

# Fair Feature Distillation for Visual Recognition

Sangwon Jung<sup>1\*</sup>, Donggyu Lee<sup>1\*</sup>, Taeon Park<sup>1\*</sup> and Taesup Moon<sup>2†</sup>

<sup>1</sup>Department of Electrical and Computer Engineering, Sungkyunkwan University, Suwon, Korea

<sup>2</sup>Department of Electrical and Computer Engineering, Seoul National University, Seoul, Korea

{s.jung, ldk308, pte1236}@skku.edu, tsmoon@snu.ac.kr

## Abstract

Fairness is becoming an increasingly crucial issue for computer vision, especially in the human-related decision systems. However, achieving algorithmic fairness, which makes a model produce indiscriminate outcomes against protected groups, is still an unresolved problem. In this paper, we devise a systematic approach which reduces algorithmic biases via feature distillation for visual recognition tasks, dubbed as MMD-based Fair Distillation (MFD). While the distillation technique has been widely used in general to improve the prediction accuracy, to the best of our knowledge, there has been no explicit work that also tries to improve fairness via distillation. Furthermore, We give a theoretical justification of our MFD on the effect of knowledge distillation and fairness. Throughout the extensive experiments, we show our MFD significantly mitigates the bias against specific minorities without any loss of the accuracy on both synthetic and real-world face datasets.

## 1. Introduction

Based on the remarkable performance of deep neural networks, computer vision has become one of the core technologies in many applications that affect various aspects of society; e.g., facial recognition [24], AI-assisted hiring [25], healthcare diagnostics [13], and law enforcement [11]. Due to these social applications of computer vision algorithms, it is becoming increasingly essential for them to be *fair*; namely, the outcomes of systems should not be discriminative against any certain groups on the basis of sensitive attributes. For example, any automated system that incorporates photographs into a decision process (e.g., job interview) should not rely on certain sensitive attributes, such as race or gender [29]. However, recent studies demonstrate that commercial API systems for facial analysis expose the gender/race bias in widely used face datasets [6, 34].

In this work, we are interested in the setting in which an

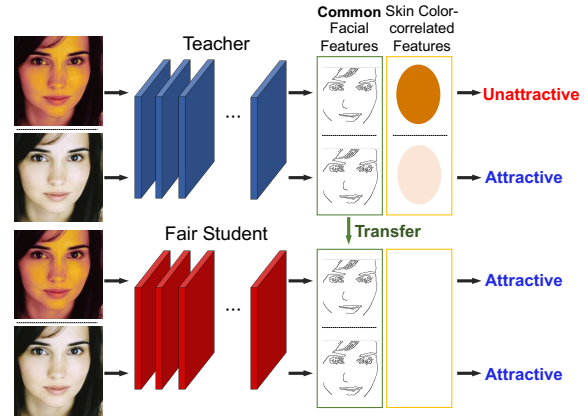


Figure 1. An illustrative example of motivation to our work. The “teacher” model may depend heavily on the skin color when deciding whether the face is attractive, while it may also have learned useful common (unbiased) facial features. To train a fair “student” model via feature distillation, only the unbiased common features from the teacher should be transferred to the student so that both high accuracy and fairness can be achieved.

already deployed model has been identified as unfair. The usual approach of the so-called *in-processing* methods to mitigate the unfair bias is to re-train the model from scratch with an additional fairness constraint [1, 18, 37]. However, such approaches typically do not utilize any predictive information already learned out by the deployed model, and hence, would lead to sacrificing the accuracy for the improved fairness. To address above limitation, the *knowledge distillation (KD)* [16] technique can be considered as a potential tool for leveraging the deployed model’s predictive power while re-training with fairness constraints. Nonetheless, the typical existing KD methods [16, 30, 38, 31, 17] focused only on improving the accuracy, and considering *both* the accuracy and fairness during the process of KD is not straightforward. We aim to resolve this challenge by proposing a new fairness-aware feature distillation scheme.

Figure 1 illustrates the key idea of our work. We assume that even when the original deployed model, the “teacher” model, may be heavily biased (i.e., heavily use the sensi-

\*The first three authors have contributed equally.

†Corresponding author (E-mail: tsmoon@snu.ac.kr)

tive “skin color” attribute), it could also have learned useful group-indistinguishable common (unbiased) features that are effective for achieving high prediction accuracy (*e.g.*, “face shape”, etc.). Our intuition is that when training a “student” model, if only those common unbiased features can be transferred from the teacher, the student should be able to achieve higher accuracy, compared to the ordinary *in-processing* methods that re-train from scratch, as well as better fairness, compared to the original teacher.

In order to realize above intuition, we propose a fair feature distillation technique by utilizing the *maximum mean discrepancy* (MMD), dubbed as MMD-based Fair Distillation (MFD); this is, to the best of our knowledge, the first approach to improve both accuracy and fairness via distillation. More concretely, we devise a regularization term for training a student that enforces the distribution of the group-conditioned features of the student to get closer to the distribution of the group-averaged features of the teacher in the MMD sense. We further provide a theoretical understanding that our MFD regularization can indeed lead to improving both the accuracy and fairness of the student in a principled way. Namely, we show our regularization term induces the distributions of the group-conditioned features of the student to get close to each other across all the sensitive groups (*i.e.*, promotes fairness), while making all those distributions also get close to the distribution of the group-averaged features of the highly accurate teacher (*i.e.*, improves accuracy via the distillation effect).

As a result, we convincingly show through extensive experiments that our MFD can simultaneously improve the accuracy as well as considerably mitigate the unfair bias of a model. Firstly, we construct a synthetic dataset, CIFAR-10S [35], and systematically validate our motivation illustrated in Figure 1. Then, with additional experiments on two real-world datasets, UTKFace [42] and CelebA [22], we identify that our MFD is the only method that can *consistently* improve both accuracy and fairness of the original unfair teacher on all three datasets, compared to the three types of baselines: ordinary KD methods, representative *in-processing* methods that re-train from scratch, and methods that naively combine the *in-processing* methods with KD methods. Finally, we demonstrate the validity of our theoretical bound via systematic ablation studies.

## 2. Related Works

**Algorithmic fairness** Recently, a number of studies have focused on mitigating unfairness, as exhaustively surveyed in [4]. Fairness algorithms are mainly divided into three categories depending on the training pipelines they apply: *pre-processing* methods [7, 23, 29, 39] that refine a dataset to remove the source of unfairness before training a model, *in-processing* methods [1, 9, 18, 19, 37, 40] that take the fairness constraints into account when training a model, and

*post-processing* methods [3, 14] that modify the predicted labels after training.

