

FAIRNESS-AWARE CONTRASTIVE LEARNING WITH PARTIALLY ANNOTATED SENSITIVE ATTRIBUTES

Anonymous authors

Paper under double-blind review

ABSTRACT

1 Learning high-quality representation is important and essential for visual recogni-
 2 tion. Unfortunately, traditional representation learning suffers from fairness issues
 3 since the model may learn information of sensitive attributes. Recently, a series of
 4 studies have been proposed to improve fairness by explicitly decorrelating target
 5 labels and sensitive attributes. Most of these methods, however, rely on the as-
 6 sumption that fully annotated labels on target variable and sensitive attributes are
 7 available, which is unrealistic due to the expensive annotation cost. In this paper,
 8 we investigate a novel and practical problem of **Fair Unsupervised Representation**
 9 **Learning with Partially annotated Sensitive labels (FURL-PS)**. FURL-PS has two
 10 key challenges: 1) how to make full use of the samples that are not annotated with
 11 sensitive attributes; 2) how to eliminate bias in the dataset without target labels.
 12 To address these challenges, we propose a general **Fairness-aware Contrastive**
 13 **Learning (FairCL)** framework consisting of two stages. Firstly, we generate con-
 14 trastive sample pairs, which share the same visual information apart from sensitive
 15 attributes, for each instance in the original dataset. In this way, we construct a bal-
 16 anced and unbiased dataset. Then, we execute fair contrastive learning by closing
 17 the distance between representations of contrastive sample pairs. Besides, we also
 18 propose an unsupervised way to balance the utility and fairness of learned rep-
 19 resentations by feature reweighting. Extensive experimental results illustrate the
 20 effectiveness of our method in terms of fairness and utility, even with very limited
 21 sensitive attributes and serious data bias.

22 1 INTRODUCTION

23 Learning powerful representation takes an important role in visual recognition, and there are a lot of
 24 works proposed to learn visual representations (Bengio et al., 2013; Kolesnikov et al., 2019; Wang
 25 et al., 2020a). Among them, contrastive learning achieves state-of-the-art performance on various
 26 computer vision tasks (Tian et al., 2020; Chuang et al., 2020). Contrastive learning first generates
 27 views (patches) from original images by random data augmentation, and the views from the same
 28 image are defined as positive samples. Then the model can learn effective representations by closing
 29 the distance between representations of positive samples, while being protected from mode collapse
 30 via an additional module such as negative samples (Chen et al., 2020a; He et al., 2020; Chen et al.,
 31 2020b), momentum update (Grill et al., 2020), and stopping gradient (Chen & He, 2021).

32 Unfortunately, traditional representation learning methods ignore potential fairness issues, which
 33 becomes an increasing concern as recognition systems are widely used in the real world (Zemel
 34 et al., 2013; Madras et al., 2018; Creager et al., 2019). For example, the model trained by con-
 35 trastive learning may learn the information of sensitive attributes (*e.g.*, gender, race) by using it as a
 36 shortcut to minimize the distance between representations of positive samples in the training stage,
 37 since the positive samples have the same sensitive attributes. As a result, decisions based on biased
 38 representation models may discriminate against certain groups or individuals in practice, by using
 39 spurious correlations between predictive target and sensitive attributes (Wang et al., 2020b; Park
 40 et al., 2021). Therefore, how to develop a fair representation model is of paramount importance for
 41 both academic research and real applications.

42 Most of existing works achieve fairness via decorrelating target labels and sensitive attributes explic-
 43 itly, which rely on the data annotations (Mehrabi et al., 2021; Zhu et al., 2022). However, assuming

44 that all data have fully annotated labels can be unrealistic (Jung et al., 2022; Zhang et al., 2020a;b).
45 In many real scenarios, the target labels and even downstream tasks are not provided, and all we have
46 are images and limited annotations of sensitive attributes. Data labels require additional expensive
47 cost of human annotations, which naturally leads us to ask the following question: *Can we train a*
48 *fair unsupervised representation model with only partially annotated sensitive attributes?*

49 In this paper, we investigate a practical and novel problem of Fair Unsupervised Representation
50 Learning with Partially annotated Sensitive attributes (FURL-PS). Our goal is to utilize the images
51 and limited sensitive labels to learn visual representations that can be used for various downstream
52 tasks of visual recognition, while achieving fairness by being minimally correlated with sensitive at-
53 tributes. It is challenging to solve the proposed problem. Firstly, most samples are not labeled with
54 sensitive attributes. A natural idea is to pseudo-label the unlabeled data by a sensitive attribute clas-
55 sifier. However, it is not advisable to train a representation model on the data with pseudo-sensitive
56 labels, since the noises in pseudo labels may severely affect the fairness performance. Secondly,
57 there may be data imbalance between demographic groups. Assuming that the female group has a
58 large proportion of samples of blond hair, while the male group has the opposite proportion. As a
59 result, the models trained on the above biased data may learn spurious correlation between gender
60 and blond hair. Unfortunately, it is difficult to balance the data distribution of different groups with-
61 out the prior of downstream tasks or annotated target labels. Generally, FURL-PS problem has two
62 main challenges: 1) How to make full use of the data that are not annotated with sensitive attributes?
63 2) How to balance the possible agnostic bias in data without target labels?

64 To address these challenges, our idea is to construct a balanced dataset annotated with sensitive
65 labels based on the original dataset, and then train a representation model with fair contrastive learn-
66 ing on the unbiased dataset. We propose a two-stage Fairness-aware Contrastive Learning (*FairCL*)
67 framework to implement the above idea. In the first stage, we design a semi-supervised learning
68 algorithm to train the image attribute editor with limited sensitive labels, which is used to edit the
69 pre-defined sensitive attributes of a given image. In the second stage, we train a representation model
70 by fair contrastive learning with balanced augmentation. Specifically, based on the image attribute
71 editor, we can generate contrastive sample pairs, which share the same visual information apart from
72 sensitive attributes (e.g., male and female), for each sample in the original dataset. By closing the
73 distance between representations of contrastive sample pairs, the model can learn powerful and fair
74 representations. Our approach has two advantages: 1) we can get the utmost out of unlabeled im-
75 ages by generating samples with given sensitive attributes from them; 2) the augmented dataset is
76 unbiased, since it consists of contrastive sample pairs and thus the data proportions are naturally
77 balanced for different demographic groups. Furthermore, we also develop an unsupervised way to
78 balance the utility and fairness of learned representations by feature reweighting.

