

Fair Contrastive Learning for Facial Attribute Classification

Sungho Park¹ Jewook Lee¹ Pilhyeon Lee¹ Sunhee Hwang² Dohyung Kim³ Hyeran Byun^{1*}

¹Yonsei University ²LG Uplus ³SK Inc.

✉ qkrtjdgh18@yonsei.ac.kr

Abstract

Learning visual representation of high quality is essential for image classification. Recently, a series of contrastive representation learning methods have achieved preeminent success. Particularly, SupCon [25] outperformed the dominant methods based on cross-entropy loss in representation learning. However, we notice that there could be potential ethical risks in supervised contrastive learning. In this paper, we for the first time analyze unfairness caused by supervised contrastive learning and propose a new Fair Supervised Contrastive Loss (FSCL) for fair visual representation learning. Inheriting the philosophy of supervised contrastive learning, it encourages representation of the same class to be closer to each other than that of different classes, while ensuring fairness by penalizing the inclusion of sensitive attribute information in representation. In addition, we introduce a group-wise normalization to diminish the disparities of intra-group compactness and inter-class separability between demographic groups that arouse unfair classification. Through extensive experiments on CelebA and UTK Face, we validate that the proposed method significantly outperforms SupCon and existing state-of-the-art methods in terms of the trade-off between top-1 accuracy and fairness. Moreover, our method is robust to the intensity of data bias and effectively works in incomplete supervised settings. Our code is available at <https://github.com/sungho-CoolG/FSCL>.

1. Introduction

Learning powerful visual representation is important for reliable performance in image classification. For a long time, most work has relied on cross-entropy loss to learn the representation due to its strong performance [7, 15, 45, 55]. Meanwhile, recent studies based on contrastive learning have been bringing a new paradigm for representation learning [5, 14, 16, 48, 54]. They effectively learn visual representation by drawing positive pairs and pushing away negative

ones in the high-dimensional space. Despite being originally introduced for unsupervised learning, the contrastive learning strategy proves to be effective in various vision fields [21, 29, 30]. Particularly, SupCon [25] achieves better top-1 accuracy than the state-of-the-art methods based on the cross-entropy loss on ImageNet [41] by simply grafting the contrastive loss to the supervised representation learning.

In this paper, we point out that the contrastive loss may pose potential ethical risks. Despite exhibiting strong performance, it has been underexplored in consideration of fairness which means that the outputs from a model should not be discriminatory in terms of sensitive attributes, such as ethnicity, gender, and age. It is a crucial ethical topic and should be diagnosed in order for the model to be leveraged in the real world [8, 33]. To this end, we analyze the representative contrastive learning model (SupCon) from two major perspectives causing unfairness and propose a new contrastive loss to address both of them.

Learning sensitive attribute information is one of the principal causes of unfairness [6, 19, 32, 46]. It incurs unfair classification by inducing a classifier to determine a decision boundary based on undesirable grounds (*i.e.*, sensitive attributes) [26, 36]. From this point of view, we demonstrate that learning sensitive attribute information leads to the decrease of SupCon loss on the biased dataset, although the desired behavior is to exclusively learn target class information. Consequently, a model learns both kinds of information to minimize the loss, which eventually aggravates unfairness.

To solve the problem, we propose a Fair Supervised Contrastive Loss (FSCL) which prevents encoder networks from learning sensitive attribute information. Basically, it inherits the spirit of supervised contrastive learning that encourages an anchor to be more similar to samples of the same class (*i.e.*, positive samples) than those of other classes (*i.e.*, negative samples). Simultaneously, we limit negative samples to only those having the same sensitive attribute with the anchor among them. In this way, we ensure that learning sensitive attribute information no longer helps the contrastive learning. Rather, it hinders optimizing the loss by increasing the similarity between the anchor and negative samples.

*Corresponding author.

On top of that, we analyze *SupCon* in terms of data imbalance between demographic groups, which is another causal factor of unfairness [39]. Concretely, we identify that the imbalanced number of anchors and positive samples between the demographic groups encourages the *SupCon* loss to put more weight on majority groups. As a result, samples from the majority groups generally have higher similarity to the other samples within the same group and lower similarity to samples having different target classes compared to those from the minority groups. We call the former intra-group compactness and the latter inter-class separability. Since their disparities between the groups result in imbalanced classification performances [13, 57, 59], we introduce a group-wise normalization that reduces the gaps by balancing the loss based on the cardinality of anchors and positive samples between the groups. In the experiments, we demonstrate that it further improves fairness with little damage to the classification performance.

To validate the effectiveness of our method, we perform facial attribute classification on CelebA [31] and UTK Face [58] datasets. In various scenarios, the proposed method significantly ameliorates fairness over *SupCon* and outperforms the state-of-the-art methods in terms of the trade-off between classification accuracy and fairness. Besides, our method is robust to the intensity of data bias and effectively improves fairness even in incompletely supervised settings (e.g., without target class labels or with only a few sensitive attribute labels). Furthermore, we show the extensibility of our method to general bias mitigation through experiments on Dogs and Cats dataset [23].

Main contributions. Our main contributions are summarized as follows. 1) We analyze the causes of unfairness in contrastive learning and propose a Fair Supervised Contrastive Loss that improves fairness by penalizing the inclusion of sensitive attribute information in representation. 2) We introduce a group-wise normalization, which mitigates the group-wise disparities of intra-group compactness and inter-class separability that exacerbate unfairness of representation. 3) Through extensive experiments, we validate that our method learns fair representation under various environments. It achieves the best trade-off performances between top-1 accuracy and fairness on CelebA and UTK Face.

2. Related Work

2.1. Fair Representation Learning

Several studies [26, 38, 51, 56] tried to learn fair representation through adversarial training. They adversarially train the encoder network and the classification head for sensitive attributes so that the encoder network is agnostic to sensitive attribute information. [38] learned fair representation by reversing gradients of classification loss for sensitive attributes through gradient reversal layer (GRL). [26]

further minimized the mutual information between representation and sensitive attribute labels to eliminate their correlations. [51, 56] designed structures in which the outputs from the classification head for target classes are fed into that for sensitive attributes. The latter head removes bias for sensitive attributes in the intermediate outputs through GRL.

Disentangled representation learning [6, 37, 43] is another mainstream for fair representation learning. [43] enforced two types of representation respectively for target classes and sensitive attributes to be orthogonal to each other by maximizing the entropy of the opposite information in each representation. [6] leveraged the disentanglement loss [27] to separate the representation space into sensitive latents and non-sensitive latents without target class labels. Both methods improved fairness by discarding representation containing sensitive attribute information in downstream classification. Moreover, [37] pointed out the shortcoming of [6] that information related to both target and sensitive attributes is contained in sensitive latents and discarded. They introduced an additional subspace for the intersected information.

Recently, [22] made a fresh attempt to improve fairness without compromising performance through fair knowledge distillation. Based on MMD [11], they encourage the feature distribution of the student model conditioned by sensitive attributes to get close to that of the teacher model averaged over the sensitive attributes. With an oversampling strategy, they ameliorate both classification accuracy and fairness on the balanced test set. Meanwhile, [39] proposed a perturbation method which decorrelates the target and sensitive attributes in the latent space of a pre-trained GAN. Then, they generate a balanced dataset with it and utilized the dataset for a fair training of a classification network.

2.2. Contrastive Representation Learning

Contrastive learning [5, 14, 35, 48, 54] has become a dominant approach to learning visual representation in a self-supervised manner. Without class labels, they learned outstanding representation by pulling samples from the same image together and pushing away those from different images. [35, 48, 54] indicated that the number of negative samples is important for contrastive learning and introduced memory banks to increase it without exploding GPU memory consumption. To solve the inconsistency problem between the updated encoder networks and outdated memory bank, [14] utilized a dynamic memory queue as a memory bank and updated it with a slowly moving momentum encoder. Furthermore, [5] proposed a simple architecture for contrastive learning (i.e., *SimCLR*) that outperforms previous methods without the memory bank and specialized architectures. Based on it, [25] proposed a supervised version of contrastive loss (i.e., *SupCon*). Unlike the previous methods, they set all samples having the same class with an anchor to positive samples and pull them to the anchor.

3. Method

In this section, we first analyze the causes of unfairness in supervised contrastive learning, and then describe the proposed method to solve them. Our method is based on a simple framework for contrastive learning similar to previous works [5, 25]. We note that our key contributions lie in not introducing a specific framework but designing a new general loss for learning fair and informative representation.

3.1. Preliminaries

3.1.1 Overall flow

Assume that we have randomly sampled N data pairs in a batch, $\{x_k, y_k, s_k\}_{k=1\dots N}$. Here, $x_k \in X$, $y_k \in Y$, and $s_k \in S$ respectively denote an input image, its target class label out of N_y classes, and its sensitive attribute label out of N_s classes. Following the prior works [5, 14], we randomly crop each image x_k to generate two independent patches (*i.e.*, views), $\hat{x}_{2k-1}, \hat{x}_{2k} \in \tilde{X}$, resulting in a multi-view batch, $\{\hat{x}_l, \hat{y}_l, \hat{s}_l\}_{l=1\dots 2N}$, where $\hat{y}_{2k-1} = \hat{y}_{2k}$ and $\hat{s}_{2k-1} = \hat{s}_{2k}$ for $k \in [1, N]$. An encoder network $\mathcal{F}(\cdot)$ maps the image patches into representation $H = \{h_l\}_{l=1\dots 2N}$, then a projection network $\mathcal{G}(\cdot)$ in turn maps h_l into another representation $Z = \{z_l\}_{l=1\dots 2N}$ for contrastive learning. The encoding networks (*i.e.*, \mathcal{F} and \mathcal{G}) are jointly optimized with contrastive objectives and this process is called representation learning.

After the representation learning process, we freeze the encoder network and throw away the projection network. The frozen encoder network produces representation h_k from the input image x_k instead of the cropped views. Taking the representation as input, a classifier is trained to predict the target class label y_k using the cross-entropy loss.

3.1.2 Self-supervised and supervised contrastive losses

Both self-supervised and supervised contrastive loss enforce an anchor to be more similar to positive samples than negative samples. The major difference is the way positive and negative samples are defined. In the self-supervised version [5], for an anchor \hat{x}_i , the other view from the same image is defined as the positive sample. Meanwhile, in the supervised version [25], all patches sharing the same target class labels with the anchor \hat{x}_i are assigned to positive samples, *i.e.*, $\tilde{X}_p(i) = \{\hat{x}_p \in \tilde{X} | \hat{y}_p = \hat{y}_i, \hat{x}_p \neq \hat{x}_i\}$. In both settings, patches that are neither positive samples nor the anchor are set to negative samples, *i.e.*, $\tilde{X}_n(i) = \{\hat{x}_n \in \tilde{X} | \hat{x}_n \notin \tilde{X}_p(i), \hat{x}_n \neq \hat{x}_i\}$.

In the latent space of z_l , the self-supervised loss maximizes the log-softmax of the similarity between z_i and z_p for the similarity between z_i and representation of all the other samples, $\tilde{X}_a(i) = \tilde{X}_p(i) \cup \tilde{X}_n(i)$. The supervised loss calculates the normalized summation of the multiple log-softmax for all z_p and maximizes it. The self-supervised contrastive

loss (L^{SS}) and supervised contrastive loss (L^{Sup}) are formulated as follows.

$$L^{SS} = - \sum_{z_i \in Z} \log \frac{\phi_p}{\sum_{z_a \in Z_a(i)} \phi_a}, \quad (1)$$

$$L^{Sup} = - \sum_{z_i \in Z} \frac{1}{|Z_p(i)|} \sum_{z_p \in Z_p(i)} \log \frac{\phi_p}{\sum_{z_a \in Z_a(i)} \phi_a}, \quad (2)$$

where ϕ_x denotes $\exp(z_i \cdot z_x / \tau)$, $x \in \{a, p\}$. τ is a temperature parameter, which is set to lower than 1 for sharper distribution of the softmax scores. $|Z_p(i)|$ is the number of positive samples for an anchor z_i . In L^{Sup} , the cardinality of positive samples varies from anchor to anchor and the factor $\frac{1}{|Z_p(i)|}$ normalizes it. The multiple positive samples and normalization factor ensure that L^{Sup} achieves better classification performances than L^{SS} .

