensi-
 92 tive penalty. Similarly, Equation 1b ensures the adversarial
 93 weights W^* optimize the adversary’s objective V_a .

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

104 Target loss function: V_t

105 This loss function represents the performance of the classi-
 106 fier on the target class. It is a supervised loss $V_t(\theta_e, \theta_t) =$
 107 $L_t(t, \hat{t}(z(x; \theta_e); \theta_t))$ with L_t being a differentiable super-
 108 vised loss function such as cross-entropy loss.

109 Adversary loss function: V_a

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

Adversarial penalty: V_p

120 Lastly, given our flexible formulation, we can consider both
 121 formulations of adversarial representation learning (ARL)
 122 presented in [Roy and Boddeti, 2019], one penalizing the em-
 123 bedding based on the entropy of the optimal adversary (re-
 124 ferred to as MaxEnt-ARL), and another based on adversary’s
 125 classification performance (referred to as ML-ARL).
 126

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

138 To encode ML-ARL in our formulation, we can consider
 139 the adversarial penalty V_p to be the negative of the classifica-
 140 tion performance of the worst-case adversary. In this case
 141 we would have $V_p(\theta_e, W^*) = L_p(s, \hat{s}(z(x; \theta_e); W^*)) =$
 142 $-CE(s, \hat{s}(z(x; \theta_e); W^*))$, or the negative of the cross
 143 entropy between the sensitive labels and the adversary’s pre-
 144 dictions of the sensitive labels.

Evaluating the objective: Equation 1a

145 Given this problem formulation, we can clearly evaluate the
 146 objective function we are trying to minimize given embed-
 147 ding and target parameters θ_e, θ_t .
 148

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

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**

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

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

167 Differentiating Through Adversary 168 Optimization

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

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

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

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

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

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

182 where δ is the Kronecker delta.

Experiments

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

Methods

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

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

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

Datasets

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

Evaluation

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

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

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

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

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 fully-trained 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.

237 procedure is able to improve the desired effect of adversarial
238 representation learning. In addition, we observe that our
239 methods provide better target accuracy at most Δ_{DP} cutoffs,
240 with the exception of the Adult and Health datasets only at
241 the lowest Δ_{DP} threshold.

242 Related Work

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

Discussion

264 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 The evaluation using four different datasets, spanning
282 criminal risk assessment, healthcare, and finance, showed that
283 our optimal adversary approach improves the performance of
284 both adversarial representation learning baselines. In particu-
285 lar, we showed we are able to (almost always) provide better
286 target accuracy at different thresholds on fairness in terms of
287 both sensitive accuracy and demographic parity.

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

292 References

- 293 [Angwin *et al.*, 2016] Julia Angwin, Jeff Larson, Surya
294 Mattu, and Lauren Kirchner. Machine bias. *ProPublica*,
295 *May*, 23:2016, 2016.
- 296 [Bertinetto *et al.*, 2019] Luca Bertinetto, Joao F. Henriques,
297 Philip Torr, and Andrea Vedaldi. Meta-learning with dif-
298 ferentiable closed-form solvers. In *International Confer-*
299 *ence on Learning Representations*, 2019.
- 300 [Boyd and Vandenberghe, 2004] Stephen P Boyd and Lieven
301 Vandenberghe. *Convex optimization*. Cambridge univer-
302 sity press, 2004.
- 303 [Calmon *et al.*, 2017] Flavio Calmon, Dennis Wei, Bhanuki-
304 ran Vinzamuri, Karthikeyan Natesan Ramamurthy, and
305 Kush R Varshney. Optimized pre-processing for discrimi-
306 nation prevention. In I. Guyon, U. V. Luxburg, S. Bengio,
307 H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett,
308 editors, *Advances in Neural Information Processing Sys-*
309 *tems 30*, pages 3992–4001. Curran Associates, Inc., 2017.
- 310 [Dwork *et al.*, 2011] Cynthia Dwork, Moritz Hardt, Toniann
311 Pitassi, Omer Reingold, and Richard Zemel. Fairness
312 Through Awareness, apr 2011.
- 313 [Frank and Asuncion, 2010] Andrew Frank and Arthur
314 Asuncion. Uci machine learning repository [[http://archive.](http://archive.ics.uci.edu/ml)
315 [ics.uci.edu/ml](http://archive.ics.uci.edu/ml)]. irvine, ca: University of california.
316 *School of information and computer science*, 213(11),
317 2010.
- 318 [Lee *et al.*, 2017] Hsin-Ying Lee, Jia-Bin Huang, Maneesh
319 Singh, and Ming-Hsuan Yang. Unsupervised representa-
320 tion learning by sorting sequences. In *Proceedings of the*
321 *IEEE International Conference on Computer Vision*, pages
322 667–676, 2017.
- 323 [Madras *et al.*, 2018] David Madras, Elliot Creager, Toniann
324 Pitassi, and Richard Zemel. Learning adversarially fair
325 and transferable representations. In *International Confer-*
326 *ence on Machine Learning*, pages 3384–3393, 2018.
- 327 [Mehrabi *et al.*, 2019] Ninareh Mehrabi, Fred Morstatter,
328 Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A
329 survey on bias and fairness in machine learning. *arXiv*
330 *preprint arXiv:1908.09635*, 2019.
- 331 [Roy and Boddeti, 2019] Proteek Chandan Roy and
332 Vishnu Naresh Boddeti. Mitigating information leakage
333 in image representations: A maximum entropy approach.
334 *Proceedings of the IEEE Conference on Computer Vision*
335 *and Pattern Recognition*, 2019.
- 336 [Song *et al.*, 2018] Jiaming Song, , Aditya Grover, Shengjia
337 Zhao, and Stefano Ermon. Learning controllable fair rep-
338 resentations. *arXiv preprint arXiv:1812.04218*, 2018.
- 339 [Xie *et al.*, 2017] Qizhe Xie, Zihang Dai, Yulun Du, Ed-
340 uard Hovy, and Graham Neubig. Controllable invariance
341 through adversarial feature learning. In *Advances in Neu-*
342 *ral Information Processing Systems*, pages 585–596, 2017.
- 343 [Zemel *et al.*, 2013] Rich Zemel, Yu Wu, Kevin Swersky,
344 Toni Pitassi, and Cynthia Dwork. Learning fair repre-
345 sentations. In Sanjoy Dasgupta and David McAllester,

editors, *Proceedings of the 30th International Conference* 346
on Machine Learning, volume 28 of *Proceedings of Ma-* 347
chine Learning Research, pages 325–333, Atlanta, Geor- 348
gia, USA, 17–19 Jun 2013. PMLR. 349