79 We validate the effectiveness of our method on two facial attribute recognition datasets: CelebA (Liu
80 et al., 2018) and UTK-Face (Zhang et al., 2017). Extensive experimental results show that the pro-
81 posed method outperforms the existing unsupervised learning methods in terms of both classifica-
82 tion accuracy and fairness, and even achieves comparable performance with the semi-supervised
83 methods that require annotations on the target labels. Besides, our method is robust to the ratio of
84 sensitive labels and severity of data bias. Furthermore, we also show the extensibility of our general
85 framework to different contrastive learning algorithms through experiments.

86 **Main Contributions:** 1) To the best our knowledge, we are the first one to propose the practical and
87 challenging problem of Fair Unsupervised Representation Learning with only Partially annotated
88 Sensitive attributes (FURL-PS). 2) We develop the Fairness-aware Contrastive Learning (*FairCL*)
89 framework to solve the proposed problem, which can be compatible with all of contrastive learning
90 algorithms to learn a fair and powerful representation model. 3) Extensive experiments illustrate the
91 effectiveness of our proposed method in terms of fairness and utility.

92 2 RELATED WORK

93 **Fairness in Unsupervised and Semi-supervised Learning.** There are three branches to guarantee
94 fairness in unsupervised learning. Firstly, fair feature selection, a pre-processing paradigm, finds
95 a subset of features that preserve the original information as much as possible while being mini-
96 mally correlated with sensitive attributes (Grgic-Hlaca et al., 2016; Grgić-Hlača et al., 2018; Xing
97 et al., 2021). However, this kind of methods is designed for structured data, and cannot be applied

98 to images data and deep models. Secondly, fair clustering balances the distribution of different
 99 subgroups formed by sensitive attributes in each cluster, but it cannot yield a model for various
 100 downstream tasks (Chierichetti et al., 2017; Kleindessner et al., 2019; Li et al., 2020). At last, some
 101 studies based on fair representation learning have been bringing a paradigm for fair unsupervised
 102 learning (Louizos et al., 2015; Raff & Sylvester, 2018). As for fair semi-supervised learning, most
 103 existing methods first pseudo-label the unlabeled data via a classifier, and then train a model on these
 104 data with fairness constraints (Jung et al., 2022; Zhang et al., 2020a;b). However, the pseudo-label
 105 noise may exacerbate model unfairness in turn. Instead, our proposed method does not directly use
 106 pseudo-labels when training the fair representation model.

107 **Contrastive Learning.** Self-supervised contrastive learning provides a representation learning
 108 paradigm without target labels, and achieves better accuracy than the state-of-the-art methods on
 109 various tasks (Xiao et al., 2020). Some methods such as *SimCLR* (Chen et al., 2020a) and *MoCo* (He
 110 et al., 2020) first define the positive/negative samples as patches generated from the same/different
 111 images via random data augmentation, and then train a representation model by closing/pushing
 112 away the distance between representations of positive/negative samples. Recent studies argue that
 113 negative samples are not necessary for contrastive learning, and they use some techniques, *e.g.*, mo-
 114 mentum update (Grill et al., 2020) and stopping gradient (Chen & He, 2021) to protect the model
 115 from mode collapse instead of negative samples. Afterwards, supervised contrastive learning out-
 116 performs other state-of-the-art methods based on traditional cross-entropy loss (Khosla et al., 2020).
 117 However, existing self-supervised contrastive learning methods ignore potential fairness issues. To
 118 this end, *FSCL* proposes a fair supervised contrastive loss to train a fair representation model (Park
 119 et al., 2022). However, *FSCL* relies on target labels and sensitive attributes. Besides, *FSCL* is based
 120 on supervised contrastive learning which needs negative samples, while our proposed framework is
 121 general to be applied to any contrastive learning algorithm to improve fairness.

122 **Image Generation.** Our proposed method involves the task of image attribute editing, which takes
 123 an image as input and aims to generate a new image with desired attributes while preserving other
 124 details (Liu et al., 2019; He et al., 2019; Dogan & Keles, 2020; Wang et al., 2022). We emphasize
 125 that advances in the field of image attribute editing can help improve the performance of our work,
 126 since the subsequent methods can also be used here. Some studies aim to construct a balanced and
 127 unbiased dataset by data augmentation (Ramaswamy et al., 2021). However, they need the prior of
 128 downstream task to generate new samples. Recent works have proposed to evaluate counterfactual
 129 fairness by generating counterfactual samples (Denton et al., 2019; Joo & Kärkkäinen, 2020; Dash
 130 et al., 2022). Different from them, we consider a more challenging problem to train a fair represen-
 131 tation model. Moreover, we consider a more practical setting where there are no target labels and
 132 fully annotated sensitive attributes.

133 3 METHOD

134 In this section, we start with a brief introduction of the problem formulation of FURL-PS and overall
 135 flow of our proposed method in Sec. 3.1. Then we display how to generate augmented samples of
 136 different sensitive attributes with limited annotated sensitive labels in Sec. 3.2. We elaborate on how
 137 to execute fair contrastive learning with balanced augmentation in Sec. 3.3. Lastly, to balance the
 138 trade-off between utility and fairness of learned representations, we propose a feature reweighting
 139 module for those sensitive attribute-dependent sub-features in Sec. 3.4.

140 3.1 PROBLEM FORMULATION AND OVERALL FLOW

141 Assume that we have n original images $\{x_k\}_{k=1,2,\dots,n}$, where $x_k \in \mathcal{X} \subset \mathbb{R}^d$. Labeled dataset is
 142 denoted as $D_l = \{x_k, s_k\}_{k=1}^{n_l}$, where n_l is the number of images with annotated sensitive labels, and
 143 $s_k \in \{0, 1, \dots, M_S - 1\}$ represents the sensitive attribute label (*e.g.*, male and female). Unlabeled
 144 dataset is denoted as $D_u = \{x_k\}_{k=n_l+1}^n$. The target labels $\{y_k\}_{k=1,2,\dots,n}$ are not available in the
 145 training state, where $y_k \in \{0, 1, \dots, M_Y - 1\}$. In this paper, we assume that both target labels
 146 and sensitive attributes are binary variables for convenience, *i.e.*, $M_S = M_Y = 2$. We emphasize,
 147 however, that our proposed problem and framework can be easily generalized to multivariate setting.
 148 The goal is to train an effective and fair encoder network $F(\cdot)$ that maps the image $x_k \in \mathcal{X}$ into
 149 representation $h_k \in \mathcal{H}$, where the representations can be used for various downstream tasks while
 150 not discriminating against demographic groups with given sensitive attributes.