3.2. Unfairness in Supervised Contrastive Loss

3.2.1 Learning of sensitive attribute information

As revealed in the literature [6, 32, 46], learning sensitive attribute information is one of the key factors causing unfair classification. Therefore, to analyze the unfairness of L^{Sup} , we start by exploring whether the loss encourages encoder networks to learn the malignant information.

Specifically, we define learning of sensitive attribute information as increasing $I(Z; S) = \mathbb{E}_{P(z,s)} \log \frac{P(z,s)}{P(z)P(s)}$ [6, 26], which is mutual information between Z and S . Subsequently, we suppose two random points, t_l and t_m , in training time, where $I(Z; S)$ is higher at t_m than at t_l (Assumption 1). In addition, to simplify a wide variety of and complicated data bias, we defined an ideally biased dataset, $\{\tilde{X}, \tilde{Y}, \tilde{S}\}$, where each target attribute is correlated with one different sensitive attribute in equal intensity. We provide further details on it in Appendix. Then we demonstrate that L^{Sup} will lead the encoding networks to learn sensitive attribute information by proving the theorem below.

Theorem 1 Given \tilde{X}, \tilde{Y} , and \tilde{S} , for all $t_l, t_m, V^{t_l} > V^{t_m}$.

Here, V^{t_l} and V^{t_m} denote the values of L^{Sup} at t_l and t_m , respectively. Theorem 1 represents that the value of L^{Sup} is always larger at t_l than at t_m . In other words, L^{Sup} is inversely proportional to the $I(Z; S)$. Therefore, it results in the following Corollary. Due to the space limit, we provide the mathematical proof in Appendix.

Corollary 1 Learning sensitive attribute information decreases L_{sup} , given \tilde{X}, \tilde{Y} , and \tilde{S} .

In conclusion, both learning the target attribute and sensitive attribute information reduce L^{Sup} . Since the encoding networks do not have the intrinsic ability to distinguish them, they will learn both kinds of information to optimize L^{Sup} , which eventually aggravates unfairness.

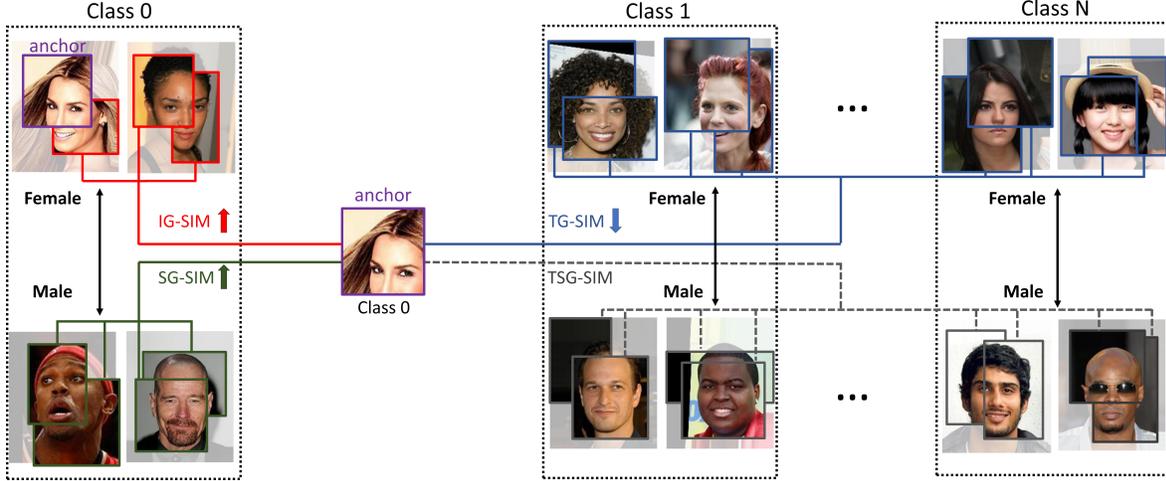


Figure 1. **The concept of fair supervised contrastive loss (FSCL).** It increases the similarity between an anchor and samples of the same target class (IG-SIM and SG-SIM). On the contrary, it decreases the similarity between the anchor and samples having different target classes and the same sensitive attribute, *e.g.*, gender, (TG-SIM). Here TSG-SIM is not directly affected by the loss.

3.2.2 Data imbalance between demographic groups

In addition to discovering the malignant information learning, we explore another cause of unfairness with respect to the imbalanced number of data between data groups. Concretely, we define a data group as a bundle of data having common target classes and sensitive attributes. Based on the group definition, we reformulate L^{Sup} as follows.

$$L^{Sup} = - \sum_{\forall j,k} \sum_{z_i \in Z^{j,k}} \frac{1}{|Z_p(i)|} \sum_{\forall k} \sum_{z_p \in Z_p^k(i)} \log \frac{\phi_p}{\sum_{z_a \in Z_a(i)} \phi_a}, \quad (3)$$

where $Z^{j,k} = \{z_i \in Z | \hat{y}_i = j, \hat{s}_i = k\}$ and $Z_p^k(i) = \{z_p \in Z_p(i) | \hat{s}_p = k\}$ for $j \in [1, N_y]$, and $k \in [1, N_s]$. We note that Eq. 2 and 3 are mathematically identical, but the latter reveals that the imbalanced number of anchors between $Z^{j,k}$ and that of positive samples between $Z_p^k(i)$ are not normalized by the existing factor $1/|Z_p(i)|$. It results in the loss putting more weight on the majority groups, and thus the loss encourages the majority groups to have better intra-group compactness and inter-class separability compared to the minority groups. Consequently, as indicated in [13, 57, 59], their disparities between the groups incur unfair classification performances.

3.3. Fair Supervised Contrastive Loss

To resolve the problem of learning sensitive attribute information (Sec. 3.2.1), we propose a Fair Supervised Contrastive Loss (FSCL) which explicitly penalizes that the encoding networks learn the unwanted information. For brief

and clear explanations, we define the following similarities in consideration of the relationship between an anchor and the other samples.

- **IG-SIM** (Intra-Group Similarity) is the similarity between an anchor and samples within the same group, which have the *same* target class and the *same* sensitive attribute with the anchor. The sample set is defined as $Z_{ig}(i) = \{z_{ig} \in Z_p(i) | \hat{s}_{ig} = \hat{s}_i, \hat{y}_{ig} = \hat{y}_i\}$.
- **SG-SIM** (Sensitive Inter-Group Similarity) is the similarity between an anchor and samples that belong to groups with the *same* target class and *different* sensitive attributes with the anchor. The sample set is defined as $Z_{sg}(i) = \{z_{sg} \in Z_p(i) | \hat{s}_{sg} \neq \hat{s}_i, \hat{y}_{sg} = \hat{y}_i\}$.
- **TG-SIM** (Target Inter-Group Similarity) is the similarity between an anchor and samples that belong to groups with *different* target classes and the *same* sensitive attribute with the anchor. The sample set is defined as $Z_{tg}(i) = \{z_{tg} \in Z_n(i) | \hat{s}_{tg} = \hat{s}_i, \hat{y}_{tg} \neq \hat{y}_i\}$.
- **TSG-SIM** (Target & Sensitive Inter-Group Similarity) is the similarity between an anchor and samples that belong to groups with *different* target classes and *different* sensitive attributes with the anchor. The sample set is defined as *i.e.*, $Z_{tsg}(i) = \{z_{tsg} \in Z_n(i) | \hat{s}_{tsg} \neq \hat{s}_i, \hat{y}_{tsg} \neq \hat{y}_i\}$.

Our key idea is to define the negative sample set as the samples with the same sensitive attributes and different target classes with the anchor (*i.e.*, Z_{tg}). Based on this, we design FSCL that encourages IG-SIM and SG-SIM to be higher than TG-SIM, as illustrated in Figure 1. It is formulated as follows.

$$FSCL = - \sum_{z_i \in Z} \frac{1}{|Z_p(i)|} \sum_{z_p \in Z_p(i)} \log \frac{\phi_p}{\sum_{z_{tg} \in Z_{tg}(i)} \phi_{tg}}, \quad (4)$$

where $Z_p(i) = Z_{ig}(i) \cup Z_{sg}(i)$ and $\phi_{tg} = \exp(z_i \cdot z_{tg} / \tau)$.

On a case-by-case basis, we explain how our *FSCL* addresses the problem of learning unwanted information. In a case of $z_p \in Z_{ig}(i)$, the positive samples and negative samples (*i.e.*, $Z_{tg}(i)$) all have the same sensitive attributes with the anchor. Therefore, the encoding networks no longer consider the sensitive attribute information to be a valuable feature for contrasting an anchor with the negative samples more than with the positive samples.

In the other case of $z_p \in Z_{sg}(i)$, the positive samples have different sensitive attributes from the anchor and negative samples (*i.e.*, $Z_{tg}(i)$). If the encoding networks learn the sensitive attribute information, the similarity between the positive samples and the anchor (*i.e.*, ϕ_p) will decrease and the similarity between the negative samples and it (*i.e.*, ϕ_{tg}) will increase, which is contrary to the objective of the loss. As a result, minimizing the loss inhibits learning the unwanted information in this case.

3.4. Group-wise Normalization

As aforementioned in Sec. 3.2.2, the imbalanced number of anchors and positive samples between data groups causes the group-wise disparities in terms of intra-group compactness and inter-class separability. To alleviate the unfairness brought by the disparities, we introduce group-wise normalization as follows.

$$FSCL+ = - \sum_{\forall j,k} \frac{1}{|Z^{j,k}|} \sum_{z_i \in Z^{j,k}} \sum_{\forall k} \frac{1}{|Z_p^k(i)|} \sum_{z_p \in Z_p^k(i)} \log \frac{\phi_p}{\sum_{z_{tg} \in Z_{tg}(i)} \phi_{tg}}, \quad (5)$$

where $1/|Z^{j,k}|$ and $1/|Z_p^k(i)|$ are the group-wise normalization factors. Different from the existing factor in L^{Sup} (*i.e.*, $1/|Z_p(i)|$), they normalize the cardinality of anchors and positive samples within each group. On an experimental basis, we demonstrate that the proposed normalization mitigates the group-wise imbalances in terms of intra-group compactness and inter-class separability.

4. Experiment

4.1. Datasets

CelebA [31] contains about 200k facial images with 40 binary attribute annotations. We set *male* (*m*) and *young* (*y*) to sensitive attributes and select target attributes having the highest Pearson correlation with the sensitive attributes [3, 49]. Amongst, we manually excluded the extremely correlated attributes for reliable evaluation. For

heavy-makeup as example, there are only 22 males with heavy-makeup in test set. As a result, we exploit three single target attributes: *attractiveness* (*a*), *bignose* (*b*), and *bags-under-eyes* (*e*) as well as two pairs of target attributes: $\{\textit{bignose}, \textit{bags-under-eyes}\}$ and $\{\textit{attractiveness}, \textit{mouth-slightly-open}\}$.

UTK Face [58] consists of about 20k facial images with three kinds of annotations: *gender*, *age*, and *ethnicity*. To evaluate fairness in varied levels of data imbalance, we design several imbalanced versions for the training set. Note that the standard protocol on data splits is not provided in this dataset. Concretely, we set *age* and *ethnicity* to the sensitive attributes and *gender* to the target attribute. *Age* and *ethnicity* are reformed to binary attributes based on whether *age* is under 35 or not and *ethnicity* is Caucasian or not, respectively. A sensitive group (*e.g.*, Caucasian) has male data α times as much as female data and the other sensitive group has the opposite gender ratio. α is set to 2, 3, and 4 to simulate varying bias levels. Unlike the training set, we organize completely balanced validation and test sets for a fair evaluation.

Dogs and Cats [23] has 38,500 dog or cat images. In addition to the original species labels (dog or cat), LNL [26] further annotated color labels (bright or dark). We set *color* to the sensitive attribute and *species* to the target attribute. We compose a *color* biased training set that contains 5 times more black cats than white cats, while 5 times more white dogs than black dogs. For a fair evaluation, we compose the test set to be completely balanced. Note that we utilize this dataset to examine the extensibility of the proposed method to general bias mitigation (*i.e.*, color) beyond its fairness.