In this paper, we focus on the *in-processing* methods, since they can be particularly useful for the circumstances in which controlling the model itself is possible. Among them, some researches formulate optimization problems with a fairness constraint indicating statistical independence between the model’s outputs and groups [19, 37]. On the other hand, Zhang *et al.* [40] adopted a simple adversarial debiasing (AD) technique for a model to give outputs from which the sensitive attribute is not predictable by an adversary. Moreover, controlling the contribution of data points to a loss function during training can also help obtain an unbiased machine. It can be done by strategically sampling (SS) the data [9], *e.g.*, oversampling, or assigning unequal weights to the training samples while performing a sequence of classification [18]. In the computer vision domain, the discrimination problem has usually been tackled in facial analysis, such as face recognition [33, 34]. Wang *et al.* [34] mitigated racial bias using the domain adaptation technique. Wang and Deng [33] utilized reinforcement learning. Their algorithms, however, have been specific only to the face recognition tasks.

**Knowledge distillation** For the purpose of knowledge transfer and model compression, diverse approaches to distill helpful information from a learned model have been proposed for deep neural networks. After the original work by Hinton *et al.* [16] (HKD), which matches the softmax output distribution of the teacher to that of the student, various extensions have focused on how to exploit the learned features. The work of Romero *et al.* [30] (FitNet) made the student mimic the features of the teacher through linear regression. Zagoruyko *et al.* [38] proposed attention transfer (AT) which transfers the knowledge using the attention map. Further, Yim *et al.* [36] and Park *et al.* [26] studied approaches using gram matrix and relation map respectively. Unlike the previous methods, several approaches proposed feature distillation algorithms devised from the statistical point of views [2, 17, 27, 31]. In particular, Passalis *et al.* [27] suggested methods to reduce the distance between the teacher and the student feature distributions measured via *Kullback-Liebler* divergence, and Huang *et al.* [17] invented neuron selectivity transfer (NST) utilizing MMD. Although numerous distillation methods have been proposed, none of them explicitly considered the fairness issue during distillation.

## 3. Fairness Criterion

A lot of fairness criteria have been introduced including statistical parity [8], equalized odds [14], overall accuracy equality [5], fairness through awareness [8] and counterfactual fairness [21]. Although each of them tries to tackle the fairness problem from various aspects, choosing the most proper one is still an open question since the notion of fair-

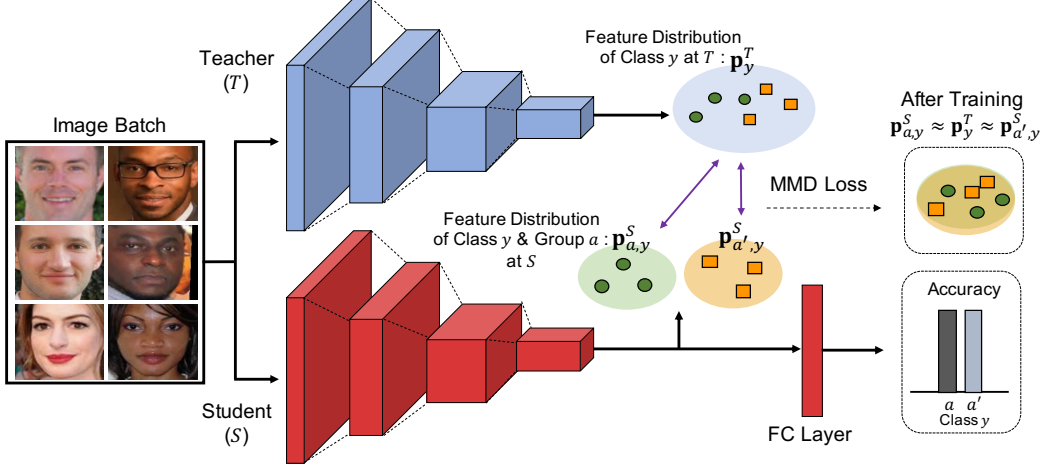


Figure 2. The illustrative concept of MFD. The student treats all groups fairly while learning the teacher’s knowledge by minimizing our MMD-based loss between  $\mathbf{p}_y^T$  and  $\mathbf{p}_{a,y}^S$  for all  $a \in \mathcal{A}$ . Sample images are drawn from UTKFace dataset [42].

ness can differ according to the social, cultural background and the application scenario.

In this work, we consider the *equalized odds* [14], which can be naturally adapted to  $M$ -ary classification and measure per-class accuracy discrepancies between groups. Equalized odds was originally defined in the binary case; given the target variable  $Y = y \in \{-1, 1\}$ , the equalized odds requires the predictor  $\tilde{Y}$  (e.g., the decision of a neural network) and the sensitive attribute  $A \in \mathcal{A}$  to be conditionally independent given  $y$ , i.e.,  $\tilde{Y} \perp A | Y = y$ . For non-binary  $Y$ , equalized odds can also be used to measure the fairness of a model by requiring that  $\forall a, a' \in \mathcal{A}$ ,  $y \in \mathcal{Y} = \{1, \dots, M\}$ ,  $\Pr(\tilde{Y} = y | A = a, Y = y) = \Pr(\tilde{Y} = y | A = a', Y = y)$ . Then, as the equalized odds-based metrics, two types of *difference of equalized odds* (DEO) are defined upon taking the maximum or the average over  $y$  as follows, respectively:

$$\text{DEO}_M \triangleq \max_y \left( \max_{a, a'} \left( \left| \Pr(\tilde{Y} = y | A = a, Y = y) - \Pr(\tilde{Y} = y | A = a', Y = y) \right| \right) \right), \quad (1)$$

$$\text{DEO}_A \triangleq \frac{1}{|\mathcal{Y}|} \sum_y \left( \max_{a, a'} \left( \left| \Pr(\tilde{Y} = y | A = a, Y = y) - \Pr(\tilde{Y} = y | A = a', Y = y) \right| \right) \right). \quad (2)$$

We note that DEO is equivalent to the class-wise accuracy difference between groups over all classes. When there is a considerable discrepancy for a specific class,  $\text{DEO}_M$  is more useful than simply measuring the group accuracy difference. On the contrary, since  $\text{DEO}_M$  only focuses on the worst unfairness,  $\text{DEO}_A$  is also a crucial measure to check the overall fairness across all classes.

## 4. Main Method

In this section, we describe our MFD in details. Our aim is to train a fair *student* model  $S$ , given a *teacher* model  $T$ , which is trained merely considering the accuracy of the given task and could be unfairly biased. Moreover, similarly as in [10], we only consider the case that the network structures of  $S$  and  $T$  are the same. As mentioned in the Introduction, our underlying assumption is that despite being biased,  $T$  could have also learned group-indistinguishable predictive features, hence, distilling those features to  $S$  could achieve higher accuracy than the model re-trained from scratch with fairness constraints, while also improving the fairness over  $T$ .