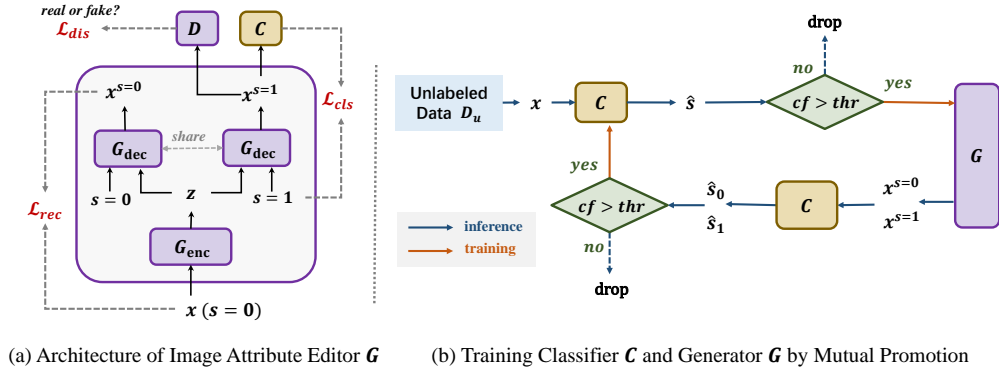


Figure 1: Semi-supervised Learning of Image Sensitive Attribute Editor.

151 Our proposed method consists of two stages: 1) **Contrastive Sample Generation** and 2) **Fairness-**
 152 **aware Contrastive Learning**. We first define the *contrastive samples* as a pair of images that share
 153 the same visual information except for sensitive attributes. In the first stage, our goal is to prepare
 154 contrastive samples $\{(x_k^0, x_k^1)\}_{k=1,2,\dots,n}$ based on the original dataset. To achieve it, we design a
 155 semi-supervised algorithm to train an image sensitive attribute editor $G(\cdot, \cdot)$, which takes an image
 156 x_k and sensitive attribute $s \in \{0, 1\}$ as input and can map the original image to a new image x_k^s
 157 with given sensitive attribute s while keeping other information as unchanged as possible. In the
 158 second stage, we execute fairness-aware contrastive learning on the augmented dataset to train a
 159 representation model $F(\cdot)$ without any target label or sensitive label.

160 3.2 CONTRASTIVE SAMPLE GENERATION WITH LIMITED SENSITIVE ATTRIBUTE LABELS

161 We start with training a generative model used to generate contrastive samples. We implement it
 162 based on *AttGAN* (He et al., 2019), but we emphasize that any approach designed for image attribute
 163 editing can be adapted to be a backbone method. The training architecture of generative model G
 164 is shown in Figure 1(a). The encoder G_{enc} maps an input image x to latent representation z . Then
 165 the decoder G_{dec} takes z and sensitive attributes as input, and generates images with corresponding
 166 sensitive attributes. There are three loss optimized jointly: 1) discriminative loss l_{dis} given by the
 167 discriminator D guaranteeing that the generated images look realistic enough, 2) classification loss
 168 l_{cls} guaranteeing that the generated images have given sensitive attributes, and 3) reconstruction loss
 169 l_{rec} given by the classifier C encourages the generator to preserve the sensitive attribute-excluding
 170 information as much as possible.

171 Note that the sensitive labels are needed in the training stage of the image attribute editor G , while
 172 we only have limited annotated sensitive sensitive labels. To this end, we develop a semi-supervised
 173 learning algorithm to train the sensitive attribute classifier and image editor by making them mutu-
 174 ally promote each other. As shown in Figure 1(b), we first train a classifier C and image editor G
 175 on the labeled data. Then we assign the predictive sensitive labels to the unlabeled data by C . We
 176 only select pseudo-labeled samples with confidence cf above a certain threshold thr , and train the
 177 generative model G on them. Afterwards, we use G to generate additional images with sensitive
 178 labels, of which high-confidence images are used to train the classifier C . We repeat the above steps
 179 until there is no new high-confidence data. Note that the classifier C is used to pseudo-label the
 180 unlabeled data, thereby providing the generator G with more annotated training data. Meanwhile,
 181 the generator G also generates additional high-quality training data for the classifier C . In this way,
 182 the two can promote each other.

183 Here we get an image attribute editor G , based on which we can generate an augmented dataset
 184 $\{(x_k^0, x_k^1)\}_{k=1,2,\dots,n}$ from original dataset, even without the knowledge of the sensitive attributes of
 185 the original images. It is worth mentioning that we only use pseudo-sensitive labels in the contrastive
 186 sample generation stage, which avoids the direct impact of pseudo-label noise on learning fair rep-
 187 resentation. Furthermore, the augmented dataset is unbiased and balanced for sensitive attributes.

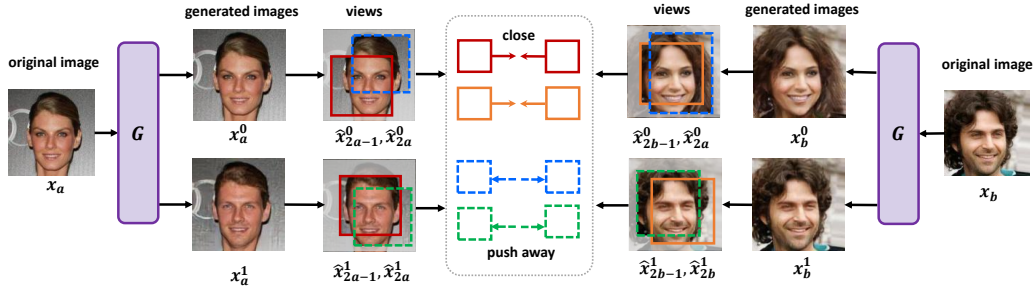


Figure 2: Training Flow of Fairness-aware Contrastive Learning (taking *SimCLR* as an example). Our proposed *FairCL* closes the distance between representations of *positive samples* (i.e., the views from the same original image but have different sensitive attributes), and pushes away the distance between representations of *negative samples* (i.e., the views from different original images but have the same sensitive attributes).