We provide more details on the datasets in Appendix.

4.2. Fairness Metrics

A variety of fairness notions are exploited to measure fairness in classification tasks (*e.g.*, demographic parity [28], equal opportunity, and equalized odds [12]). Demographic parity means that the proportion of positive outcomes in each sensitive group should be equal. Although it may be used as reliable metrics in situations where equality of outcome has to be guaranteed, there is a drawback in that a classifier should deliberately misclassify some labels to satisfy it if the proportion of positive outcomes is not equal in the ground truth (GT) [9, 12]. Equal opportunity solves this issue by pursuing the equal true positive rate (TPR) between sensitive groups. However, it does not address unfairness in negative outcomes. In many real-world applications such as facial attribute classification, fairness of positive and negative outcomes is equivalently important. Therefore, equalized odds, which demands both the equal TPR and false positive rate (FPR), are the most suitable to measure fairness in our experiments. Following the definition in [12], we measure the degree of equalized odds (EO) in various settings (*e.g.*, mul-

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=a & o / S=m		T=e & b / S=m		T=a / S=m & y	
	EO	Acc.	EO	Acc.	EO	Acc.	EO	Acc.	EO	Acc.	EO	Acc.	EO	Acc.	EO	Acc.	EO	Acc.
<i>CE</i> [15]	27.8	79.6	16.8	79.8	17.6	84.0	14.7	84.5	15.0	83.9	12.7	83.8	30.0	73.9	12.9	72.6	31.3	79.5
<i>GRL</i> [38]	24.9	77.2	14.7	74.6	14.0	82.5	10.0	83.3	6.7	81.9	5.9	82.3	17.8	73.1	9.4	71.4	22.9	78.6
<i>LNL</i> [26]	21.8	79.9	13.7	74.3	10.7	82.3	6.8	82.3	5.0	81.6	3.3	80.3	16.7	72.9	7.4	70.8	20.7	77.7
<i>FD-VAE</i> [37]	15.1	76.9	14.8	77.5	11.2	81.6	6.7	81.7	5.7	82.6	6.2	84.0	18.2	73.4	8.2	70.2	19.9	78.0
<i>MFD</i> [22]	7.4	78.0	14.9	80.0	7.3	78.0	5.4	78.0	8.7	79.0	5.2	78.0	8.7	74.0	9.0	70.0	19.4	76.1
<i>SupCon</i> [25]	30.5	80.5	21.7	80.1	20.7	84.6	16.9	84.4	20.8	84.3	10.8	84.0	22.8	74.0	12.5	72.7	24.4	81.7
<i>FSCL</i>	11.5	79.1	13.0	79.1	7.0	82.1	6.4	83.8	3.8	82.7	1.8	82.0	8.1	74.1	6.8	71.1	19.9	79.4
<i>FSCL+</i>	6.5	79.1	12.4	79.1	4.7	82.9	4.8	84.1	3.0	83.4	1.6	83.5	3.6	74.8	2.5	70.8	17.0	77.2

Table 1. **Classification results on CelebA.** We measure classification accuracy (ACC.) and equalized odds (EO) in various scenarios. Here $a, b, e, o, m,$ and y respectively denote *attractiveness, bignose, bags-under-eyes, mouth-slightly-open, male,* and *young*. On the other hand, T and S represent target and sensitive attributes, respectively. All the results are the averaged scores over three independent runs. The standard deviations are provided in Appendix.

tuple classes or sensitive attributes) as follows.

$$\overline{\sum_{\forall y, c, \{s^0, s^1\} \subset S} |P_{s^0}(C = c | Y = y) - P_{s^1}(C = c | Y = y)|}, \quad (6)$$

where $\overline{\sum}$ is the averaged sum. $y \in Y$ and $c \in C$ are target labels and outputs from a classifier, respectively, and $\{s^0, s^1\}$ is a two-element subset of sensitive attribute groups S .

4.3. Implementation Details

For contrastive learning, we utilize ResNet-18 [15] for the encoder network \mathcal{F} and a MLP with two hidden layers for the projection network \mathcal{G} . The dimensions of latent spaces are set to 256 and 128, respectively. We augment two cropped patches per image following the augmentation strategy in [5] and resize them to 128×128 . We set the temperature parameter τ to 0.1 based on the analysis in [25]. We train the encoding networks for 100 epochs in the representation learning stage, and subsequently train the classifier, which is a MLP with one hidden layer, for 10 epochs using the cross-entropy loss. For the experiments with multiple target or sensitive attributes, we combine multiple binary attribute labels into a multi-class label. All comparative models share the same structures of the encoder network and classifier as ours for a fair comparison. The results reported in this paper are averaged over three independent runs. More details for the augmentation strategy, structure of networks, and experiment settings are provided in Appendix.

4.4. Classification Results on CelebA

Table 1 shows the classification results on CelebA. For diverse combinations of target and sensitive attributes, we measure classification performances and fairness with top-1 accuracy and equalized odds (EO), respectively. In all the experiments, *Cross-Entropy* (CE) [15] and *SupCon* [25] record excellent top-1 accuracy but suffer from severe unfairness. Notably, the proposed methods (*FSCL* and *FSCL+*) significantly improve EO over them while preserving the competitive performances. Particularly, the comparison between

Method	Adversarial Training [38]	EO (\downarrow)	Acc. (\uparrow)
<i>SupCon</i> [25]	\times	$30.5_{\pm 1.3}$	$80.5_{\pm 0.7}$
	\checkmark	$21.0_{\pm 0.5}$	$76.6_{\pm 0.3}$
<i>FSCL+</i>	\times	$6.5_{\pm 0.4}$	$79.1_{\pm 0.1}$
	\checkmark	$9.0_{\pm 0.5}$	$79.2_{\pm 0.1}$

Table 2. **Effect of adversarial training in contrastive learning on CelebA dataset.** We utilize *GRL* [38] for the adversarial training. Here *attractiveness* and *male* are set to the target class and sensitive attribute, respectively.

FSCL (blue) and *FSCL+* (red) shows that the group-wise normalization brings about better fairness while well preserving the performance or even improving it. Furthermore, we compare ours with various state-of-the-art approaches for fairness such as adversarial training (*GRL* [38] and *LNL* [26]), disentangled representation learning (*FD-VAE* [37]), and fair distillation (*MFD* [22]). *FSCL+* substantially outperforms all the state-of-the-art methods in terms of the trade-off between top-1 accuracy and EO in all the settings. For a clearer comparison of the trade-off performances, we also provide the experimental results in figure form in Appendix.

4.5. Adversarial Training in Contrastive Learning

Intuitively, to mitigate the unfairness of *SupCon*, one may imagine simply combining adversarial training with it. In Table 2, we demonstrate the effect of adversarial training by applying *GRL* [38] to *SupCon* and *FSCL+* in the representation learning. For *SupCon*, while improving fairness to an extent, it largely damages the classification performance. In addition, *FSCL+* achieves much better EO and top-1 accuracy than *SupCon* combined with *GRL*. This indicates that the simple graft of adversarial training to contrastive learning does not sufficiently improve fairness and designing a new method seamlessly integrated into contrastive learning is more effective. We do not see further improvements in EO when applying *GRL* to *FSCL+*.

Method	Ramaswamy <i>et al.</i> [39]	EO (\downarrow)	Acc. (\uparrow)
<i>Cross-Entropy</i> [15]	\times	27.8 ± 0.2	79.6 ± 0.5
	\checkmark	24.1 ± 0.5	79.6 ± 0.2
<i>FSCL+</i>	\times	6.5 ± 0.4	79.1 ± 0.1
	\checkmark	4.2 ± 0.4	79.6 ± 0.1

Table 3. **Compatibility with fair data augmentation [39] on CelebA dataset.** We set *attractiveness* and *male* to the target class and sensitive attribute, respectively.

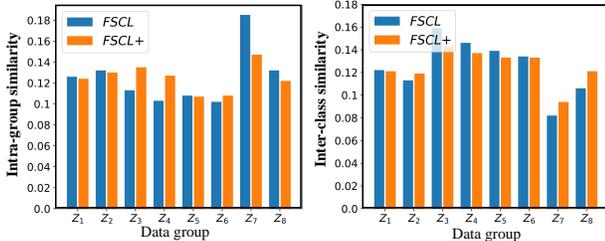


Figure 2. **Effectiveness of group-wise normalization.** The group-wise normalization (*i.e.*, *FSCL+*) significantly mitigates the group-wise disparities in terms of intra-group and inter-class similarities.

4.6. Compatibility with Fair Data Augmentation

We incorporate our method with Ramaswamy *et al.* [39], one of the state-of-the-art pre-processing methods for fair classification. It generates a de-biased dataset through a Progressive GAN [24] and augments the original dataset with the generated one. In Table 3, we report the performances of the baseline (*i.e.*, *Cross-Entropy*) and *FSCL+* trained on the original/augmented dataset. The results show that *FSCL+* outperforms Ramaswamy *et al.* (2nd row) in terms of both EO and top-1 accuracy. Besides, the fairness of ours is further enhanced when adopting the fair data augmentation, which indicates its compatibility.

4.7. Effectiveness of Group-wise Normalization

To analyze the effectiveness of the group-wise normalization, we compare intra-group compactness and inter-class separability between *FSCL* and *FSCL+*. To this end, we first divide the test set into 8 groups with respect to one target class, *attractiveness*, and two sensitive attributes, *male* and *young*, and then calculate them as follows. The former is measured by averaging the similarities between representation within a group (*i.e.*, intra-group similarity) and the latter is measured by averaging the similarity between representation in a group and representation having different class labels with it (*i.e.*, inter-class similarity). For easier comparison, the values are normalized to sum to unity, as shown in Figure 2. The plots demonstrate that the group-wise normalization significantly diminishes the group-wise disparities. In specific, *FSCL* has the standard deviations of 0.084 and 0.031 in intra-group and inter-class similarities,

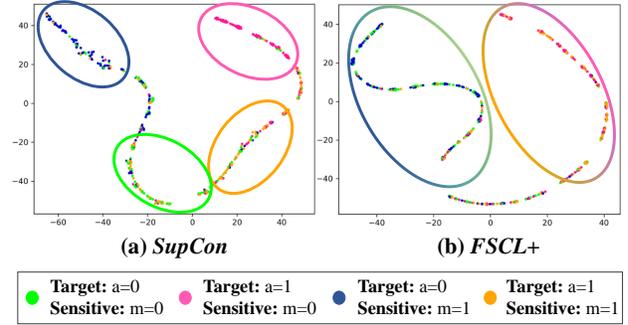


Figure 3. **Qualitative comparison using t-SNE visualizations.** It is clearly shown that *FSCL+* (b) learns representation more independent of the sensitive attribute than *SupCon* (a).

respectively, while *FSCL+* has lower standard deviations of 0.038 and 0.024.

4.8. Qualitative Analysis with t-SNE Visualization

In Figure 3, we provide t-SNE plots [50] of representation from *SupCon* and *FSCL+* on CelebA dataset. The representation is divided into 4 groups in terms of the target class (*i.e.*, *attractiveness*) and sensitive attribute (*i.e.*, *male*), which are visualized in different colors. In *SupCon*, the representation is divided by both the target class and sensitive attribute, suggesting that the encoding networks learn information for the sensitive attribute as well as the target class. Consequently, the representation of minority groups (*i.e.*, green and orange colors) is more similar to the representation of the counterpart class than that of majority groups (*i.e.*, blue and pink colors). In contrast, in *FSCL+*, the representation is divided by only the target class, that is, it is more agnostic to the sensitive attribute. Accordingly, majority groups can no longer have more privileges than minority groups, which explains why our loss can achieve fairer performance than *SupCon* in image classification. Details of experimental settings are provided in Appendix.

4.9. Robustness to Severity of Data Bias

In Figure 4, we present the trends of EO and top-1 accuracy according to the intensity shift of data imbalance (α) on UTK Face dataset. It can be clearly noticed that our loss best prevents the degradation of fairness caused by an increase in α , achieving the fairest performance at all the intensities. In the figure, as α increases, the EO gaps between ours and the others become larger, which manifests the robustness of the proposed methods against the severity of data bias. Moreover, at all the intensities, our loss successfully maintains the top-1 accuracy, which is close to *SupCon*. Experimental results on another sensitive attribute (*i.e.*, *age*) draw similar conclusions and are provided in Appendix.