One straightforward method to achieve our goal is to simply introduce two regularization terms associated with KD and fairness, respectively. The typical regularization for KD employs the difference between the softmax outputs or features of  $T$  and  $S$ , e.g., *Kullback-Leibler divergence* [16] or point-wise  $L_2$  distance [30]. The regularization for fairness can be often specified by the correlation [37] or mutual information [19] of a model’s output and the sensitive attribute, which targets the statistical independence. However, naively combining these two terms could lead to an additional trade-off between the knowledge distillation and fairness, which requires additional hyperparameter tuning between the terms. In contrast, we devise a novel *single* regularizer that can simultaneously implement the knowledge distillation and fairness.

### 4.1. MMD-based Regularization for MFD

We approach distillation by matching feature distributions of each model as in [27], rather than minimizing instance-wise distances, since the distributional perspective is more proper for considering the group fairness at the same time. To formulate our regularization, we use the *max*-

imum mean discrepancy (MMD) [12], which measures the largest difference in expectations over functions in the unit ball  $\mathcal{F}$  of a reproducing kernel Hilbert space (RKHS)  $\mathcal{H}$ . For some distributions  $\mathbf{p}$  and  $\mathbf{q}$ , MMD is defined as follows:

$$D(\mathbf{p}, \mathbf{q}) \triangleq \sup_{f \in \mathcal{F}} (\mathbb{E}_{\mathbf{p}}[f(x)] - \mathbb{E}_{\mathbf{q}}[f(x')]) \quad (3)$$

$$= \|\mu_{\mathbf{p}} - \mu_{\mathbf{q}}\|_{\mathcal{H}}, \quad (4)$$

where  $\mu_{\mathbf{p}} \triangleq \mathbb{E}_{\mathbf{p}}[\phi(x)]$ ,  $\phi : \mathcal{X} \rightarrow \mathbb{R}$  is defined as  $\phi(\cdot) \triangleq k(\cdot, x)$  and  $k(\cdot, \cdot)$  is a kernel inducing  $\mathcal{H}$ . Under universal RKHS  $\mathcal{H}$ , MMD is a well-defined metric since it is proven that MMD value is zero if and only if  $\mathbf{p} = \mathbf{q}$  [12]. In this work, we use the Gaussian RBF kernel, well-known as the kernel that induces a universal RKHS, i.e.,  $k(x, x') = \exp(-\frac{1}{2\sigma^2} \|x - x'\|^2)$ .

The key for accomplishing our goal is to exploit and distill the group-indistinguishable features of the teacher  $T$ . It, however, is challenging to explicitly identify them in practice, and thus, we instead adopt a trick to use the per-class feature distribution of  $T$  as a target to distill while learning the group-conditioned feature distribution of  $S$ . Namely, we define our regularization term as

$$\mathcal{L}_{MFD} \triangleq \sum_y \sum_a D^2(\mathbf{p}_y^T, \mathbf{p}_{a,y}^S), \quad (5)$$

in which  $\mathbf{p}_y^T = \mathbb{E}_A[\mathbf{p}_{A,y}^T]$  is the group-averaged feature distribution of  $T$  for class  $y$ , and  $\mathbf{p}_{a,y}^S$  is the group-conditioned feature distribution of  $S$  for class  $y$  and the *sensitive* group (attribute)  $a$ . The rationale behind using  $\mathbf{p}_y^T$  as a target is that, by taking average across the groups, we expect the group-specific features would wash out while the common, group-agnostic predictive features would remain.

In Section 4.3, we give a theoretical analysis that minimizing  $\mathcal{L}_{MFD}$  can simultaneously have the knowledge distillation effect and promote fairness of the student  $S$ . Namely, we show that it leads to assimilating  $\mathbf{p}^T$  and  $\mathbf{p}^S$  (thus, KD effect) and reduces the distances among  $\mathbf{p}_{a,y}^S$  for all  $a \in \mathcal{A}$  (thus, fairness effect) by having the common distillation target  $\mathbf{p}_y^T$ . Furthermore, we note that considering the *class-wise* MMD in (5) fits well with the equalized odds metric that we consider.

## 4.2. Objective Function

Based on the rationale on  $\mathcal{L}_{MFD}$  described above, we design the final objective for training  $S$  as follows:

$$\min_{\theta} \mathcal{L}_{CE}(\theta) + \frac{\lambda}{2} \hat{\mathcal{L}}_{MFD}(\theta), \quad (6)$$

where  $\theta$  is the model parameter for the student  $S$ . In Eq.(6),  $\mathcal{L}_{CE}(\theta)$  denotes the ordinary cross entropy loss, and  $\lambda$  is a tunable hyperparameter that sets the trade-off between accuracy and fairness.  $\hat{\mathcal{L}}_{MFD}(\theta) \triangleq \sum_y \sum_a \hat{D}^2(\mathbf{p}_y^T, \mathbf{p}_{a,y}^S(\theta))$

is the empirical estimate of  $\mathcal{L}_{MFD}$ , in which the summand is defined as

$$\begin{aligned} \hat{D}^2(\mathbf{p}_y^T, \mathbf{p}_{a,y}^S(\theta)) &= \frac{1}{N_1^2} \sum_{i=1}^{N_1} \sum_{j=1}^{N_1} k(x_i, x_j) \\ &+ \frac{1}{N_2^2} \sum_{i=1}^{N_2} \sum_{j=1}^{N_2} k(x'_i(\theta), x'_j(\theta)) - \frac{2}{N_1 N_2} \sum_{i=1}^{N_1} \sum_{j=1}^{N_2} k(x_i, x'_j(\theta)), \end{aligned}$$

where  $x, x'(\theta)$  are the feature vectors sampled according to  $\mathbf{p}_y^T$  and  $\mathbf{p}_{a,y}^S(\theta)$ , respectively. Here,  $\hat{\mathcal{L}}_{MFD}$  can be applied to several layers of deep neural network, but we study only the case of applying to the penultimate layer throughout our work. In summary, (6) looks similar to the typical in-processing methods that employ additional fairness regularization, but our MFD also utilizes the information from the teacher  $T$ . Finally, we give a pictorial summary of our method in Figure 2.