188 3.3 FAIRNESS-AWARE CONTRASTIVE LEARNING WITH BALANCED AUGMENTATION

189 Based on the balanced augmented dataset generated by the image attribute editor G , our goal is to
 190 train a fair and powerful representation model. To this end, we develop a fairness-aware contrastive
 191 learning framework, and elaborate how it works based on *SimCLR*. However, we emphasize that
 192 our proposed general framework is not confined to *SimCLR* but can be easily applied to any con-
 193 trastive learning algorithm. The key idea is to define the positive samples as the contrastive samples
 194 with different sensitive attributes generated by the image attribute editor G , and define the nega-
 195 tive samples as the views generated from the different images with the same sensitive attributes.
 196 Based on it, for a minibatch of contrastive image pairs $\{(x_k^0, x_k^1)\}_{k=1,2,\dots,N}$, we first generate views
 197 $\{(\hat{x}_{2k-1}^0, \hat{x}_{2k}^0, \hat{x}_{2k-1}^1, \hat{x}_{2k}^1)\}_{k=1,2,\dots,N}$ by data augmentation. An encoder $F(\cdot)$ maps the views into
 198 representations $\{(\hat{h}_{2k-1}^0, \hat{h}_{2k}^0, \hat{h}_{2k-1}^1, \hat{h}_{2k}^1)\}_{k=1,2,\dots,N}$, then a projection head $P(\cdot)$ maps h_k into an-
 199 other representations $\{(\hat{z}_{2k-1}^0, \hat{z}_{2k}^0, \hat{z}_{2k-1}^1, \hat{z}_{2k}^1)\}_{k=1,2,\dots,N}$ for contrastive learning. Then we define
 200 fairness-aware contrastive loss as:

$$L^{Fair} = - \sum_{k=1}^N \log \frac{\exp(\hat{z}_{2k-1}^0 \cdot \hat{z}_{2k-1}^1 / \tau)}{\sum_{l \neq k} \exp(\hat{z}_{2k}^0 \cdot \hat{z}_{2l}^0 / \tau) + \sum_{l \neq k} \exp(\hat{z}_{2k}^1 \cdot \hat{z}_{2l}^1 / \tau)}, \quad (1)$$

201 where τ is a temperature parameter. As shown in Figure 2, a fair and effective representation model
 202 can be trained on the balanced augmented dataset with the proposed fairness-aware contrastive loss.

203 3.4 BALANCE UTILITY AND FAIRNESS VIA FEATURE REWEIGHTING

204 There is often a trade-off between utility and fairness of representation (Zhao & Gordon, 2019).
 205 To balance them without target labels, we propose the feature reweighting module. Our idea is
 206 to identify the sensitive attribute-dependent sub-features and reweight them when computing the
 207 similarity/distance between representations in contrastive learning. Intuitively, larger weights for
 208 those sensitive attribute-dependent sub-features will result in a fairer model.

209 The challenge is to judge whether a sub-feature is related to a sensitive attribute. Suppose that
 210 we train a linear classifier C_{probe} to classify sensitive attributes, which takes the representations
 211 generated by the fixed encoder $F(\cdot)$ and projection head $P(\cdot)$ as input. An intuition is that those
 212 sensitive attribute-dependent sub-features are easier to activate when predicting sensitive attributes,
 213 and the corresponding parameters of trained classifier C_{probe} will have a larger absolute value. Based
 214 on it, we propose a simple but effective solution. We alternately optimize the representation model
 215 $F(\cdot)$, $P(\cdot)$ and an additional linear classifier C_{probe} . We can strengthen the fairness constraint by
 216 assigning more weights to those sub-features where the absolute values of corresponding classifier
 217 parameters are large. We use the soft selection which is more flexible and general. Assuming that the
 218 parameters of C_{probe} is $(\theta_1, \theta_2, \dots, \theta_d)$, where d is the number of dimensions of latent representation.

219 Then we can compute the weights $w = (w_1, w_2, \dots, w_d)$ of sub-features as:

$$w_i = \frac{\exp(\alpha \cdot |\theta_i|)}{\sum_{j=1}^d \exp(\alpha \cdot |\theta_j|)}, \quad (2)$$

220 where $\alpha \in \mathbb{R}$ is a scaling parameter, and a larger α means a stronger fairness constraint. Now we
 221 can balance the utility and fairness of learned representation in fairness-aware contrastive learning
 222 via the weights w as the following:

$$L^{FairW} = - \sum_{k=1}^N \log \frac{\exp(\hat{z}_{2k-1}^0 \cdot \hat{z}_{2k-1}^1 \cdot w/\tau)}{\sum_{l \neq k} \exp(\hat{z}_{2k}^0 \cdot \hat{z}_{2l}^0 \cdot w/\tau) + \sum_{l \neq k} \exp(\hat{z}_{2k}^1 \cdot \hat{z}_{2l}^1 \cdot w/\tau)}. \quad (3)$$

223 4 EXPERIMENTS

224 4.1 EXPERIMENTAL SETUP

225 **Datasets.** We validate our method on the following datasets: 1) **CelebA** (Liu et al., 2018) is a dataset
 226 with over 200k facial images, each with 40 binary attributes labels. In this paper, we follow the
 227 setting of the previous works (Park et al., 2022). We select *Male* (m) and *Young* (y) as the sensitive
 228 attributes, and set *Attractive* (a), *Big Nose* (b), and *Bags Under Eyes* (e) as target attributes, which
 229 have the highest Pearson correlation with the sensitive attributes. Besides, to verify the performance
 230 of our method in the setting of multi-target labels and multi-sensitive attributes, we also set $\{Male,$
 231 $Young\}$ as sensitive attribute and $\{Big Nose, Bags Under Eyes\}$ as target label. The downstream
 232 task is unknown and only 5% sensitive attribute labels are available in the training stage. 2) **UTK-**
 233 **Face** (Zhang et al., 2017) contains over 20k facial images, each with attributes labels. We first define
 234 a binary sensitive attribute *Young* based on whether age is under 35 or not and construct a task to
 235 predict whether the facial image is *Male* or not. We further validate the robustness of our method to
 236 the ratio of sensitive labels and data bias on UTK-Face dataset.

237 **Evaluation Metrics.** The goal of FURL-PS is to learn a fair and powerful representation model.
 238 To validate the fairness and utility of the learned representations, we train a linear classifier on top
 239 of the frozen representation model and then use the test performance of the classifier as a proxy for
 240 representation quality. In this paper, we use Equal Odds (EO) (Hardt et al., 2016), one of the most
 241 commonly used notion of group fairness (Dwork et al., 2012), as the fairness metric:

$$\overline{\sum_{\forall y, \hat{y}} \left| P_{s^0}(\hat{Y} = \hat{y} | Y = y) - P_{s^1}(\hat{Y} = \hat{y} | Y = y) \right|}, \quad (4)$$

242 where $\overline{\sum}$ is the averaged sum, Y is target label, \hat{Y} is predictive label given by the classifier, and
 243 $s^0, s^1 \in S$ is the value of sensitive attributes. Following (Jung et al., 2022), we extend EO to
 244 multi-sensitive attribute setting:

$$\max_{\forall s^i, s^j \in S} \overline{\sum_{\forall y, \hat{y}} \left| P_{s^i}(\hat{Y} = \hat{y} | Y = y) - P_{s^j}(\hat{Y} = \hat{y} | Y = y) \right|}, \quad (5)$$

245 where $s^i, s^j \in S$ is the value of sensitive attributes. A smaller EO means a fairer model. Besides,
 246 we use top-1 accuracy (%) to measure the effectiveness of learned representations.