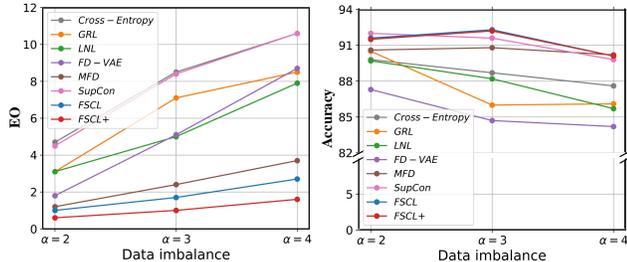


Figure 4. **Classification results on UTK Face.** We measure classification accuracy and equalized odds (EO) under different severity levels of imbalance (α). We set *gender* and *ethnicity* to the target class and sensitive attribute, respectively. Larger α indicates that the training set is more imbalanced.

Method	Target: <i>a</i> / Sensitive: <i>m</i>		Target: <i>b</i> / Sensitive: <i>m</i>	
	EO (\downarrow)	Acc. (\uparrow)	EO (\downarrow)	Acc. (\uparrow)
<i>SimCLR</i> [5]	29.4 \pm 2.5	75.7 \pm 0.2	16.4 \pm 0.4	82.0 \pm 0.1
<i>SimCLR</i> [5] + <i>GRL</i> [38]	21.9 \pm 0.9	72.3 \pm 0.4	13.7 \pm 0.3	82.3 \pm 0.0
<i>FSCL</i> [†]	14.8 \pm 0.9	74.6 \pm 0.4	6.1 \pm 0.6	80.8 \pm 0.2

Table 4. **Classification results on CelebA in the absence of target class labels during representation learning.** *FSCL*[†] is a modified version of *FSCL* that does not use the target class labels.

4.10. Results in Incomplete Supervised Settings

We explore a more challenging problem setting, where target class labels are unavailable during the representation learning process. To this end, we introduce a modified version of *FSCL* that does not exploit target class labels, which is denoted by *FSCL*[†]. Similar to *SimCLR* [5], it uses only a single positive sample that comes from the same image with an anchor. As shown in Table 4, ours significantly improves fairness at the acceptable cost of top-1 accuracy, compared to *SimCLR* and *SimCLR*+*GRL*. Details of the modification are provided in Appendix.

Moreover, we conduct experiments under another challenging environment where only a small portion of data have sensitive attribute labels. One of our simple strategies to handle this task is to generate pseudo-labels for applying *FSCL*+ loss. Specifically, we train a classifier to predict sensitive attribute labels only with the samples having sensitive attribute labels, and then generate the pseudo labels of sensitive attributes for the other samples. Another strategy is to apply *FSCL*+ loss only to data with sensitive labels and *SupCon* to the other data. Table 5 shows that *FSCL*+ effectively ameliorates EO over *SupCon* even under the incomplete supervision of sensitive attributes. Surprisingly, *FSCL*+ with only 5% of labels is able to outperform *SupCon*+*GRL* using all the labels.

4.11. Extensibility to General Bias Mitigation

To verify the efficacy of the proposed methods in a general bias type, we conduct experiments on Dogs and Cats [23]

Method	# of Sensitive	Pseudo-labeling	EO (\downarrow)	Acc. (\uparrow)
<i>SupCon</i> [25]	0	-	30.5 \pm 1.3	80.5 \pm 0.7
<i>SupCon</i> [25] + <i>GRL</i> [38]	1	-	21.0 \pm 0.5	76.6 \pm 0.3
<i>FSCL</i> +	1	-	6.5 \pm 0.4	79.1 \pm 0.1
	1/2	✗	13.4 \pm 0.1	79.3 \pm 0.3
		✓	12.8 \pm 1.2	79.4 \pm 0.3
	1/4	✗	18.7 \pm 0.3	80.0 \pm 0.3
		✓	13.4 \pm 0.1	79.5 \pm 0.5
	1/10	✗	20.7 \pm 0.5	80.2 \pm 0.1
	✓	16.5 \pm 0.5	79.6 \pm 0.4	
1/20	✗	23.4 \pm 0.0	80.6 \pm 0.1	
	✓	18.8 \pm 1.1	78.5 \pm 0.2	

Table 5. **Classification results on CelebA under incomplete supervision of sensitive attribute labels.** “# of Sensitive” denotes the ratio of data having sensitive attribute labels. We set *attractiveness* and *male* to the target class and sensitive attribute, respectively.

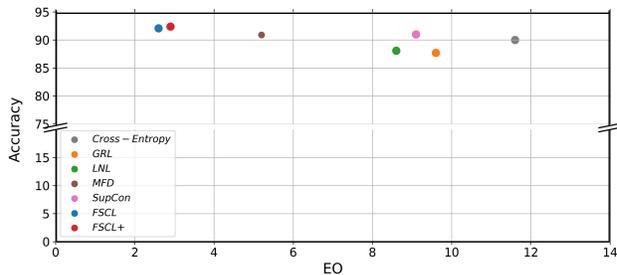


Figure 5. **Classification results on Dogs and Cats.** We set *species* and *color* to the target and sensitive attributes, respectively.

with *color* bias. The results are shown in Figure 5, where our models (*FSCL* and *FSCL*+) best eliminate the color bias, which implies that they are generalizable to various types of bias. Note that *FSCL*, *FSCL*+, and *MFD* show higher top-1 accuracy than their baselines since fairness improves the performance on the completely balanced test set.

5. Conclusion

In this paper, we addressed fairness in contrastive learning. We first analyzed the causative factors of unfairness in the supervised contrastive loss. Then we proposed the fair supervised contrastive loss and introduced the group-wise normalization into the loss. Through extensive experiments, we validated that our loss effectively improves fairness with little degradation of the classification performance.

Acknowledgements

This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No. 2022R1A2B5B02001467) and Institute for Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No. 2019-0-01396: Development of Framework for Analyzing, Detecting, Mitigating of Bias in AI model and Training Data, No. 2020-0-01361: Artificial Intelligence Graduate School Program (Yonsei University)).

References

- [1] Martim Brandao. Age and gender bias in pedestrian detection algorithms. *arXiv e-prints*, page arXiv:1906.10490, June 2019. [19](#)
- [2] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In Sorelle A. Friedler and Christo Wilson, editors, *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, volume 81 of *Proceedings of Machine Learning Research*, pages 77–91, New York, NY, USA, 23–24 Feb 2018. PMLR. [19](#)
- [3] Luigi Celona, S. Bianco, and R. Schettini. Fine-grained face annotation using deep multi-task cnn. *Sensors (Basel, Switzerland)*, 18, 2018. [5](#)
- [4] Bor-Chun Chen, Chu-Song Chen, and Winston H. Hsu. Cross-age reference coding for age-invariant face recognition and retrieval. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *Computer Vision – ECCV 2014*, pages 768–783, Cham, 2014. Springer International Publishing. [16](#)
- [5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR, 13–18 Jul 2020. [1](#), [2](#), [3](#), [6](#), [8](#), [18](#), [19](#)
- [6] Elliot Creager, David Madras, Joern-Henrik Jacobsen, Marissa Weis, Kevin Swersky, Toniann Pitassi, and Richard Zemel. Flexibly fair representation learning by disentanglement. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1436–1445, Long Beach, California, USA, 09–15 Jun 2019. PMLR. [1](#), [2](#), [3](#), [19](#)
- [7] Ekin Dogus Cubuk, Barret Zoph, Jon Shlens, and Quoc Le. Randaugment: Practical automated data augmentation with a reduced search space. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 18613–18624. Curran Associates, Inc., 2020. [1](#)
- [8] Conor Dougherty. Google photos mistakenly labels black people gorillas. Twitter, 2015. [1](#)
- [9] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, ITCS ’12, page 214–226, New York, NY, USA, 2012. Association for Computing Machinery. [5](#)
- [10] Sixue Gong, Xiaoming Liu, and Anil K. Jain. Jointly debiasing face recognition and demographic attribute estimation. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XXIX*, volume 12374 of *Lecture Notes in Computer Science*, pages 330–347. Springer, 2020. [19](#)
- [11] Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *J. Mach. Learn. Res.*, 13(null):723–773, mar 2012. [2](#)
- [12] Moritz Hardt, Eric Price, and Nathan Srebro. Equality of opportunity in supervised learning. In *Proceedings of the 30th International Conference on Neural Information Processing Systems, NIPS’16*, page 3323–3331, Red Hook, NY, USA, 2016. Curran Associates Inc. [5](#)
- [13] Munawar Hayat, Salman Khan, Syed Waqas Zamir, Jianbing Shen, and Ling Shao. Gaussian affinity for max-margin class imbalanced learning. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6468–6478, 2019. [2](#), [4](#)
- [14] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [1](#), [2](#), [3](#)
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [1](#), [6](#), [7](#), [17](#), [18](#)
- [16] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. In *International Conference on Learning Representations*, 2019. [1](#)
- [17] Sunhee Hwang and Hyeran Byun. Unsupervised image-to-image translation via fair representation of gender bias. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1953–1957, 2020. [19](#)
- [18] Sunhee Hwang, Sungho Park, Dohyung Kim, Mirae Do, and Hyeran Byun. Fairfacegan: Fairness-aware facial image-to-image translation. In *BMVC*, volume 2020, 2020. [19](#)
- [19] Sunhee Hwang, Sungho Park, Pilhyeon Lee, Seogkyu Jeon, Dohyung Kim, and Hyeran Byun. Exploiting transferable knowledge for fairness-aware image classification. In *Proceedings of the Asian Conference on Computer Vision (ACCV)*, November 2020. [1](#)
- [20] S. M. Seitz I. Kemelmacher-Shlizerman, S. Suwajanakorn. Illumination-aware age progression. In *CVPR*, 2014. [16](#)
- [21] Seogkyu Jeon, Kibeom Hong, Pilhyeon Lee, Jewook Lee, and Hyeran Byun. Feature stylization and domain-aware contrastive learning for domain generalization. In *The 29th ACM International Conference on Multimedia*, pages 22–31, 2021. [1](#)
- [22] Sangwon Jung, Donggyu Lee, Taeon Park, and Taesup Moon. Fair feature distillation for visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12115–12124, June 2021. [2](#), [6](#), [17](#)
- [23] Kaggle. Dogs vs. cats. 2013. [2](#), [5](#), [8](#)
- [24] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of GANs for improved quality, stability, and variation. In *International Conference on Learning Representations*, 2018. [7](#)