**Mini-batch optimization** For the mini-batch stochastic descent, a standard optimization method for neural networks, we calculate the  $(a, y)$ -pairwise MMD using data points in a mini-batch. But, for a certain group-label pair  $(a, y)$ , the mean of the pair's conditional distribution in MMD can be biased if a mini-batch has few points for the pair. Hence, we strategically sample the data points with replacement to make a mini-batch in which the data points for each pair are contained with the same proportion. Furthermore, we set the kernel parameter  $\sigma^2$  as the mean of squared distance between all data points for each pair to maintain the stability.

## 4.3. Analysis

In this section, we give a theoretical justification of our MFD. We first show that minimizing  $\mathcal{L}_{MFD}$  leads to distributional matching of  $T$  and  $S$ .

**Lemma 1 (Knowledge Distillation)**

$$\sum_y \sum_a p(a, y) D^2(\mathbf{p}_y^T, \mathbf{p}_{a,y}^S) \geq D^2(\mathbf{p}^T, \mathbf{p}^S). \quad (7)$$

**Proof :** The proof follows from the following chain of inequalities

$$\begin{aligned} &\sum_y \sum_a p(a, y) D^2(\mathbf{p}_y^T, \mathbf{p}_{a,y}^S) \\ &= \sum_y \sum_a p(a, y) \left( \sup_{f \in \mathcal{F}} (\mathbb{E}_{\mathbf{p}_y^T}[f(x)] - \mathbb{E}_{\mathbf{p}_{a,y}^S}[f(x')]) \right)^2 \\ &\geq \left( \sum_y \sum_a p(a, y) \sup_{f \in \mathcal{F}} (\mathbb{E}_{\mathbf{p}_y^T}[f(x)] - \mathbb{E}_{\mathbf{p}_{a,y}^S}[f(x')]) \right)^2 \\ &\geq \left( \sup_{f \in \mathcal{F}} \left( \sum_y \sum_a p(a, y) (\mathbb{E}_{\mathbf{p}_y^T}[f(x)] - \mathbb{E}_{\mathbf{p}_{a,y}^S}[f(x')]) \right) \right)^2 \\ &= D^2(\mathbf{p}^T, \mathbf{p}^S), \end{aligned} \quad (8)$$



in which the first inequality follows from  $x^2$  being an increasing convex function for  $x \geq 0$  and using Jensen's inequality, and the second inequality follows from the subadditivity of supremum. ■

Note the LHS of Eq.(7) is equivalent to  $\mathcal{L}_{MFD}$  when  $p(a, y)$  is a uniform distribution. Therefore, the lemma shows that minimizing  $\mathcal{L}_{MFD}$  would also lead to the feature distribution of  $S$  get close to that of  $T$ , which is the knowledge distillation process.

Next, we investigate the relation between the  $\mathcal{L}_{MFD}$  and the equalized odds by introducing the following lemma.

**Lemma 2 (Fairness Constraints)** For every  $y \in \mathcal{Y}$ ,

$$\sum_a D^2(\mathbf{p}_y^T, \mathbf{p}_{a,y}^S) \geq \frac{1}{2|\mathcal{A}|} \sum_{a,a'} D^2(\mathbf{p}_{a,y}^S, \mathbf{p}_{a',y}^S). \quad (9)$$

**Proof :** Consider the followings:

$$\sum_a D^2(\mathbf{p}_y^T, \mathbf{p}_{a,y}^S) = \sum_a \left\| \mu_{\mathbf{p}_y^T} - \mu_{\mathbf{p}_{a,y}^S} \right\|_{\mathcal{H}}^2 \quad (10)$$

$$\geq \sum_a \left\| \mu_y^* - \mu_{\mathbf{p}_{a,y}^S} \right\|_{\mathcal{H}}^2 \quad (11)$$

$$= \frac{1}{2|\mathcal{A}|} \sum_{a,a'} \left\| \mu_{\mathbf{p}_{a,y}^S} - \mu_{\mathbf{p}_{a',y}^S} \right\|_{\mathcal{H}}^2 \quad (12)$$

$$= \frac{1}{2|\mathcal{A}|} \sum_{a,a'} D^2(\mathbf{p}_{a,y}^S, \mathbf{p}_{a',y}^S),$$

in which Eq.(11) follows from the fact that each  $\mu_y^* \triangleq \frac{1}{|\mathcal{A}|} \sum_a \mu_{\mathbf{p}_{a,y}^S}$  is the minimizer of each summand of Eq.(10). Eq.(12) follows from the equivalence between the sum of pairwise distance and the sum of distance to their mean. ■

From Lemma 2, we have that  $\sum_y \sum_{a,a'} D^2(\mathbf{p}_{a,y}^S, \mathbf{p}_{a',y}^S)$  is upper bounded by  $\mathcal{L}_{MFD}$  and equality holds when  $\mu_{\mathbf{p}_y^T} = \mu_y^*$  for all  $y$ . When the global optimum is achieved, *i.e.*,  $\mathcal{L}_{MFD} = 0$ , we get that  $\mathbf{p}_{a,y}^S$  is the same as  $\mathbf{p}_{a',y}^S$  for all  $a, a' \in \mathcal{A}, y \in \mathcal{Y}$ , which implies the independence between feature distribution of groups for given  $y$ , leading to the equalized odds condition,  $\tilde{Y} \perp A | Y = y$ .

## 5. Experimental Results

In the following section, we investigate our MFD can indeed reduce per-class accuracy discrepancy and improve accuracy in various object classification scenarios. We first consider a toy dataset, CIFAR-10S [35], and then experiment on two real-world datasets; age classification using UTKFace [42] and face attribute recognition using CelebA [22]. We describe the detailed experimental settings in the corresponding subsections.

**Baselines.** We compare our MFD with three classes of baselines. The first class is the ordinary KD methods, HKD

[16], FitNet [30], AT [38], and NST [17], that purely focus on improving the prediction accuracy via distillation. The second class is the state-of-the-art in-processing methods, AD [40] and SS [9], that explicitly take the fairness criterion into account while re-training the model. As described in Section 4.2, MFD also uses the same sampling strategy as SS, but we show through our experiments that merely controlling the ratio of group data points in a mini-batch can fail to reduce the unwanted discrimination of a model. The third class of baselines is the simple combination of the first two classes; namely, we combine the in-processing methods, AD or SS, with the KD methods, HKD or FitNet, by simply adding the distillation regularization terms to the objective functions of the in-processing methods. We show in our experiments that our MFD shows *consistent* improvements in *both* accuracy and fairness across *all three datasets*, while all other baselines cannot always improve both criteria on all datasets.