247 **Baselines.** To our best knowledge, there is no existing work focusing on dealing with the problem
 248 of FURL-PS. Therefore, we construct some powerful baselines by combining the SOTA methods
 249 to solve partially annotated sensitive labels and the advanced methods designed for fair unsuper-
 250 vised representation learning. Specifically, *CGL* (Jung et al., 2022) is the SOTA method to solve
 251 the problem of partially annotated sensitive attribute labels by assigning pseudo sensitive labels
 252 based on confidence. *VFAE* (Louizos et al., 2015) is a fair representation learning method based
 253 on a variational autoencoding architecture with priors that encourage latent factors of variation
 254 to be independent of sensitive attribute. We implement its unsupervised version and combine it
 255 with *CGL* as a baseline (*CGL+VFAE*). We also compare our method with the combination of *CGL*
 256 and *GRL* (Raff & Sylvester, 2018), which is an adversarial method used for fair representations
 257 (*CGL+GRL*). Since there are few existing methods for fair unsupervised representation learning, we
 258 also implement some fair supervised representation learning methods which rely on target labels.
 259 Group DRO (*G-DRO*) (Sagawa et al., 2019) is a classical and powerful method for robust and fair



Figure 3: Illustration of contrastive samples generated by sensitive attribute editor.

260 learning by learning a set of weights for different data subgroups. *FSCL* (Park et al., 2022) learns
 261 fair representations based on supervised contrastive learning. For those 5% samples annotated with
 262 sensitive attributes, their target labels are available for *G-DRO* and *FSCL*. The *CGL* strategy is also
 263 used for them (*CGL+G-DRO*, *CGL+FSCL*).

264 **Implementation Details.** We resize the images of CelebA and UTK-Face to 128×128 , and use a
 265 5-layer CNN (Krizhevsky et al., 2017) as the encoder of generative model. Besides, the decoder also
 266 has 5 layers. Since the quality of some generated samples may not be good enough, we select images
 267 based on confidence. Specifically, we use the trained classifier to predict the sensitive attributes of
 268 the generated images, and then remove some low-confidence samples. Besides, to further improve
 269 the quality of training data, we use high-confidence of original images instead of corresponding
 270 generated samples. We use the same random data augmentation strategy as (Chen et al., 2020a).
 271 The projection head $P(\cdot)$ is only used in contrastive learning, and then we remove it. We use the
 272 ResNet-18 (He et al., 2016) as encoder model and a MLP as projection head, and train them via
 273 weighted fairness-aware contrastive loss for 100 epochs. Afterwards, we train a linear classifier on
 274 top of the frozen representation given by encoder $F(\cdot)$ for 10 epochs on the training dataset.

275 4.2 CONTRASTIVE SAMPLE GENERATION EXPERIMENTS

276 We set *Male* as the sensitive attribute, and select some contrastive samples generated by the sensitive
 277 attribute editor, which is trained on CelebA dataset with partially annotated sensitive attributes, as
 278 shown in Figure 3. As can be seen, the generated contrastive samples remain most of the visual
 279 details of the original images, but have different sensitive attributes.

280 We next discuss the issues of proxy variables. First, for those proxy variables which have the
 281 causal/stable correlation with the sensitive attributes, e.g. *Beard* and *Moustache*, we note that the
 282 sensitive attribute editor can learn the correlation between them as we expected. For example, the
 283 output female does not have a beard, even if the input image is a male face with a beard. Unfortu-
 284 nately, for those proxy variables having the extreme spurious correlation with the sensitive attributes,
 285 e.g., *Heavy Makeup* and *Lipstick*, the sensitive attribute editor also learns it, since almost all male
 286 images have no lipstick or heavy makeup. This will result in the representation model not being
 287 able to learn information about these proxy variables and thus unable to make accurate predictions
 288 about these attributes. We would like to emphasize that this issues is an open problem and, to our
 289 best knowledge, there is no existing method that can solve it without any prior. However, we find
 290 that our method exhibits robustness to data bias caused by spurious correlations. Specifically, the
 291 sensitive attribute editor would not change the attributes which have the high Pearson correlation
 292 with the sensitive attributes, e.g., *Big Nose*, due to the reconstruction loss.

293 4.3 RESULTS OF FAIRNESS AND ACCURACY ON CELEBA DATASET

294 We report the classification results of unsupervised methods on CelebA dataset in Table 1. We use
 295 equalized odds (EO) and top-1 accuracy (Acc.) of trained linear classifier to evaluate the fairness and
 296 utility of learned representations, respectively. A smaller EO means a fairer model. We find that the
 297 models trained via *SimCLR* achieves the best accuracy, but suffer from fairness issues. Our proposed
 298 *FairCL* based on *SimCLR* improves the fairness of learned representations. *FairCL* outperforms
 299 other unsupervised baselines (*CGL+VFAE* and *CGL+GRL*) in terms of EO and accuracy.

Table 1: **Unsupervised Attribute Classification Results on CelebA.** To evaluate the quality of representation, we measure equalized odds (EO) and top-1 accuracy (Acc.) of trained linear classifier on CelebA dataset with 5% annotated sensitive attributes. A smaller EO means a fairer model. T and S represent target and sensitive attributes, respectively.

| Method | T=a / S=m | | T=a / S=y | | T=b / S=m | | T=b / S=y | | T=e / S=m | | T=e / S=y | | T=b&e / S=y | | T=a / S=y&m | |
|----------------------|-------------|-------------|-------------|-------------|------------|-------------|------------|-------------|------------|-------------|------------|-------------|-------------|-------------|-------------|-------------|
| | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. |
| <i>CGL+VFAE</i> | 19.1 | 72.7 | 16.2 | 74.0 | 15.5 | 78.3 | 10.6 | 78.8 | 7.6 | 79.6 | 6.9 | 79.5 | 10.2 | 68.7 | 28.7 | 72.7 |
| <i>CGL+GRL</i> | 21.3 | 73.4 | 15.6 | 74.4 | 13.1 | 79.6 | 10.9 | 79.5 | 7.4 | 79.9 | 6.1 | 80.2 | 9.8 | 69.6 | 26.9 | 73.8 |
| <i>SimCLR</i> | 36.2 | 77.1 | 22.5 | 77.1 | 26.9 | 81.5 | 19.2 | 81.5 | 19.7 | 81.2 | 10.6 | 81.2 | 12.8 | 71.5 | 39.6 | 77.7 |
| <i>FairCL (Ours)</i> | 16.8 | 75.3 | 13.1 | 76.9 | 8.4 | 80.0 | 9.2 | 80.3 | 4.2 | 80.8 | 4.5 | 80.5 | 7.8 | 71.3 | 24.5 | 74.1 |