- [25] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 18661–18673. Curran Associates, Inc., 2020. [1](#), [2](#), [3](#), [6](#), [8](#), [17](#), [18](#), [19](#)
- [26] Byungju Kim, Hyunwoo Kim, Kyungsu Kim, Sungjin Kim, and Junmo Kim. Learning not to learn: Training deep neural networks with biased data. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. [1](#), [2](#), [3](#), [5](#), [6](#), [16](#), [17](#)
- [27] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2649–2658, Stockholm, Sweden, 10–15 Jul 2018. PMLR. [2](#)
- [28] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4066–4076. Curran Associates, Inc., 2017. [5](#)
- [29] Pilhyeon Lee and Hyeran Byun. Learning action completeness from points for weakly-supervised temporal action localization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 13648–13657, 2021. [1](#)
- [30] Pilhyeon Lee, Sunhee Hwang, Jewook Lee, Minjung Shin, Seogkyu Jeon, and Hyeran Byun. Inter-subject contrastive learning for subject adaptive eeg-based visual recognition. In *The 10th International Winter Conference on Brain-Computer Interface (BCI)*, 2022. [1](#)
- [31] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015. [2](#), [5](#), [16](#)
- [32] Francesco Locatello, Gabriele Abbati, Thomas Rainforth, Stefan Bauer, Bernhard Schölkopf, and Olivier Bachem. On the fairness of disentangled representations, 2019. [1](#), [3](#)
- [33] Natasha Lomas. Faceapp apologizes for building a racist ai. TechCrunch, 2018. [1](#)
- [34] David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. Learning adversarially fair and transferable representations. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 3384–3393. PMLR, 10–15 Jul 2018. [19](#)
- [35] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [2](#)
- [36] Junhyun Nam, Hyuntak Cha, Sungsoo Ahn, Jaeho Lee, and Jinwoo Shin. Learning from failure: De-biasing classifier from biased classifier. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 20673–20684. Curran Associates, Inc., 2020. [1](#)
- [37] Sungho Park, Sunhee Hwang, Dohyung Kim, and Hyeran Byun. Learning disentangled representation for fair facial attribute classification via fairness-aware information alignment. *Proceedings of AAAI-2021*, 2021. [2](#), [6](#), [16](#), [17](#), [19](#)
- [38] Edward Raff and Jared Sylvester. Gradient reversal against discrimination: A fair neural network learning approach. In *2018 IEEE 5th International Conference on Data Science and Advanced Analytics (DSAA)*, pages 189–198, 2018. [2](#), [6](#), [8](#), [15](#), [17](#)
- [39] Vikram V. Ramaswamy, Sunnie S. Y. Kim, and Olga Russakovsky. Fair attribute classification through latent space de-biasing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9301–9310, June 2021. [2](#), [7](#)
- [40] Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016. [18](#)
- [41] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. [1](#)
- [42] Shiori Sagawa*, Pang Wei Koh*, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks. In *International Conference on Learning Representations*, 2020. [18](#)
- [43] Mhd Hasan Sarhan, Nassir Navab, Abouzar Eslami, and Shadi Albarqouni. Fairness by learning orthogonal disentangled representations. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XXIX*, volume 12374 of *Lecture Notes in Computer Science*, pages 746–761. Springer, 2020. [2](#)
- [44] P. Sattigeri, S. C. Hoffman, V. Chenthamarakshan, and K. R. Varshney. Fairness gan: Generating datasets with fairness properties using a generative adversarial network. *IBM Journal of Research and Development*, 63(4/5):3:1–3:9, 2019. [19](#)
- [45] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. [1](#)
- [46] Jiaming Song, Pratyusha Kalluri, Aditya Grover, Shengjia Zhao, and Stefano Ermon. Learning controllable fair representations. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 2164–2173. PMLR, 16–18 Apr 2019. [1](#), [3](#)
- [47] Yi Sun, Xiaogang Wang, and Xiaoou Tang. Hybrid deep learning for face verification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(10):1997–2009, 2016. [16](#)

- [48] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, pages 776–794, Cham, 2020. Springer International Publishing. [1](#), [2](#), [19](#)
- [49] Robert Torfason, Eiríkur Agustsson, Rasmus Rothe, and Radu Timofte. From face images and attributes to attributes. In *Asian Conference on Computer Vision*, pages 313–329. Springer, 2016. [5](#)
- [50] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, June 2015. [7](#), [16](#)
- [51] Christina Wadsworth, Francesca Vera, and Chris Piech. Achieving fairness through adversarial learning: an application to recidivism prediction. *arXiv preprint arXiv:1807.00199*, 2018. [2](#)
- [52] Tianlu Wang, Jieyu Zhao, Mark Yatskar, Kai-Wei Chang, and Vicente Ordonez. Balanced datasets are not enough: Estimating and mitigating gender bias in deep image representations. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5309–5318, 2019. [19](#)
- [53] Zeyu Wang, Klint Qinami, Ioannis Christos Karakozis, Kyle Genova, Prem Nair, Kenji Hata, and Olga Russakovsky. Towards fairness in visual recognition: Effective strategies for bias mitigation. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8916–8925, 2020. [19](#)
- [54] Zhirong Wu, Yuanjun Xiong, Stella X. Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. [1](#), [2](#), [19](#)
- [55] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V. Le. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. [1](#)
- [56] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating unwanted biases with adversarial learning. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, AIES ’18, page 335–340, New York, NY, USA, 2018. Association for Computing Machinery. [2](#)
- [57] Xiao Zhang, Zhiyuan Fang, Yandong Wen, Zhifeng Li, and Yu Qiao. Range loss for deep face recognition with long-tailed training data. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 5419–5428, 2017. [2](#), [4](#)
- [58] Zhifei Zhang, Yang Song, and Hairong Qi. Age progression/regression by conditional adversarial autoencoder. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. [2](#), [5](#), [16](#)
- [59] He Zhu and Shan Yu. Intra-class uncertainty loss function for classification. *CoRR*, abs/2104.05298, 2021. [2](#), [4](#)

A. Definition of Ideally Biased Dataset

To confine a wide variety of data bias, we first define an ideally biased dataset that satisfies the following conditions.

1. The dataset has m target and sensitive classes (i.e., $N_t = N_s = m$). Each target and sensitive class contains the same number of data.
2. Target classes are biased to sensitive classes with a one-to-one mapping. That is, each target class has only one *biased sensitive class*, and no more than one target class has the same *biased sensitive class*.
3. In each target class, *biased sensitive class* has r times more data than other sensitive classes.
4. Target classes are highly biased to sensitive classes (i.e., $r \geq m^2$).

We illustrated it in Figure 6, where the number of data in non-biased classes is set to C . All the proof below is based on this ideally biased dataset.

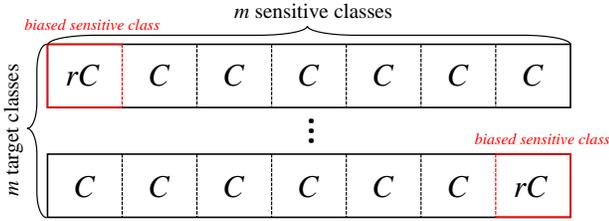


Figure 6. **The composition of the ideally biased dataset.** It has m target and sensitive classes, and each target class has only one *biased sensitive class*. C represents the number of data in non-biased classes and r is the bias for *biased sensitive classes*.

B. Mathematical Proof on Theorem 1

In the main paper, we demonstrated that *SupCon* will lead the encoding networks to learn sensitive attribute information based on Theorem 1. We provide the mathematical proof for the theorem below.

Assumption 1

Let input data come from the ideally biased dataset (refer to Sec. A), where \tilde{X} , \tilde{Y} , \tilde{S} denote input images, target class labels, and sensitive attribute labels, respectively. We note that target classes are highly correlated with sensitive attributes in the dataset ($r \geq m^2$).

Definition 1

Learning of sensitive attribute information indicates an increase of $I(Z; \tilde{S})$, where $I(Z; \tilde{S}) = \mathbb{E}_{P(z, \tilde{s})} \log \frac{P(z, \tilde{s})}{P(z)P(\tilde{s})}$ and Z denotes the visual representation.

Assumption 2

Let t_l, t_m be random points in training time when $I(Z^{t_l}; \tilde{S}) < I(Z^{t_m}; \tilde{S})$.

Axiom 1

Given \tilde{X} , \tilde{Y} , and \tilde{S} , for all z_i , $|Z_p(i)| = Cr + (m - 1)C - 1$, which is a constant.

Definition 2

$$L_a^{sup} = \sum_{z_i \in Z} \sum_{z_p \in Z_p(i)} \left[\log \left(\sum_{z_a \in Z_a(i)} \phi_a \right) \right].$$

Definition 3

$$L_p^{sup} = \sum_{z_i \in Z} \sum_{z_p \in Z_p(i)} \log \phi_p.$$

Proposition 1

$L^{Sup} = \hat{C}(-L_p^{sup} + L_a^{sup})$, where \hat{C} is a constant.

Proof.

$$\begin{aligned} L^{Sup} &= - \sum_{z_i \in Z} \frac{1}{|Z_p(i)|} \sum_{z_p \in Z_p(i)} \log \frac{\phi_p}{\sum_{z_a \in Z_a(i)} \phi_a} \\ &= - \sum_{z_i \in Z} \frac{1}{|Z_p(i)|} \sum_{z_p \in Z_p(i)} \log \phi_p \\ &\quad + \sum_{z_i \in Z} \frac{1}{|Z_p(i)|} \sum_{z_p \in Z_p(i)} \log \sum_{z_a \in Z_a(i)} \phi_a \quad (7) \\ &= \frac{1}{|Z_p(i)|} \left(-L_p^{sup} + L_a^{sup} \right) \\ &= \hat{C} \left(-L_p^{sup} + L_a^{sup} \right) (\because \text{Axiom1}). \end{aligned}$$

Definition 4

Let $V_x^{t_l}$ and $V_x^{t_m}$ be the values of L_x^{sup} , $x \in \{p, a\}$, at t_l and t_m , respectively.

For example, the value of L_a^{sup} at a certain point in training time, t_l , can be represented as:

$$V_a^{t_l} = \sum_{z_i \in Z} \sum_{z_p \in Z_p(i)} \left[\log \left(\sum_{z_a \in Z_a(i)} \phi_a^{t_l} \right) \right], \quad (8)$$

where $\phi_x^{t_k} = \exp(z_i^{t_k} \cdot z_x^{t_k} / \tau)$, $x \in \{p, a\}$, $k \in \{l, m\}$.

Definition 5

Let $Z_x(i) = Z_x^s(i) \cup Z_x^d(i)$, where $Z_x^s(i) = \{z_x \in Z_x(i) | \tilde{s}_x = \tilde{s}_i\}$ and $Z_x^d(i) = \{z_x \in Z_x(i) | \tilde{s}_x \neq \tilde{s}_i\}$, $x \in \{p, a\}$.

Proposition 2

From Definition 2, 3, 4 and 5,

$$\begin{aligned} V_a^{t_k} &= \sum_{z_i \in Z} \sum_{z_p \in Z_p(i)} \left[\log \left(\sum_{z_a \in Z_a(i)} \phi_a^{t_k} \right) \right] \\ &= \sum_{z_i \in Z} \sum_{z_p \in Z_p(i)} \left[\log \left(\sum_{z_a \in Z_a^s(i)} \phi_a^{t_k} + \sum_{z_a \in Z_a^d(i)} \phi_a^{t_k} \right) \right]. \end{aligned} \quad (9)$$

$$\begin{aligned} V_p^{t_k} &= \sum_{z_i \in Z} \left[\sum_{z_p \in Z_p(i)} \log \phi_p^{t_k} \right] \\ &= \sum_{z_i \in Z} \left[\sum_{z_p \in Z_p^s(i)} \log \phi_p^{t_k} + \sum_{z_p \in Z_p^d(i)} \log \phi_p^{t_k} \right]. \end{aligned} \quad (10)$$

Conjecture 1

- a) $\sum_{z_a \in Z_a^s(i)} \phi_a^{t_l} < \sum_{z_a \in Z_a^s(i)} \phi_a^{t_m}$
- b) $\sum_{z_a \in Z_a^d(i)} \phi_a^{t_l} > \sum_{z_a \in Z_a^d(i)} \phi_a^{t_m}$
- c) $\sum_{z_p \in Z_p^s(i)} \log \phi_p^{t_l} < \sum_{z_p \in Z_p^s(i)} \log \phi_p^{t_m}$
- d) $\sum_{z_p \in Z_p^d(i)} \log \phi_p^{t_l} > \sum_{z_p \in Z_p^d(i)} \log \phi_p^{t_m}$

From Assumption 2, $I(Z^{t_l}; \tilde{S}) < I(Z^{t_m}; \tilde{S})$, hence the similarity between z_i and $Z_k^s(i)$ is larger at t_m than t_l . Meanwhile, the similarity between z_i and $Z_k^d(i)$ is smaller at t_m than at t_l .

Proposition 3

Let $\alpha_{z_x}, \beta_{z_x} \in R^+$, $x \in \{p, a\}$, then

$$\begin{aligned} V_a^{t_m} &= \sum_{z_i \in Z} \sum_{z_p \in Z_p(i)} \left[\log \left(\sum_{z_a \in Z_a^s(i)} (1 + \alpha_{z_a}) \phi_a^{t_l} \right. \right. \\ &\quad \left. \left. + \sum_{z_a \in Z_a^d(i)} (1 - \beta_{z_a}) \phi_a^{t_l} \right) \right], \end{aligned} \quad (11)$$

$$\begin{aligned} V_p^{t_m} &= \sum_{z_i \in Z} \left[\sum_{z_p \in Z_p^s(i)} \log(1 + \alpha_{z_p}) \phi_p^{t_l} \right. \\ &\quad \left. + \sum_{z_p \in Z_p^d(i)} \log(1 - \beta_{z_p}) \phi_p^{t_l} \right], \end{aligned} \quad (12)$$

where α_{z_x} is an increasing rate of similarity between an anchor and each sample from t_l to t_m . Conversely, β_{z_x} is the decreasing rate of similarity.

proof.