**Implementation details** For CIFAR-10S, we employed a simple convolutional neural network. Details of the network architecture is described in the Supplementary Material. For UTKFace [42] and CelebA [22], we adopted ImageNet-pretrained ResNet18 [15] and ShuffleNet [41], respectively. All algorithms were reproduced following the original papers using PyTorch [28]. Feature distillation was applied at the penultimate layer for all methods except for AT. Since AT is originally designed to transfer attention maps, we applied it to the feature after the last convolutional layer of the networks in each experiment. For AD, we omitted the loss projection in their original work due to the training instability in our experiments. We did the grid search for the hyperparameters of all methods sufficiently and chose the best one in terms of accuracy for the first class of baselines and DEO<sub>M</sub> for the second and third classes of baselines. We excluded the cases for which models achieve severely low accuracy despite their low DEO<sub>M</sub>. More details on training schemes and full hyperparameter settings are given in the Supplementary Material.

### 5.1. Synthetic Dataset

**Dataset** We adopted the CIFAR-10 Skewed (CIFAR-10S) dataset [35] which is a modified version of CIFAR-10 [20] in order to study bias mitigation in object classification. CIFAR-10 is a 10-way image classification dataset composed of  $32 \times 32$  images. In [35], the images of each class in the dataset are divided into two new domains (*i.e.*, two groups) of color and grayscale with a fixed ratio. They make the images of the first 5 classes be skewed towards color domain and the others towards grayscale, so that the total number of images belonging to each domain is balanced. Since each class data is skewed towards a specific domain, the extent of bias can be easily controlled by the skew ratio. Based on their protocol, we built CIFAR-10S

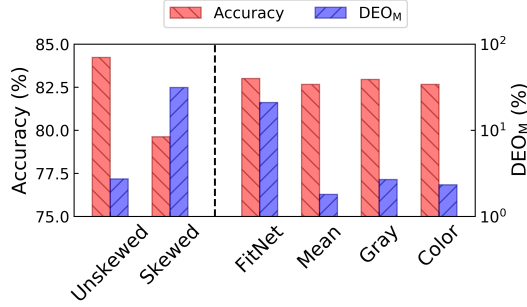


Figure 3. The effect of distillation using different feature information. Transferring the mean features of two domains helps the most in achieving fairness.  $DEO_M$  (in %) is reported in log scale.

dataset with the skewed ratio of 0.8. More specifically, from 5,000 per-class training images of the standard CIFAR-10, we set 4,000 images to grayscale and 1,000 to color for the first five classes and vice versa for the others. For the test set, we doubled the original CIFAR-10 test set by converting the images to grayscale and combining with the colored set, thereby we balanced the test set between two domains and have the pair of the same images; one in color, one in grayscale.

**Validating our motivation** Before testing the performance of our MFD on CIFAR-10S, we first carry out a toy experiment to validate our motivation described in the Introduction; namely, only distilling the group-indistinguishable predictive features from the teacher  $T$  to the student  $S$  should help improve both the accuracy and fairness of  $S$ .

To that end, we implemented an ideal *Unskewed* teacher and tested with a different choice of distilling features. For more details, we first constructed a Composite CIFAR-10 dataset that contains the original CIFAR-10 train set and the its grayscale version. Then, we trained two models from scratch with the Composite CIFAR-10 and CIFAR-10S, denoted as *Unskewed* and *Skewed*, respectively, in Figure 3. Note that the *Unskewed* model does not suffer from unfairness since it is trained on the balanced training set, while the *Skewed* model suffers from very high  $DEO_M$  due to the imbalance in CIFAR-10S described above. Now, setting the *Unskewed* model as the teacher  $T$  for the knowledge distillation, we trained four students  $S$  by fixing each student’s input as CIFAR-10S and changing the teacher’s input to the following four choices: 1) providing the same images as the student’s input (*FitNet*), 2) providing the color and grayscale image pair that corresponds to the student’s input, then distilling the mean of the two features (*Mean*), 3) providing only the grayscale version of the image that corresponds to the student’s input (*Gray*), and 4) providing only the original color image that corresponds to the student’s input (*Color*). We note that for the knowledge distillation, all approaches minimize  $L_2$  distance between features of the teacher and the student as in FitNet [30]. Here, *Mean* is intended to approximate the distillation with the group-

Table 1. The comparison of algorithms on CIFAR-10S dataset. The red and green arrows indicate that the performance got worse and better compared to the teacher, respectively. The numbers in parentheses represent how much they are changed from the value of the teacher, i.e., relative change in percentage (%).

Model	Accuracy ( $\uparrow$ )	$DEO_A$ ( $\downarrow$ )	$DEO_M$ ( $\downarrow$ )
Teacher	79.62	15.63	31.32
HKD [16]	80.34 (0.90 $\uparrow$ )	15.54 (0.58 $\downarrow$ )	34.12 (8.94 $\uparrow$ )
FitNet [30]	81.66 (2.56 $\uparrow$ )	14.83 (5.12 $\downarrow$ )	32.28 (3.07 $\uparrow$ )
AT [38]	79.00 (0.78 $\downarrow$ )	15.57 (0.38 $\downarrow$ )	31.25 (0.22 $\downarrow$ )
NST [17]	79.70 (0.10 $\uparrow$ )	15.11 (3.33 $\downarrow$ )	30.87 (1.44 $\downarrow$ )
SS [9]	82.69 (3.86 $\uparrow$ )	3.29 (78.95 $\downarrow$ )	7.13 (77.23 $\downarrow$ )
AD [40]	62.49 (21.51 $\downarrow$ )	11.59 (25.85 $\downarrow$ )	23.07 (26.34 $\downarrow$ )
SS+HKD	82.27 (3.33 $\uparrow$ )	10.15 (35.06 $\downarrow$ )	20.37 (34.96 $\downarrow$ )
SS+FitNet	81.73 (2.65 $\uparrow$ )	10.35 (33.78 $\downarrow$ )	20.92 (33.21 $\downarrow$ )
AD+HKD	79.27 (0.44 $\downarrow$ )	16.19 (3.58 $\uparrow$ )	33.25 (6.16 $\uparrow$ )
AD+FitNet	79.59 (0.04 $\downarrow$ )	15.90 (1.73 $\uparrow$ )	32.47 (3.67 $\uparrow$ )
MFD	<b>82.77 (3.96 <math>\uparrow</math>)</b>	<b>2.73 (82.53 <math>\downarrow</math>)</b>	<b>6.08 (80.59 <math>\downarrow</math>)</b>

indistinguishable informative features. *Gray* and *Color* are meant to further identify the effects on the student following the teacher’s features of one specific domain.