Table 2: **Comparison with the methods relying on target labels on CelebA.** 5% samples are annotated with sensitive attributes, and corresponding target labels are available for those target label-dependent methods. Notably, our proposed *FairCL* does not rely on target labels.

| Method | T=a / S=m | | T=a / S=y | | T=b / S=m | | T=b / S=y | | T=e / S=m | | T=e / S=y | | T=b&e / S=y | | T=a / S=y&m | |
|----------------------|-------------|------|-------------|------|------------|------|------------|------|------------|------|------------|------|-------------|------|-------------|------|
| | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. |
| <i>CE</i> | 27.8 | 77.9 | 19.2 | 77.9 | 20.3 | 82.1 | 14.2 | 82.1 | 16.4 | 82.0 | 11.7 | 82.0 | 9.6 | 71.9 | 36.2 | 77.9 |
| <i>CGL+G-DRO</i> | 14.2 | 73.8 | 11.3 | 75.3 | 7.9 | 77.1 | 6.2 | 76.3 | 4.5 | 76.9 | 5.1 | 76.7 | 5.3 | 67.2 | 21.9 | 71.4 |
| <i>CGL+FSCL</i> | 17.4 | 75.3 | 13.5 | 76.2 | 9.5 | 79.7 | 9.6 | 79.1 | 5.9 | 81.1 | 5.9 | 80.4 | 8.2 | 69.8 | 25.6 | 74.0 |
| <i>FairCL (Ours)</i> | 16.8 | 75.3 | 13.1 | 76.9 | 8.4 | 80.0 | 9.2 | 80.3 | 4.2 | 80.8 | 4.5 | 80.5 | 7.8 | 71.3 | 24.5 | 74.1 |

300 We also compare our *FairCL* with the methods relying on target labels, as shown in Table 2. Notably,
 301 *FairCL* even outperforms semi-supervised methods in some cases. The reason may be that semi-
 302 supervised methods suffers from the issues of pseudo-label noise. Besides, *CGL+G-DRO* achieves
 303 excellent EO but has low accuracy due to the over pessimism problem (Hu et al., 2018).

304 4.4 T-SNE VISUALIZATION

305 To further evaluate the quality of
 306 our learned representations and ex-
 307 plain how our method works, we
 308 sample 1000 images from test
 309 dataset and visualize the represen-
 310 tations of them via t-SNE (Van der
 311 Maaten & Hinton, 2008), as shown
 312 in Figure 4. We find that both
 313 *FairCL* and *SimCLR* can assign
 314 similar features to the images with
 315 same target labels. However, *Sim-*
 316 *CLR* also learn information of sen-
 317 sitive attributes, so that the rep-
 318 resentations given by *SimCLR* can
 319 also be divided by the sensitive at-
 320 tributes. In contrast, our proposed
 321 *FairCL* focuses on both utility and fairness by closing the distance between representa-
 322 tions of contrastive samples, which have similar visual information but have different sensitive features.

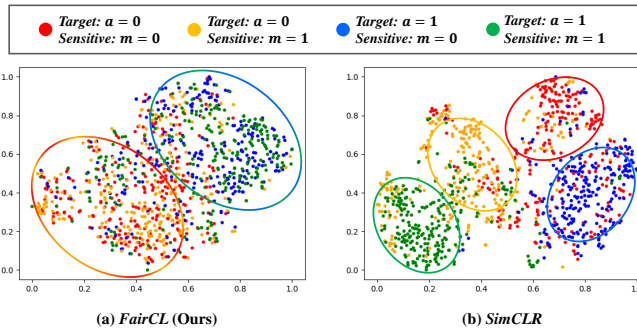


Figure 4: t-SNE visualization for the learned representations.

323 4.5 EFFECTIVENESS OF FEATURE REWEIGHTING

324 We provide an unsupervised way to balance the fairness and utility of learned representation by
 325 feature reweighting. To analyze the effectiveness of feature reweighting, we set the scaling param-
 326 eter α as different values, and train the representation models. The performance of linear classifier
 327 in terms of EO and accuracy is reported in Table 3. As can be seen, we find that a larger scaling
 328 parameter α can yield a fairer model but with lower accuracy. This is in line with our expectations,
 329 since the feature weights will be hard if the scaling parameter α is large, and then the sensitive
 330 attribute-dependent features will have a greater impact on the similarity/distance calculation.

331 4.6 COMPATIBILITY WITH CONTRASTIVE LEARNING ALGORITHMS

332 Our proposed fairness-aware contrastive learning framework is general and flexible that can be used
 333 for any contrastive learning algorithm. We have shown the effectiveness of *SimCLR*-based *FairCL* in

Table 3: **Effectiveness of feature reweighting on CelebA dataset.** We train the representation models with different scaling parameter α to analyze the effectiveness of feature reweighting.

| α | T=a / S=m | | T=a / S=y | | T=b / S=m | | T=b / S=y | | T=e / S=m | | T=e / S=y | | T=b&e / S=y | | T=a / S=y&m | |
|----------------|-----------|------|-----------|------|-----------|------|-----------|------|-----------|------|-----------|------|-------------|------|-------------|------|
| | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. |
| $\alpha = 0.5$ | 16.8 | 75.3 | 13.1 | 76.9 | 8.4 | 80.0 | 9.2 | 80.3 | 4.2 | 80.8 | 4.5 | 80.5 | 7.8 | 71.3 | 24.5 | 74.1 |
| $\alpha = 2.0$ | 14.6 | 74.2 | 11.5 | 76.3 | 6.2 | 79.1 | 6.9 | 79.3 | 3.6 | 80.5 | 2.8 | 80.1 | 5.9 | 70.4 | 22.3 | 72.9 |

Table 4: **Compatibility with contrastive learning on CelebA dataset.** We apply our general framework to *BYOL* and denote it as *FairCL**. The experimental setup is the same as before.

| Method | T=a / S=m | | T=a / S=y | | T=b / S=m | | T=b / S=y | | T=e / S=m | | T=e / S=y | | T=b&e / S=y | | T=a / S=y&m | |
|----------------|-------------|------|-------------|------|-------------|------|-------------|------|------------|------|------------|------|-------------|------|-------------|------|
| | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. | EO | Acc. |
| <i>BYOL</i> | 38.8 | 77.9 | 24.1 | 77.9 | 28.2 | 81.8 | 20.6 | 81.8 | 22.3 | 81.7 | 12.9 | 81.7 | 14.6 | 72.2 | 41.5 | 77.9 |
| <i>FairCL*</i> | 19.2 | 76.2 | 15.4 | 77.1 | 13.7 | 80.8 | 11.0 | 80.9 | 7.2 | 81.0 | 6.8 | 80.9 | 9.9 | 71.6 | 28.6 | 74.9 |