By Conjecture 1,

$$\begin{aligned} \sum_{z_a \in Z_a^s(i)} \phi_a^{t_m} &= \sum_{z_a \in Z_a^s(i)} (1 + \alpha_{z_a}) \phi_a^{t_l}, \\ \sum_{z_a \in Z_a^d(i)} \phi_a^{t_m} &= \sum_{z_a \in Z_a^d(i)} (1 - \beta_{z_a}) \phi_a^{t_l}, \\ \sum_{z_p \in Z_p^s(i)} \log \phi_p^{t_m} &= \sum_{z_p \in Z_p^s(i)} \log(1 + \alpha_{z_p}) \phi_p^{t_l}, \\ \sum_{z_p \in Z_p^d(i)} \log \phi_p^{t_m} &= \sum_{z_p \in Z_p^d(i)} \log(1 - \beta_{z_p}) \phi_p^{t_l}. \end{aligned} \quad (13)$$

Therefore,

$$\begin{aligned} V_a^{t_m} &= \sum_{z_i \in Z} \sum_{z_p \in Z_p(i)} \left[\log \left(\sum_{z_a \in Z_a^s(i)} \phi_a^{t_m} + \sum_{z_a \in Z_a^d(i)} \phi_a^{t_m} \right) \right] \\ &= \sum_{z_i \in Z} \sum_{z_p \in Z_p(i)} \left[\log \left(\sum_{z_a \in Z_a^s(i)} (1 + \alpha_{z_a}) \phi_a^{t_l} \right. \right. \\ &\quad \left. \left. + \sum_{z_a \in Z_a^d(i)} (1 - \beta_{z_a}) \phi_a^{t_l} \right) \right], \end{aligned} \quad (14)$$

$$\begin{aligned} V_p^{t_m} &= \sum_{z_i \in Z} \left[\sum_{z_p \in Z_p^s(i)} \log \phi_p^{t_m} + \sum_{z_p \in Z_p^d(i)} \log \phi_p^{t_m} \right] \\ &= \sum_{z_i \in Z} \left[\sum_{z_p \in Z_p^s(i)} \log(1 + \alpha_{z_p}) \phi_p^{t_l} \right. \\ &\quad \left. + \sum_{z_p \in Z_p^d(i)} \log(1 - \beta_{z_p}) \phi_p^{t_l} \right]. \end{aligned} \quad (15)$$

Assumption 3

Let $\overline{\alpha_{z_x}}$ be the mean increasing rate of similarity (*i.e.*, α_{z_x}) over $Z_x^s(i)$ and $\overline{\beta_{z_x}}$ be the mean decreasing rate of similarity (*i.e.*, β_{z_x}) over $Z_x^d(i)$, then $\overline{\alpha_{z_x}} \approx \overline{\beta_{z_x}}$.

Definition 6

In Eq. 11, let the mean $\phi_a^{t_l}$ over $Z_a^s(i)$ be $\overline{\phi_a^{t_l}}^s$ and that over $Z_a^d(i)$ be $\overline{\phi_a^{t_l}}^d$.

Assumption 4

Let the difference between $\overline{\phi_a^{t_l}}^s$ and $\overline{\phi_a^{t_l}}^d$ by sensitive attribute information be $\epsilon \in R^+$. Then, $\overline{\phi_a^{t_l}}^s \approx \overline{\phi_a^{t_l}}^d + \epsilon$, where $\epsilon \ll \overline{\phi_a^{t_l}}^d, \overline{\phi_a^{t_l}}^s$.

Lemma 1

Given \tilde{X}, \tilde{Y} , and \tilde{S} , for all t_l, t_m , $V_a^{t_l} \geq V_a^{t_m}$.

proof.

From Proposition 3,

$$\begin{aligned}
V_a^{t_m} &= \sum_{z_i \in Z} \sum_{z_p \in Z_p(i)} \left[\log \left(\sum_{z_a \in Z_a^s(i)} (1 + \alpha_{z_a}) \phi_a^{t_l} \right. \right. \\
&\quad \left. \left. + \sum_{z_a \in Z_a^d(i)} (1 - \beta_{z_a}) \phi_a^{t_l} \right) \right] \\
&\approx \sum_{z_i \in Z} \sum_{z_p \in Z_p(i)} \left[\log \left(\sum_{z_a \in Z_a^s(i)} (1 + \overline{\alpha_{z_a}}) \phi_a^{t_l} \right. \right. \\
&\quad \left. \left. + \sum_{z_a \in Z_a^d(i)} (1 - \overline{\beta_{z_a}}) \phi_a^{t_l} \right) \right]. \tag{16}
\end{aligned}$$

Note that $\overline{\alpha_{z_a}}$ and $\overline{\beta_{z_a}}$ are defined in Assumption 3. Then we compare $V_a^{t_m}$ and $V_a^{t_l}$ as follows.

$$\begin{aligned}
\Delta V_a &= V_a^{t_m} - V_a^{t_l} \\
&\approx \sum_{z_i \in Z} \sum_{z_p \in Z_p(i)} \left[\log \left(\frac{\sum_{z_a \in Z_a^s(i)} (1 + \overline{\alpha_{z_a}}) \phi_a^{t_l}}{\sum_{z_a \in Z_a(i)} \phi_a^{t_l}} \right. \right. \\
&\quad \left. \left. + \frac{\sum_{z_a \in Z_a^d(i)} (1 - \overline{\beta_{z_a}}) \phi_a^{t_l}}{\sum_{z_a \in Z_a(i)} \phi_a^{t_l}} \right) \right] \\
&= \sum_{z_i \in Z} \sum_{z_p \in Z_P(i)} \left[\log \left(1 + \frac{\sum_{z_a \in Z_a^s(i)} \overline{\alpha_{z_a}} \phi_a^{t_l}}{\sum_{z_a \in Z_a(i)} \phi_a^{t_l}} \right. \right. \\
&\quad \left. \left. - \frac{\sum_{z_a \in Z_a^d(i)} \overline{\beta_{z_a}} \phi_a^{t_l}}{\sum_{z_a \in Z_a(i)} \phi_a^{t_l}} \right) \right]. \tag{17}
\end{aligned}$$

By Definition 6, it is rephrased as follows.

$$\begin{aligned}
\Delta V_a &= \sum_{z_i \in Z} \sum_{z_p \in Z_P(i)} \left[\log \left(1 + \frac{\sum_{z_a \in Z_a^s(i)} \overline{\alpha_{z_a}} \overline{\phi_a^{t_l}}^s}{\sum_{z_a \in Z_a(i)} \phi_a^{t_l}} \right. \right. \\
&\quad \left. \left. - \frac{\sum_{z_a \in Z_a^d(i)} \overline{\beta_{z_a}} \overline{\phi_a^{t_l}}^d}{\sum_{z_a \in Z_a(i)} \phi_a^{t_l}} \right) \right]. \tag{18}
\end{aligned}$$

From Assumption 4, $\overline{\phi_a^{t_x}}^s \approx \overline{\phi_a^{t_x}}^d + \epsilon$. Based on this, ΔV_a is approximated as follows.

$$\begin{aligned}
\Delta V_a &\approx \sum_{z_i \in Z} \sum_{z_p \in Z_P(i)} \left[\log \left(1 + \frac{\left(\sum_{z_a \in Z_a^s(i)} \overline{\alpha_{z_a}} \right. \right. \right. \\
&\quad \left. \left. - \frac{\sum_{z_a \in Z_a^d(i)} \overline{\beta_{z_a}} \overline{\phi_a^{t_l}}^d}{\sum_{z_a \in Z_a(i)} \phi_a^{t_l}} \right) \right), \tag{19}
\end{aligned}$$

where we omit ϵ for readability since $\epsilon \ll \overline{\phi_a^{t_x}}^d, \overline{\phi_a^{t_x}}^s$. In the ideally biased dataset, regardless of z_i and z_p , $|Z_a^s(i)| = rC + (m-1)C$ and $|Z_a^d(i)| = (m-1)rC + (m-1)^2C$. Thus, we can reformulate Eq. 19 as follows.

$$\begin{aligned}
\Delta V_a &= \sum_{z_i \in Z} \sum_{z_p \in Z_P(i)} \left[\log \left(1 + \frac{\left(rC + (m-1)C \right) \overline{\alpha_{z_a}}}{\sum_{z_a \in Z_a(i)} \phi_a^{t_l}} \right. \right. \\
&\quad \left. \left. - \frac{\left((m-1)rC + (m-1)^2C \right) \overline{\beta_{z_a}} \overline{\phi_a^{t_l}}^d}{\sum_{z_a \in Z_a(i)} \phi_a^{t_l}} \right) \right] \\
&= \sum_{z_i \in Z} \sum_{z_p \in Z_P(i)} \left[\log \left(1 + \frac{(m+r-1)C \left(\overline{\alpha_{z_a}} \right. \right. \right. \\
&\quad \left. \left. - \frac{(m-1)\overline{\beta_{z_a}} \overline{\phi_a^{t_l}}^d}{\sum_{z_a \in Z_a(i)} \phi_a^{t_l}} \right) \right). \tag{20}
\end{aligned}$$

By Assumption 3, it is approximated as follows.

$$\begin{aligned}
\Delta V_a &\approx \sum_{z_i \in Z} \sum_{z_p \in Z_P(i)} \left[\log \left(1 + \frac{(m+r-1)C}{\sum_{z_a \in Z_a(i)} \phi_a^{t_l}} \right. \right. \\
&\quad \left. \left. \times \left(\frac{\overline{\alpha_{z_a}} \overline{\phi_a^{t_l}}^d}{\sum_{z_a \in Z_a(i)} \phi_a^{t_l}} \right) \right) \right] \leq 0. \tag{21}
\end{aligned}$$

Here, $m \geq 2$, $r > m^2$, $C > 0$ (\because Assumption 1), $\overline{\alpha_{z_a}} > 0$ (\because Proposition 3), and $\phi_a^{t_l} > 0$ (\because Definition 4). Therefore, $\Delta V_a \leq 0$.

Lemma 2

Given \tilde{X} , \tilde{Y} , and \tilde{S} , for all t_l, t_m , $V_p^{t_l} < V_p^{t_m}$.

proof.