In Figure 3, we observe that all four methods utilizing the teacher’s knowledge succeed in improving the accuracy compared to the *Skewed*. Interestingly, *Mean*, *Gray* and *Color* also make significant improvements in fairness compared to *Skewed*, which just trains from scratch only using CIFAR-10S. Note that *Mean* achieves the lowest  $DEO_M$ , and we believe the reason for this improvement is that the unbiased, group-indistinguishable feature obtained by the mean feature from the teacher successfully mitigates the biased information, in line with our motivation given in the Introduction. In addition, we also believe the fairness gains of *Gray* and *Color* occur because providing the images of opposite domain for the half of CIFAR-10S train set has the effect of bias mitigation through distillation, so that the group-indistinguishable feature from *Unskewed* teacher can be distilled to the student. However, we also note that the amount of fairness improvement is smaller than *Mean*. In contrast, *FitNet* still suffers from high  $DEO_M$  despite the accuracy improvement, which exemplifies that a naive knowledge distillation may not be effective in mitigating the unfairness. Encouraged by this result, we now evaluate the performance of MFD on CIFAR-10S.

**Performance comparison** Table 1 shows the accuracy,  $DEO_A$ , and  $DEO_M$  (all in %) of the teacher (which is simply trained on CIFAR-10S), the students trained with the schemes from the three classes of baselines, and the student trained with our MFD on CIFAR-10S. We can make the following observations from the table. Firstly, MFD dominates all baselines, significantly improving both the accuracy and the fairness over the teacher. Secondly, we find that the knowledge distillation from the unfair teacher can exacerbate the discrimination (e.g.,  $DEO_M$ ) of the student while the accuracy is improved, e.g., HKD and FitNet. Thirdly,

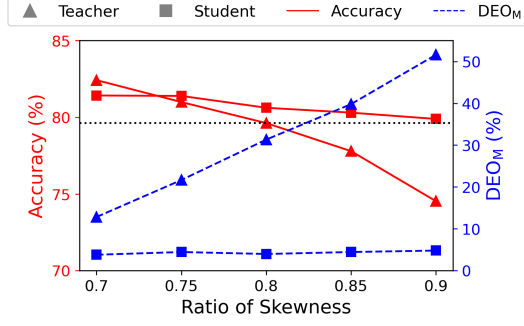


Figure 4. The level of student unfairness as the unfairness of teacher is intensified. Lower  $DEO_M$  indicates the model is fairer. Black dotted line indicates accuracy of the model trained from scratch with the skew rate of 0.8

for in-processing method baselines, we observe both AD and SS successfully improve the fairness, as expected, but AD worsens the accuracy. In case of SS, while it also improves the accuracy of the student, we observe our MFD outperforms it in terms of both accuracy and fairness. This reveals that our MFD effectively exploits the teacher by employing our regularizer  $\hat{\mathcal{L}}_{MFD}$ . Finally, we observe that a simple combination of in-processing methods with KD methods may either impair the accuracy or limit the fairness improvements.

**Distillation from unfaier teacher** As mentioned in above dataset subsection, the teacher in Table 1 was trained with the skew rate of 0.8. We now test the effect of the different level of unfairness of the teacher in the performance of the student trained with MFD.

Figure 4 shows the accuracy and  $DEO_M$  of the teachers that are trained with different skew rates (shown in the horizontal axis) of CIFAR-10S train set as well as those of the corresponding student that are trained with MFD employing each teacher. To see only the effect of differently biased teachers, we always fixed the skew rate of the train set for the student to 0.8. From the figure, we clearly see that as the skew rate increases, the teacher becomes increasingly unfair and inaccurate. In contrast, we observe that the student always achieves the higher accuracy than that of the model trained from scratch (*i.e.*, the teacher at the skew rate 0.8), even when the teacher MFD employs is heavily biased. Moreover, we observe the fairness of the student is significantly improved compared to the model trained from scratch and stays almost the same regardless of the unfairness level of the teacher. We believe this result corroborates our intuition that even when the original teacher is heavily biased, MFD can successfully distill the group-indistinguishable features from the teacher so that both the accuracy and fairness can be improved in the student.

**Feature visualization** To qualitatively investigate how MFD successfully reduces the discrimination, we visualize t-SNE embeddings of the teacher and the student trained

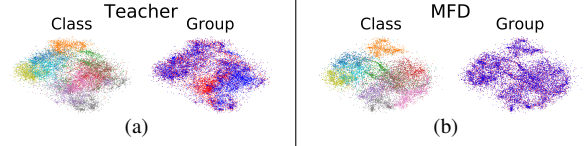


Figure 5. t-SNE [32] plots of features from CIFAR-10S test set.

Table 2. The comparison of algorithms on UTKFace dataset. Other settings are identical to Table 1.

Model	Accuracy ( $\uparrow$ )	$DEO_A$ ( $\downarrow$ )	$DEO_M$ ( $\downarrow$ )
Teacher	74.54	21.92	39.25
HKD [16]	<b>76.17 (2.19 <math>\uparrow</math>)</b>	22.50 (2.65 $\uparrow$ )	41.25 (5.10 $\uparrow$ )
FitNet [30]	75.23 (0.93 $\uparrow$ )	21.50 (1.92 $\downarrow$ )	40.00 (1.91 $\uparrow$ )
AT [38]	75.17 (0.85 $\uparrow$ )	22.67 (3.42 $\uparrow$ )	40.50 (3.18 $\uparrow$ )
NST [17]	75.10 (0.75 $\uparrow$ )	22.75 (3.79 $\uparrow$ )	42.00 (7.01 $\uparrow$ )
SS [9]	75.23 (0.93 $\uparrow$ )	24.33 (10.99 $\uparrow$ )	38.50 (1.91 $\downarrow$ )
AD [40]	74.67 (0.17 $\uparrow$ )	20.42 (6.84 $\downarrow$ )	36.00 (8.28 $\downarrow$ )
SS+HKD	76.08 (2.07 $\uparrow$ )	21.92 (0.00 $-$ )	37.50 (4.46 $\downarrow$ )
SS+FitNet	75.50 (1.29 $\uparrow$ )	21.92 (0.00 $-$ )	38.00 (3.18 $\downarrow$ )
AD+HKD	69.48 (6.79 $\downarrow$ )	18.75 (14.46 $\downarrow$ )	32.50 (17.20 $\downarrow$ )
AD+FitNet	70.23 (5.78 $\downarrow$ )	21.17 (3.42 $\downarrow$ )	33.75 (14.01 $\downarrow$ )
MFD	74.69 (0.20 $\uparrow$ )	<b>17.75 (19.02 <math>\downarrow</math>)</b>	<b>28.50 (27.39 <math>\downarrow</math>)</b>

Table 3. The comparison of algorithms on CelebA dataset. Other settings are identical to Table 1.