334 the previous subsection. To further illustrate the compatibility of *FairCL*, we apply it to *BYOL* (Grill
 335 et al., 2020), a widely used contrastive learning method without negative samples. Due to the ab-
 336 sence of target labels, negative sample-based methods (e.g., *SimCLR*) may incorrectly push the rep-
 337 resentations of semantically similar samples farther away, which may lead to compromised accuracy
 338 on downstream tasks. *BYOL* overcomes this problem by not using negative samples. Therefore, as
 339 shown in Table 4, *BYOL* achieves excellent accuracy. However, it also suffers from fairness issues.
 340 Notably, *FairCL** improves EO over it while keeping comparable accuracy.

341 4.7 ROBUSTNESS TO RATIO OF ANNOTATED SENSITIVE LABELS AND DATA BIAS

342 To further validate the robustness of our proposed
 343 method to ratio of annotated sensitive labels and un-
 344 known data bias, we run different methods on UTK-
 345 Face dataset. We set *Young* as the sensitive attribute,
 346 and the target label is *Male*. Only 5% samples are
 347 annotated with sensitive attributes. Besides, the dataset
 348 is unbalanced, where the *Young* group has 65% fe-
 349 male data and 35% male, while another sensitive group
 350 has the opposite gender ratio. As can be seen in Fig-
 351 ure 5, our proposed *FairCL/FairCL** improves fairness
 352 compared with *SimCLR/BYOL*, while keeping high ac-
 353 curacy. Our method outperforms other unsupervised
 354 methods and even achieves comparable performance
 355 with the target label-dependent methods (*CGL+G-DRO*
 356 and *CGL+FSCL*), in terms of the trade-off between ac-
 357 curacy and EO, which illustrates that our method is ro-
 358 bust to annotation ratio and data bias.

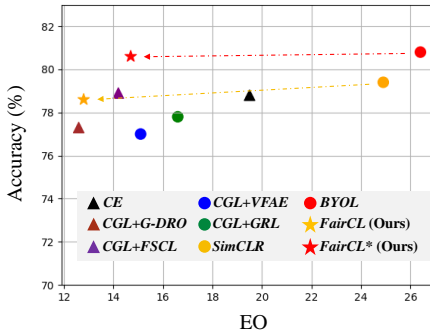


Figure 5: Robustness to ratio of sensitive labels and data bias on UTK-Face. \triangle represents methods relying on target labels.

359 5 CONCLUSIONS

360 In this paper, we investigate a novel and practical problem of which the goal is to learn a fair
 361 and powerful representation model with no target label and limited sensitive attributes. To solve
 362 this problem, we develop a general contrastive learning-based framework *FairCL* consisting of two
 363 stages: contrastive sample generation and fairness-aware contrastive learning with feature reweight-
 364 ing. Extensive experiments show that our proposed method can yield a fair representation model
 365 even with limited sensitive attributes and imbalanced data.

366 We admit that there are some limitations of our method and can be improved in the future work.
 367 Firstly, our work relies on the quality of the generated images. More effective image attribute editing
 368 methods can help improve the representation model trained by our method in terms of both fairness
 369 and utility. Secondly, our method still requires a small number of sensitive attribute labels. How to
 370 learn a fair representation model without any demographic information is still an open problem. We
 371 hope our study can bring some inspiration for future work.

372 REFERENCES

- 373 Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new
374 perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828,
375 2013.
- 376 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
377 contrastive learning of visual representations. In *International conference on machine learning*,
378 pp. 1597–1607. PMLR, 2020a.
- 379 Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of*
380 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15750–15758, 2021.
- 381 Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum
382 contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020b.
- 383 Flavio Chierichetti, Ravi Kumar, Silvio Lattanzi, and Sergei Vassilvitskii. Fair clustering through
384 fairlets. *Advances in Neural Information Processing Systems*, 30, 2017.
- 385 Ching-Yao Chuang, Joshua Robinson, Yen-Chen Lin, Antonio Torralba, and Stefanie Jegelka. De-
386 biased contrastive learning. *Advances in neural information processing systems*, 33:8765–8775,
387 2020.
- 388 Elliot Creager, David Madras, Jörn-Henrik Jacobsen, Marissa Weis, Kevin Swersky, Toniann Pitassi,
389 and Richard Zemel. Flexibly fair representation learning by disentanglement. In *International*
390 *conference on machine learning*, pp. 1436–1445. PMLR, 2019.
- 391 Saloni Dash, Vineeth N Balasubramanian, and Amit Sharma. Evaluating and mitigating bias in
392 image classifiers: A causal perspective using counterfactuals. In *Proceedings of the IEEE/CVF*
393 *Winter Conference on Applications of Computer Vision*, pp. 915–924, 2022.
- 394 Emily Denton, Ben Hutchinson, Margaret Mitchell, Timnit Gebru, and Andrew Zaldívar.
395 Image counterfactual sensitivity analysis for detecting unintended bias. *arXiv preprint*
396 *arXiv:1906.06439*, 2019.
- 397 Yahya Dogan and Hacer Yalim Keles. Semi-supervised image attribute editing using generative
398 adversarial networks. *Neurocomputing*, 401:338–352, 2020.
- 399 Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness
400 through awareness. In *Proceedings of the 3rd innovations in theoretical computer science confer-*
401 *ence*, pp. 214–226, 2012.
- 402 Nina Grgic-Hlaca, Muhammad Bilal Zafar, Krishna P Gummadi, and Adrian Weller. The case for
403 process fairness in learning: Feature selection for fair decision making. In *NIPS symposium on*
404 *machine learning and the law*, volume 1, pp. 2. Barcelona, Spain, 2016.
- 405 Nina Grgić-Hlača, Muhammad Bilal Zafar, Krishna P Gummadi, and Adrian Weller. Beyond dis-
406 tributive fairness in algorithmic decision making: Feature selection for procedurally fair learning.
407 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- 408 Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre Richemond, Elena
409 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,
410 et al. Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural*
411 *information processing systems*, 33:21271–21284, 2020.
- 412 Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. *Advances*
413 *in neural information processing systems*, 29, 2016.
- 414 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
415 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
416 770–778, 2016.
- 417 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for
418 unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on*
419 *computer vision and pattern recognition*, pp. 9729–9738, 2020.