By Proposition 3,

$$\begin{aligned}
V_p^{t_m} &= \sum_{z_i \in Z} \left[\sum_{z_p \in Z_P^s(i)} \log(1 + \alpha_{z_p}) \phi_p^{t_l} \right. \\
&\quad \left. + \sum_{z_p \in Z_P^d(i)} \log(1 - \beta_{z_p}) \phi_p^{t_l} \right] \\
&\approx \sum_{z_i \in Z} \left[\sum_{z_p \in Z_P^s(i)} \log(1 + \overline{\alpha_{z_p}}) \phi_p^{t_l} \right. \\
&\quad \left. + \sum_{z_p \in Z_P^d(i)} \log(1 - \overline{\beta_{z_p}}) \phi_p^{t_l} \right]. \tag{22}
\end{aligned}$$

Similar to Eq. 17, we compare $V_p^{t_m}$ and $V_p^{t_l}$ as follows.

$$\Delta V_p = V_p^{t_m} - V_p^{t_l} = \sum_{z_i \in Z} \left[\sum_{z_p \in Z_p^s(i)} \log \frac{(1 + \overline{\alpha_{z_p}}) \phi_p^{t_l}}{\phi_p^{t_l}} + \sum_{z_p \in Z_p^d(i)} \log \frac{(1 - \overline{\beta_{z_p}}) \phi_p^{t_l}}{\phi_p^{t_l}} \right]. \quad (23)$$

Here, $\log(1 - \overline{\alpha_{z_p}}) \approx \log(1 - \overline{\beta_{z_p}})$ (\because Assumption 3), and $\log(1 - \overline{\alpha_{z_p}}) \approx -\log(1 + \overline{\alpha_{z_p}})$ since $\log(1) = 0$ and $\frac{d \log(1)}{dx} = 1$. Therefore, $\log(1 + \overline{\alpha_{z_p}}) \approx -\log(1 - \overline{\beta_{z_p}})$. Based on this, we can approximate ΔV_p as follows.

$$\Delta V_p \approx \log(1 + \overline{\alpha_{z_p}}) \left(\sum_{z_i \in Z} \sum_{z_p \in Z_p^s(i)} \mathbb{1} - \sum_{z_i \in Z} \sum_{z_p \in Z_p^d(i)} \mathbb{1} \right), \quad (24)$$

where $\mathbb{1}$ is an indicator function. In the ideally biased dataset, $\sum_{z_i \in Z} \sum_{z_p \in Z_p^s(i)} \mathbb{1} = (rC)^2 + (m-1)C^2$ and $\sum_{z_i \in Z} \sum_{z_p \in Z_p^d(i)} \mathbb{1} = 2(m-1)rC^2 + (m-1)(m-2)C^2$. Therefore, we rephrase it as follows.

$$\begin{aligned} \Delta V_p &= \left((rC)^2 + (m-1)C^2 \right) \\ &\quad - \left(2(m-1)rC^2 + (m-1)(m-2)C^2 \right) \log(1 + \overline{\alpha_{z_p}}) \\ &= C^2 \left(r^2 + (-2m+1)r - m^2 + 4m - 3 \right) \log(1 + \overline{\alpha_{z_p}}) \\ &= C^2 \left((r + \lambda m)(r - ((2 + \lambda)m - 1)) + (4 - \lambda)m - 3 \right) \\ &\quad \times \log(1 + \overline{\alpha_{z_p}}) > 0 \quad \text{s.t.} \quad r > (2 + \lambda)m - 1 \end{aligned} \quad (25)$$

where $\lambda = -1 + \sqrt{2}$. Finally, $\Delta V_p > 0$ since $m > 2$, $r > m^2$, and $C > 0$ (\because Assumption 1).

Theorem 1

Given \tilde{X} , \tilde{Y} , and \tilde{S} , for all t_l, t_m , $V^{t_l} > V^{t_m}$.

proof.

From Lemma 1 and 2, $V_a^{t_l} \geq V_a^{t_m}$ and $V_p^{t_l} < V_p^{t_m}$ for all t_l, t_m . Since $V^{t_k} = \hat{C}(-V_p^{t_k} + V_a^{t_k})$ by Proposition 1, $V^{t_l} > V^{t_m}$ for all t_l, t_m .

Corollary 1

Learning sensitive attribute information decreases L^{Sup} , given \tilde{X} , \tilde{Y} , and \tilde{S} .

proof.

From Definition 1, learning of sensitive attribute information equals to the increase of $I(Z; \tilde{S})$. In addition, the increase of $I(Z; \tilde{S})$ corresponds to a transition from t_l to

t_m since $I(Z; \tilde{S})$ is always higher at t_m than at t_l (\because Assumption 2). Finally, V^{t_m} is always smaller than V^{t_l} (\because Theorem 1), therefore, learning sensitive attribute information decreases L^{Sup} .

Method	Adversarial Training	EO (\downarrow)	Acc. (\uparrow)
<i>SupCon</i>	\times	30.5 \pm 1.3	80.5 \pm 0.7
	\checkmark	20.0 \pm 0.3	77.2 \pm 0.1
<i>FSCL+</i>	\times	6.5\pm0.4	79.1 \pm 0.1
	\checkmark	20.5 \pm 0.4	77.8 \pm 0.2

Table 6. Effectiveness of adversarial training in classifier training stage on CelebA. We set *attractiveness* and *male* to the target class and sensitive attribute, respectively.

C. Fairness Strategy in Classifier Training Stage

In Table 6, we explore the effectiveness of applying *GRL* [38] in the classifier training stage, after finishing the representation learning with *SupCon* and *FSCL+*. To this end, we deploy an additional classifier for the sensitive attribute and do not freeze the encoder and projection networks in the second stage. As might be expected, *GRL* improves the fairness of *SupCon* by sacrificing the classification accuracy. Meanwhile, it degrades EO as well as the classification accuracy in ours. We speculate that it is because the fair representation learned by *FSCL+* becomes biased by re-training the encoding networks with the cross entropy loss and *GRL*. The similar results of EO and top-1 accuracy between *SupCon* with *GRL* and *FSCL+* with *GRL* support that the learned representation is almost renewed in the classifier training stage. In conclusion, the results show that applying the additional strategy for fairness in the classifier training stage is not effective to our method.

D. Modification for Incomplete Supervised Setting

To apply our method to the environment where target class labels are not provided, we introduce *FSCL[†]*, which is a modified version of *FSCL*. We set a positive sample to another patch from the same image with an anchor and negative samples to Z_{ig} and Z_{tg} . It is formulated as follows.

$$FSCL^\dagger = - \sum_{z_i \in Z} \log \frac{\exp(z_i \cdot z_p / \tau)}{\sum_{z_f^* \in Z_f^*(i)} \exp(z_i \cdot z_f^* / \tau)}, \quad (26)$$

where $Z_f^*(i) = \{z_f^* \in Z | \hat{s}_f = \hat{s}_i\}$. Except for the loss function, the overall structure is the same as the original.

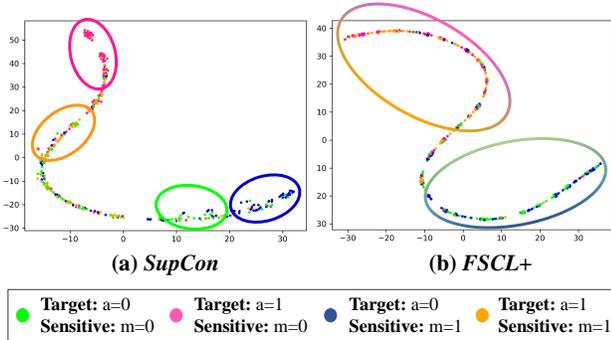


Figure 7. t-SNE visualizations with random initialization.

E. Details of t-SNE Visualization

For the t-SNE [50] visualization, we exploit the models pre-trained on CelebA dataset [31] for 100 epochs. Then we obtain 50 random samples (*i.e.*, representation) per data group with the pre-trained models. Before applying the t-SNE algorithm, we reduce the dimensionality of the samples using PCA reduction. We tune the hyperparameters in the scikit-learn implementation as follows.

- Perplexity: from 10 to 40 by 1
- Learning rate: 10 or 100
- Iteration= 100, 1000, or 10000

We set the perplexity, learning rate, and iterations 10, 10, and 10000 respectively, but in all the cases, we note that representation learned by *FSCL+* is more agnostic to the sensitive attribute than that learned by *SupCon*. Furthermore, we provide t-SNE plots without PCA reduction in Figure 7 since it considerably affects the structure of representations.

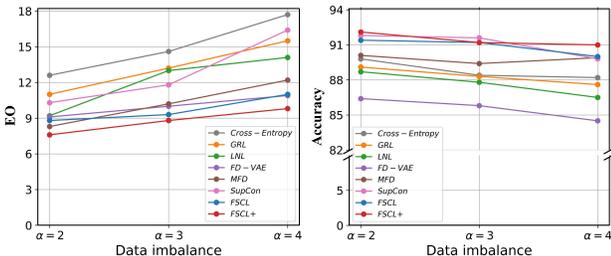


Figure 8. Classification results on UTK Face dataset. We set *gender* and *age* to the target class and sensitive attribute, respectively. It shows trends of classification accuracy and equalized odds (EO) at different α .

F. Further Experiments on UTK Face

In Figure 8, we provide experimental results on UTK Face with the other sensitive attribute, *age*. It shows that *FSCL+* maintain the fairest EO and the best top-1 accuracy at all α .

Although FD-VAE [37] achieves similar EO with *FSCL*, its accuracy is significantly inferior to ours. It indicates that ours highly outperform it in terms of the trade-off performance between fairness and accuracy.

G. Additional Experimental Results on CelebA

To clearly show the trade-off performances between classification accuracy and fairness, we plot the experimental results on CelebA in Figure 9. *FSCL+* achieves the best trade-off performances in all the results. Furthermore, we supplement the experimental results by reporting standard deviation in Table 7.

H. Dataset Composition

H.1. CelebA and UTK Face

In CelebA [31], we conduct experiments in terms of a variety of target and sensitive attribute pairs. Table 8 shows the specific composition of the training set in all the settings. In UTK Face [58], we involve 10,000, 2,400, and 2,400 data in the training, validation, and test sets, respectively. We provide the various compositions of the training set according to α in Table 9.

H.2. Dogs and Cats

Similar to UTK Face, we leverage 3,425 black cat and white dog images, and 685 white cat and black dog images for training. The test set includes 2,400 images which are completely balanced. We note that it is different from the original setting in [26]. In the study, the target attribute and bias are completely correlated in the training set. For instance, cats are always black and dogs are always white. Although they solved the task by utilizing the pixel-level of bias labels (*i.e.*, RGB values of each pixel), it is an almost unsolvable problem with only the image-level of labels since the target attribute and bias labels are always the same at the image-level. Therefore, we designed the task more reasonable to validate fairness methods which mostly exploit the image-level of labels.

H.3. Discussion on License and Data Collection

Both CelebA [31] and UTK Face [58] have a non-standard license (*i.e.*, Custom (non-commercial)), but the creators clarify the datasets are available for non-commercial research purposes only.

CelebA consists of the images collected from Celeb-Faces dataset [47] and attribute labels. According to [47], the images are collected by searching names of celebrities on the web. Also in UTK Face, the creators combine the images from CACD [4] and Morph [20] datasets with the images crawled in Bing and Google search engines. In both CACD and Morph, the images are gathered by searching on the web.

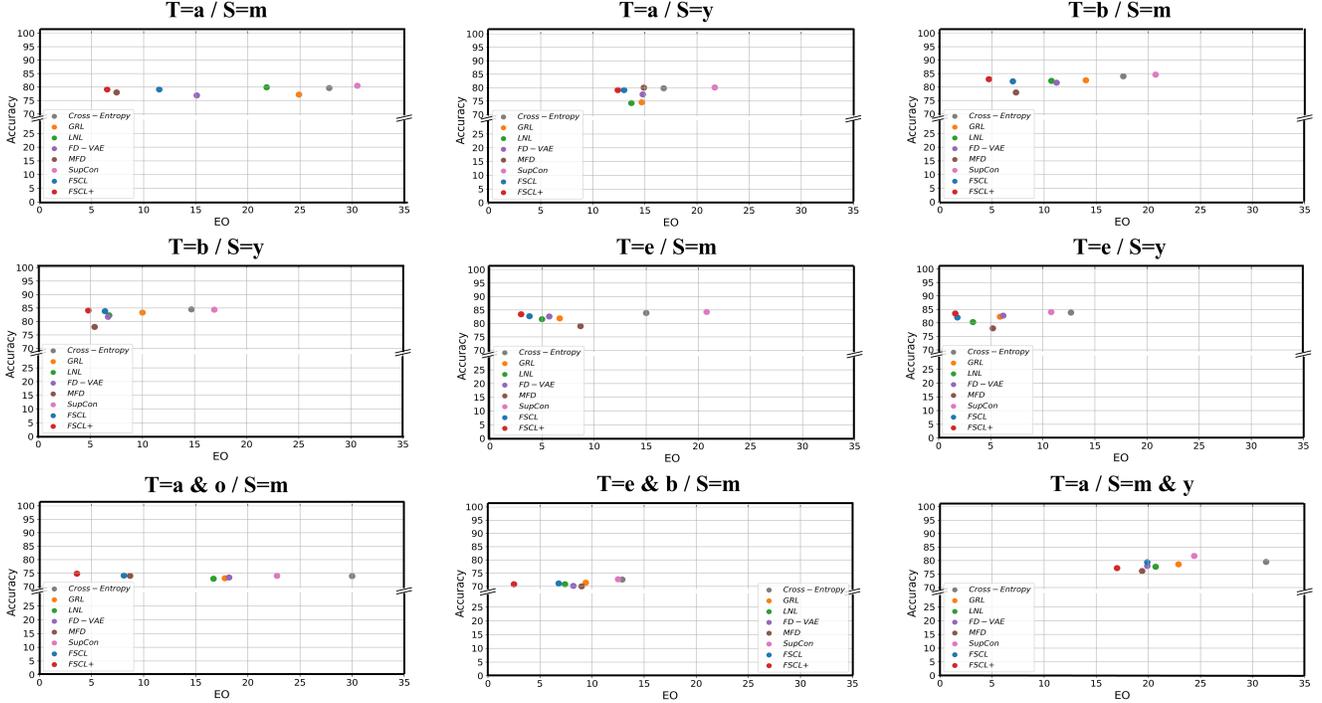


Figure 9. Experimental results in figure form on CelebA dataset. It shows the trade-off performances between ACC. and EO more clearly. The upper left corner of the plots corresponds to the optimal trade-off performance.