Model	Accuracy ( $\uparrow$ )	$DEO_A$ ( $\downarrow$ )	$DEO_M$ ( $\downarrow$ )
Teacher	78.33	21.04	21.81
HKD [16]	78.64 (0.40 $\uparrow$ )	21.56 (2.47 $\uparrow$ )	22.54 (3.35 $\uparrow$ )
FitNet [30]	78.62 (0.37 $\uparrow$ )	20.66 (1.81 $\downarrow$ )	21.70 (0.50 $\downarrow$ )
AT [38]	78.63 (0.38 $\uparrow$ )	21.28 (1.14 $\uparrow$ )	22.24 (1.97 $\uparrow$ )
SS [9]	79.67 (1.71 $\uparrow$ )	4.87 (76.85 $\downarrow$ )	5.22 (76.07 $\downarrow$ )
AD [40]	76.10 (2.85 $\downarrow$ )	<b>2.51 (88.07 <math>\downarrow</math>)</b>	<b>3.34 (84.69 <math>\downarrow</math>)</b>
SS+HKD	79.95 (2.07 $\uparrow$ )	8.41 (60.03 $\downarrow$ )	8.27 (62.08 $\downarrow$ )
SS+FitNet	79.77 (1.84 $\uparrow$ )	9.31 (55.75 $\downarrow$ )	8.61 (60.52 $\downarrow$ )
AD+HKD	80.31 (2.53 $\uparrow$ )	3.40 (83.84 $\downarrow$ )	4.05 (81.43 $\downarrow$ )
AD+FitNet	<b>80.60 (2.90 <math>\uparrow</math>)</b>	5.12 (75.67 $\downarrow$ )	5.51 (74.74 $\downarrow$ )
MFD	80.15 (2.32 $\uparrow$ )	5.46 (74.05 $\downarrow$ )	5.86 (73.13 $\downarrow$ )

with MFD in Figure 5 (a) and (b). In the figure, each point represents the feature vector of an image at the penultimate layer of the model used in Table 1. The points at the left and the right of (a) and (b) are colored according to its class (left) and group (right), respectively. Note MFD significantly reduces the distributional bias between the features for the grayscale (red) and color (blue) groups, while maintaining separability for the ten target classes. This visualization again shows that MFD can considerably mitigate the discrepancies between different groups, while maintaining information related to the classification task. Hyperparameters of t-SNE to reproduce the results in Figure 5 are provided in the Supplementary Material.

## 5.2. Real-world Datasets

We now consider two real-world scenarios; age classification and attribute recognition. For each scenario, we used UTKFace [42] and CelebA [22]; the former was used as a benchmark with multi-classes and multi-groups and the lat-

Table 4. Ablation study for MFD on all dataset. All tunable hyperparameter search proceeds the same way as Table 1. We reported MFD-K with the highest accuracy and, MFD-F and MFD with the best  $DEO_M$ .

	CIFAR-10S			UTKFace			CelebA		
	Accuracy	$DEO_A$	$DEO_M$	Accuracy	$DEO_A$	$DEO_M$	Accuracy	$DEO_A$	$DEO_M$
Teacher	79.62	15.63	31.32	74.54	21.92	39.25	78.33	21.04	21.81
MFD-K	80.13	14.70	29.83	<b>75.42</b>	21.67	38.5	78.43	21.19	20.59
MFD-F	82.45	2.98	6.18	72.42	19.50	35.00	79.84	<b>2.58</b>	<b>2.98</b>
MFD	<b>82.77</b>	<b>2.73</b>	<b>6.08</b>	74.69	<b>17.75</b>	<b>28.50</b>	<b>80.15</b>	5.46	5.86

ter was used to test on a larger scale data.

**Dataset** UTKFace is a face dataset containing more than 20,000 face images of different ethnicity over the age from 0 to 116. The ethnicity is originally composed of 5 different groups of *White*, *Black*, *Asian*, *Indian*, and *Others* including *Hispanic*, *Latino*, etc. We excluded *Others* from the dataset and used the remaining four race groups as sensitive attributes. We then divided the age range into three classes: ages between 0 to 19, 20 to 40, and ages more than 40. CelebA consists of more than 200,000 face images annotated with 40 binary attributes. Since the dataset has severe attribute imbalance, using multiple attributes would significantly reduce the test samples, hence, undermine the statistical significance of the results. Therefore, we only consider the binary group and binary class in our experiment; namely, we set *Gender* as the sensitive attribute and *Attractive* as the target variable, as in the work of Quadrianto *et al.* [29]. For unbiased evaluation of the accuracy and fairness, the test set was balanced by randomly taking the same number of images for each group and each class on both UTKFace and CelebA.

**Performance comparison** Table 2 and 3 evaluate the performance of various baselines as well as our MFD on the two real-world datasets. We omit the result for NST on CelebA due to computational limitations. We again confirm MFD considerably improves the fairness metrics,  $DEO_A$  and  $DEO_M$ , as well as the accuracy. For both datasets, we again observe the KD baselines improve the accuracy, as expected, but generally hurt the fairness of the teacher. The in-processing method baselines, SS and AD, and their KD-combined versions perform quite well on CelebA for both accuracy and fairness; however, we observe they show no or only little improvement in fairness on UTKFace, which is a multi-class, multi-group dataset. In contrast, we observe MFD *robustly* improves both the fairness and accuracy of the teacher regardless of the datasets.

### 5.3. Ablation Study

To further study the effectiveness of our regularization term,  $\hat{\mathcal{L}}_{MFD}$ , and verify our theoretical analyses, we evaluate the performance of the two variants of MFD, MFD-K and MFD-F, that only consider KD and fairness aspect, respectively. These variants utilize MMD-based regularization terms derived from our lemmas. Namely, MFD-K

adopts RHS in Lemma 1 as the regularization term to distill the knowledge from the teacher by minimizing the MMD loss between the feature distributions of teacher  $\mathbf{p}^T$  and student  $\mathbf{p}^S$  with no consideration of fairness. On the other hand, MFD-F trains a model *without* the teacher, using RHS in Lemma 2 as the regularization term, hence, no distillation. In our implementation of MFD-F, for the stable and efficient training, we substituted the class-wise, pairwise distance  $D^2(\mathbf{p}_{a,y}^S, \mathbf{p}_{a',y}^S)$  with the distance  $D^2(\mathbf{p}_y^S, \mathbf{p}_{a,y}^S)$ , and only used gradients obtained from  $\mathbf{p}_{a,y}^S$ .