- 420 Zhenliang He, Wangmeng Zuo, Meina Kan, Shiguang Shan, and Xilin Chen. Attgan: Facial attribute
421 editing by only changing what you want. *IEEE transactions on image processing*, 28(11):5464–
422 5478, 2019.
- 423 Weihua Hu, Gang Niu, Issei Sato, and Masashi Sugiyama. Does distributionally robust supervised
424 learning give robust classifiers? In *International Conference on Machine Learning*, pp. 2029–
425 2037. PMLR, 2018.
- 426 Jungseock Joo and Kimmo Kärkkäinen. Gender slopes: Counterfactual fairness for computer vision
427 models by attribute manipulation. In *Proceedings of the 2nd international workshop on fairness,
428 accountability, transparency and ethics in multimedia*, pp. 1–5, 2020.
- 429 Sangwon Jung, Sanghyuk Chun, and Taesup Moon. Learning fair classifiers with partially annotated
430 group labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
431 Recognition*, pp. 10348–10357, 2022.
- 432 Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron
433 Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in Neural
434 Information Processing Systems*, 33:18661–18673, 2020.
- 435 Matthäus Kleindessner, Pranjal Awasthi, and Jamie Morgenstern. Fair k-center clustering for data
436 summarization. In *International Conference on Machine Learning*, pp. 3448–3457. PMLR, 2019.
- 437 Alexander Kolesnikov, Xiaohua Zhai, and Lucas Beyer. Revisiting self-supervised visual repre-
438 sentation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern
439 recognition*, pp. 1920–1929, 2019.
- 440 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-
441 lutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.
- 442 Peizhao Li, Han Zhao, and Hongfu Liu. Deep fair clustering for visual learning. In *Proceedings of
443 the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9070–9079, 2020.
- 444 Ming Liu, Yukang Ding, Min Xia, Xiao Liu, Errui Ding, Wangmeng Zuo, and Shilei Wen. Stgan:
445 A unified selective transfer network for arbitrary image attribute editing. In *Proceedings of the
446 IEEE/CVF conference on computer vision and pattern recognition*, pp. 3673–3682, 2019.
- 447 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Large-scale celebfaces attributes (celeba)
448 dataset. *Retrieved August*, 15(2018):11, 2018.
- 449 Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard Zemel. The variational fair
450 autoencoder. *arXiv preprint arXiv:1511.00830*, 2015.
- 451 David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. Learning adversarially fair and
452 transferable representations. In *International Conference on Machine Learning*, pp. 3384–3393.
453 PMLR, 2018.
- 454 Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey
455 on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54(6):1–35, 2021.
- 456 Sungho Park, Sunhee Hwang, Dohyung Kim, and Hyeran Byun. Learning disentangled representa-
457 tion for fair facial attribute classification via fairness-aware information alignment. In *Proceed-
458 ings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 2403–2411, 2021.
- 459 Sungho Park, Jewook Lee, Pilhyeon Lee, Sunhee Hwang, Dohyung Kim, and Hyeran Byun. Fair
460 contrastive learning for facial attribute classification. In *Proceedings of the IEEE/CVF Conference
461 on Computer Vision and Pattern Recognition*, pp. 10389–10398, 2022.
- 462 Edward Raff and Jared Sylvester. Gradient reversal against discrimination: A fair neural network
463 learning approach. In *2018 IEEE 5th International Conference on Data Science and Advanced
464 Analytics (DSAA)*, pp. 189–198. IEEE, 2018.
- 465 Vikram V Ramaswamy, Sunnie SY Kim, and Olga Russakovsky. Fair attribute classification through
466 latent space de-biasing. In *Proceedings of the IEEE/CVF conference on computer vision and
467 pattern recognition*, pp. 9301–9310, 2021.

- 468 Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust
469 neural networks for group shifts: On the importance of regularization for worst-case generaliza-
470 tion. *arXiv preprint arXiv:1911.08731*, 2019.
- 471 Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What
472 makes for good views for contrastive learning? *Advances in Neural Information Processing*
473 *Systems*, 33:6827–6839, 2020.
- 474 Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine*
475 *learning research*, 9(11), 2008.
- 476 Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu,
477 Yadong Mu, Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learn-
478 ing for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 43
479 (10):3349–3364, 2020a.
- 480 Tengfei Wang, Yong Zhang, Yanbo Fan, Jue Wang, and Qifeng Chen. High-fidelity gan inversion
481 for image attribute editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
482 *Pattern Recognition*, pp. 11379–11388, 2022.
- 483 Zeyu Wang, Klint Qinami, Ioannis Christos Karakozis, Kyle Genova, Prem Nair, Kenji Hata, and
484 Olga Russakovsky. Towards fairness in visual recognition: Effective strategies for bias mitigation.
485 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.
486 8919–8928, 2020b.
- 487 Tete Xiao, Xiaolong Wang, Alexei A Efros, and Trevor Darrell. What should not be contrastive in
488 contrastive learning. *arXiv preprint arXiv:2008.05659*, 2020.
- 489 Xiaoying Xing, Hongfu Liu, Chen Chen, and Jundong Li. Fairness-aware unsupervised feature
490 selection. In *Proceedings of the 30th ACM International Conference on Information & Knowledge*
491 *Management*, pp. 3548–3552, 2021.
- 492 Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations.
493 In *International conference on machine learning*, pp. 325–333. PMLR, 2013.
- 494 Tao Zhang, Jing Li, Mengde Han, Wanlei Zhou, Philip Yu, et al. Fairness in semi-supervised learn-
495 ing: Unlabeled data help to reduce discrimination. *IEEE Transactions on Knowledge and Data*
496 *Engineering*, 2020a.
- 497 Tao Zhang, Tianqing Zhu, Mengde Han, Jing Li, Wanlei Zhou, and Philip S Yu. Fairness constraints
498 in semi-supervised learning. *arXiv preprint arXiv:2009.06190*, 2020b.
- 499 Zhifei Zhang, Yang Song, and Hairong Qi. Age progression/regression by conditional adversarial
500 autoencoder. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
501 pp. 5810–5818, 2017.
- 502 Han Zhao and Geoff Gordon. Inherent tradeoffs in learning fair representations. *Advances in neural*
503 *information processing systems*, 32, 2019.
- 504 Jianggang Zhu, Zheng Wang, Jingjing Chen, Yi-Ping Phoebe Chen, and Yu-Gang Jiang. Balanced
505 contrastive learning for long-tailed visual recognition. In *Proceedings of the IEEE/CVF Confer-*
506 *ence on Computer Vision and Pattern Recognition*, pp. 6908–6917, 2022.