Attributes	CE [15]		GRL [38]		LNL [26]		FD-VAE [37]		MFD [22]		SupCon [25]		FSCL		FSCL+	
	EO	Acc.	EO	Acc.	EO	Acc.	EO	Acc.	EO	Acc.	EO	Acc.	EO	Acc.	EO	Acc.
T=a / S=m	27.8±0.2	79.6±0.5	24.9±0.3	77.2±0.5	21.8±0.4	79.9±0.5	15.1±0.1	76.9±0.0	7.4±0.3	78.0±0.3	30.5±1.3	80.5±0.7	11.5±0.3	79.1±0.4	6.5±0.4	79.1±0.4
T=a / S=y	16.8±0.3	79.8±0.4	14.7±0.4	74.6±0.4	13.7±0.3	74.3±0.4	14.8±0.2	77.5±0.1	14.9±0.4	80.0±0.3	21.7±1.0	80.1±0.8	13.0±0.6	79.1±0.5	12.4±0.5	79.1±0.5
T=b / S=m	17.6±0.3	84.0±0.3	14.0±0.3	82.5±0.5	10.7±0.2	82.3±0.4	11.2±0.1	81.6±0.3	7.3±0.2	78.0±0.3	20.7±0.5	84.6±0.6	7.0±0.4	82.1±0.3	4.7±0.5	82.9±0.4
T=b / S=y	14.7±0.1	84.5±0.3	10.0±0.2	83.3±0.5	6.8±0.3	82.3±0.5	6.7±0.2	81.7±0.0	5.4±0.1	78.0±0.2	16.9±0.9	84.4±0.8	6.4±0.4	83.8±0.4	4.8±0.3	84.1±0.5
T=e / S=m	15.0±0.3	83.9±0.2	6.7±0.4	81.9±0.6	5.0±0.3	81.6±0.3	5.7±0.0	82.6±0.1	8.7±0.3	79.0±0.4	20.8±1.1	84.3±0.5	3.8±0.3	82.7±0.3	3.0±0.4	83.4±0.6
T=e / S=y	12.7±0.2	83.8±0.3	5.9±0.4	82.3±0.4	3.3±0.4	80.3±0.6	6.2±0.1	84.0±0.2	5.2±0.2	78.0±0.2	10.8±1.0	84.0±0.7	1.8±0.3	82.0±0.4	1.6±0.3	83.5±0.3
T=a & o / S=m	30.0±0.2	73.9±0.5	17.8±0.2	73.1±0.5	16.7±0.4	72.9±0.5	18.2±0.1	73.4±0.1	8.7±0.4	74.0±0.3	22.8±0.7	74.0±0.5	8.1±0.3	74.1±0.3	3.6±0.3	74.8±0.4
T=b & e / S=m	12.9±0.2	72.6±0.4	9.4±0.3	71.4±0.4	7.4±0.2	70.8±0.5	8.2±0.1	70.2±0.2	9.0±0.1	70.0±0.1	12.5±0.8	72.7±0.9	6.8±0.4	71.1±0.2	2.5±0.6	70.8±0.5
T=a / S=m & y	31.3±0.3	79.5±0.4	22.9±0.4	78.6±0.5	20.7±0.3	77.7±0.5	19.9±0.0	78.0±0.1	19.4±0.2	76.1±0.3	24.4±1.3	81.7±0.7	19.9±0.5	79.4±0.3	17.0±0.5	77.2±0.5

Table 7. Classification results on CelebA. We further specify the standard deviation in this table.

I. Implementation Details

I.1. Structure of Comparable Models

Cross-Entropy [15], **GRL** [38], **LNL** [26]: The models utilize ResNet-18 [15] for backbone networks and a MLP with one hidden layer for classifiers. The dimensions of representation are the same as ours. **GRL** and **LNL** are reproduced based on [26, 38], and the hyperparameter to determine a weight for the reversed gradient is searched in the range from 0.01 to 0.1 in each experiment. For **LNL**, hyperparameter λ for regularization loss is searched in the range from 0.01 to 0.1 in each experiment. For all the models, we train them in an end-to-end manner for 100 epochs.

FD-VAE [37]: We build the model with the same structure as the original paper [37] without the encoder network.

For a fair comparison, we substitute the encoder network to ResNet-18 and obtain better reproduction performances. Following the paper, we separate each latent space to have the same dimensions to each other and set hyperparameter β to 1. The other hyperparameters are found by grid searching and set to $\alpha = 1$, $\gamma = 5$, and $\lambda = 1$ for all the experiments. For representation learning, we train the encoder networks for 100 epochs. After that, we train the classifiers for downstream tasks for 10 epochs.

MFD [22]: We implement the model with source code released by the authors. The teacher and student models both leverage ResNet-18 for backbone networks and a MLP with one hidden layer for a classifier. Following the original paper, we train the models for 50 epochs and set hyperparameter λ to 7 and 5 for CelebA and UTK Face, respectively. For Dogs

CelebA								
	$a=0$	$a=1$		$b=0$	$b=1$		$e=0$	$e=1$
$m=0$	29,920	64,589	$m=0$	84,954	9,555	$m=0$	84,963	9,546
$m=1$	49,247	19,014	$m=1$	39,475	28,786	$m=1$	44,527	23,734
	$a=0$	$a=1$		$b=0$	$b=1$		$e=0$	$e=1$
$y=0$	30,618	5,364	$m=0$	19,164	16,818	$m=0$	22,146	13,836
$y=1$	48,549	78,239	$m=1$	105,265	21,523	$m=1$	107,344	19,444
	$a=0$	$a=1$	$m=0$	$m=1$		$m=0$	$m=1$	
$m=0, y=0$	7,522	3,645	$a=0, o=0$	13,995	27,966	$b=0, e=0$	78,613	30,481
$m=1, y=0$	23,096	1,719	$a=1, o=0$	30,943	11,380	$b=1, e=0$	6,350	14,046
$m=0, y=1$	22,398	60,944	$a=0, o=1$	15,925	21,281	$b=0, e=1$	6,341	8,994
$m=1, y=1$	26,151	17,295	$a=1, o=1$	33,646	7,634	$b=1, e=1$	3,205	14,740

Table 8. **Composition of the training set of CelebA.** $a, b, e, o, m,$ and y denote *attractiveness, bignose, bags-under-eyes, mouth-slightly-open, male,* and *young,* respectively.

UTK Face					
$\alpha = 2 / \alpha = 3 / \alpha = 4$					
	Ethnicity			Age	
	Caucasian		Others	More than 35	Others
Female	1,666 / 1,250 / 1,000		3,334 / 3,750 / 4,000	1,666 / 1,250 / 1,000	3,334 / 3,750 / 4,000
Male	3,334 / 3,750 / 4,000		1,666 / 1,250 / 1,000	3,334 / 3,750 / 4,000	1,666 / 1,250 / 1,000

Table 9. **Composition of the training set of UTK Face.** α denotes the intensities of data imbalance.

and Cats, λ is determined as 7 through grid searching.

SupCon [25], **SimCLR** [5], **FSCL (ours)**: We implement *SupCon* and *SimCLR* with source code released by the authors of [25], and *FSCL* is also based on the code (which is licensed under the terms of the MIT license). The models use ResNet-18 [15] for the encoder network and a MLP with two hidden layers for the projection network, which have 256 hidden nodes.

I.2. Augmentation Strategy and Experimental Setup

For the models based on contrastive loss, we augment two patches per image. Except for this, we use the same augmentation strategy [5] for all the models. Specifically, we sequentially and randomly apply cropping and resizing, horizontal flipping, color jittering, and gray scaling.

For all the models, we set the identical environments of SGD optimizer with momentum [40], batch sizes of 128, and learning rate of 0.1. All the experiments are based on the PyTorch library and are conducted in a Linux environment with 4 NVIDIA Titan Xp GPUs with 12GB of memory.

J. Comparison with GDRO

GDRO [42] is one of the state-of-the-art methods to minimize the performance gaps between data groups and has a

Method	Regularization	EO (\downarrow)	Acc. (\uparrow)
GDRO	Standard	21.3 \pm 1.0	76.3 \pm 0.2
	Early Stopping	4.0 \pm 0.1	74.7 \pm 0.1
	Strong L_2 (lr=0.1)	8.7 \pm 2.6	76.3 \pm 0.1
	Strong L_2 & Group adjustments (C=5)	8.0 \pm 2.0	77.1 \pm 0.2
FSCL+	Standard	6.5 \pm 0.4	79.1 \pm 0.1

Table 10. **Comparison with GDRO on CelebA.** We set *attractiveness* and *male* to the target class and sensitive attribute, respectively.

goal similar to our group-wise normalization. Thus, we report comparison results with GDRO in Table 10. Following the original paper, we search for the best C in the range of [0, 5]. The results show that ours achieves a better trade-off performance than GDRO.

K. Two kinds of Supervised Contrastive Losses

In this section, we summarize two kinds of supervised contrastive losses (*i.e.*, L_{out}^{sup} and L_{in}^{sup}) proposed in [25] and why we leverage L_{out}^{sup} as our baseline. Unlike L_{out}^{sup} (*i.e.*, L^{Sup} in the main paper), L_{in}^{sup} places the summation over positive samples and the normalization factor inside the log

as follows.

$$L_{in}^{Sup} = - \sum_{z_i \in Z} \log \left(\frac{1}{|Z_p(i)|} \sum_{z_p \in Z_p(i)} \frac{\phi_p}{\sum_{z_a \in Z_a(i)} \phi_a} \right). \quad (27)$$

In the loss, the normalization factor works as a constant (*i.e.*, $-\sum_{z_i \in Z} \log \frac{1}{|Z_p(i)|}$), so it cannot normalize the imbalance in the positive samples. As the result, L_{in}^{sup} is more vulnerable to the data bias and shows inferior classification performances to L_{out}^{sup} . For these reasons, we utilize the latter as our baseline.

L. Discussion on Limitations

In this section, we discuss two limitations of our study. The first one is that our work is confined to the image classification task. We discuss it by explaining why we cover the task in this paper. One reason is that the superior performance of our baselines (*i.e.*, SupCon and SimCLR) has been experimentally validated in the image classification task [5, 25]. Therefore, through the task, we can make a fair comparison with the models and convincingly demonstrate our improvement over them. The other reason is that image classification is a fundamental and common task not only in contrastive representation learning [5, 25, 48, 54] but in fairness studies in the field of computer vision [6, 37, 44, 53]. Although fair visual representation can be exploited in other tasks, such as object recognition [52], image-to-image translation [17, 18], face recognition [2, 10], and object detection [1], each of them requires a suitable notion of fairness [1, 52] and specialized architectures [10, 17, 18]. Therefore, to achieve the best performance on the tasks, we also need to modify the proposed loss more appropriately for them. We leave the extension of *FSCl* to broader tasks for future work.

Second, our method essentially requires sensitive attribute labels to improve fairness. Even though supervision of the sensitive attribute labels is common in the literature on fair classification [6, 34, 37, 44], sometimes we cannot access the labels and it is laborious and expensive to annotate them. Although we show that our method can reduce such costs by effectively improving fairness using only a few labels, it cannot be utilized in the complete absence of the labels. Therefore, future works that develop a fair contrastive loss free of the sensitive attribute labels would make a significant contribution to the research community. We expect our study to be a bridgehead for them.