Table 4 reports the accuracy and DEO metrics evaluated on all our datasets, for teacher, MFD-K, MFD-F and MFD. We also used the same mini-batch technique for MFD-F as MFD. Followings are our observations. Firstly, we observe MFD-K indeed improves the accuracy of the teacher, hence, it can be used as a yet another KD scheme. Secondly, we note that MFD-F creates fairer models than teacher, but it may lead to a slight loss of accuracy as in UTKFace. This implies that utilizing the teacher has a critical role in maintaining or improving the accuracy while training a fairer model. Finally, we clearly see that MFD is the only method that consistently makes fairer models than the teachers while improving accuracy over all datasets. Thus, we conclude that  $\hat{\mathcal{L}}_{MFD}$  is very effective in building a fair model via knowledge distillation, as verified in our lemmas.

## 6. Conclusion

We proposed a novel in-processing method, MFD, that can both improve the accuracy and fairness of an already deployed, unfair model via feature distillation. Namely, our novel MMD-based regularizer utilizes the group-indistinguishable predictive features from the teacher while promoting the student model to not discriminate against any protected groups. Throughout the theoretical justification and extensive experimental analyses, we showed that our MFD is very effective and robust across diverse datasets.

## Acknowledgment

This work was supported in part by NRF Mid-Career Research Program [NRF-2021R1A2C2007884] and IITP grant [No.2019- 0-01396, Development of framework for analyzing, detecting, mitigating of bias in AI model and training data], funded by the Korean government.



## References

- [1] Alekh Agarwal, Alina Beygelzimer, Miroslav Dudik, John Langford, and Hanna Wallach. A reductions approach to fair classification. In *International Conference on Machine Learning (ICML)*, 2018. 1, 2
- [2] Sungsoo Ahn, Shell Xu Hu, Andreas Damianou, Neil D Lawrence, and Zhenwen Dai. Variational information distillation for knowledge transfer. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [3] Wael Alghamdi, Shahab Asoodeh, Hao Wang, Flavio P Calmon, Dennis Wei, and Karthikeyan Natesan Ramamurthy. Model projection: Theory and applications to fair machine learning. In *IEEE International Symposium on Information Theory (ISIT)*, 2020. 2
- [4] Rachel KE Bellamy, Kuntal Dey, Michael Hind, Samuel C Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilovic, et al. Ai fairness 360: An extensible toolkit for detecting, understanding, and mitigating unwanted algorithmic bias. *arXiv preprint arXiv:1810.01943*, 2018. 2
- [5] Richard Berk, Hoda Heidari, Shahin Jabbari, Michael Kearns, and Aaron Roth. Fairness in criminal justice risk assessments: The state of the art. *arXiv preprint arXiv:1703.09207*, 2017. 2
- [6] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on Fairness, Accountability and Transparency (FAccT)*, 2018. 1
- [7] Elliot Creager, David Madras, Joern-Henrik Jacobsen, Marissa Weis, Kevin Swersky, Toniann Pitassi, and Richard Zemel. Flexibly fair representation learning by disentanglement. In *International Conference on Machine Learning (ICML)*, 2019. 2
- [8] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Conference on Innovations in Theoretical Computer Science (ITCS)*, 2012. 2
- [9] Charles Elkan. The foundations of cost-sensitive learning. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2001. 2, 5, 6, 7
- [10] Tommaso Furlanello, Zachary Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar. Born again neural networks. In *International Conference on Machine Learning (ICML)*, 2018. 3
- [11] Clare Garvie. *The perpetual line-up: Unregulated police face recognition in America*. Georgetown Law, Center on Privacy & Technology, 2016. 1
- [12] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(1):723–773, 2012. 4
- [13] Varun Gulshan, Lily Peng, Marc Coram, Martin C Stumpe, Derek Wu, Arunachalam Narayanaswamy, Subhashini Venugopalan, Kasumi Widner, Tom Madams, Jorge Cuadros, et al. Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. *Jama*, 316(22):2402–2410, 2016. 1
- [14] Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016. 2, 3
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 5
- [16] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 1, 2, 3, 5, 6, 7
- [17] Zehao Huang and Naiyan Wang. Like what you like: Knowledge distill via neuron selectivity transfer. *arXiv preprint arXiv:1707.01219*, 2017. 1, 2, 5, 6, 7
- [18] Heinrich Jiang and Ofir Nachum. Identifying and correcting label bias in machine learning. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2020. 1, 2
- [19] Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. Fairness-aware classifier with prejudice remover regularizer. In *European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD)*, 2012. 2, 3
- [20] Alex Krizhevsky. Learning multiple layers of features from tiny images. *Master's thesis, University of Toronto*, 2009. 5
- [21] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 2
- [22] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *IEEE International Conference on Computer Vision (ICCV)*, 2015. 2, 5, 7
- [23] Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard S Zemel. The variational fair autoencoder. In *International Conference on Learning Representations (ICLR)*, 2016. 2
- [24] Iacopo Masi, Yue Wu, Tal Hassner, and Prem Natarajan. Deep face recognition: A survey. In *SIBGRAPI Conference on Graphics, Patterns and Images*, 2018. 1
- [25] Laurent Son Nguyen and Daniel Gatica-Perez. Hirability in the wild: Analysis of online conversational video resumes. *IEEE Transactions on Multimedia*, 18(7):1422–1437, 2016. 1
- [26] Wonpyo Park, Dongju Kim, Yan Lu, and Minsu Cho. Relational knowledge distillation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [27] Nikolaos Passalis and Anastasios Tefas. Learning deep representations with probabilistic knowledge transfer. In *European Conference on Computer Vision (ECCV)*, 2018. 2, 3
- [28] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 5
- [29] Novi Quadrianto, Viktoriia Sharmanska, and Oliver Thomas. Discovering fair representations in the data domain. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8227–8236, 2019. 1, 2, 8

- [30] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fit-nets: Hints for thin deep nets. In *International Conference on Learning Representations (ICLR)*, 2015. 1, 2, 3, 5, 6, 7
- [31] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. In *International Conference on Learning Representations (ICLR)*, 2020. 1, 2
- [32] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(86):2579–2605, 2008. 7
- [33] Mei Wang and Weihong Deng. Mitigating bias in face recognition using skewness-aware reinforcement learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [34] Mei Wang, Weihong Deng, Jiani Hu, Xunqiang Tao, and Yaohai Huang. Racial faces in the wild: Reducing racial bias by information maximization adaptation network. In *IEEE International Conference on Computer Vision (CVPR)*, 2019. 1, 2
- [35] 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 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 5
- [36] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [37] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rogriguez, and Krishna P. Gummadi. Fairness constraints: Mechanisms for fair classification. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2017. 1, 2, 3
- [38] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *International Conference on Learning Representations (ICLR)*, 2017. 1, 2, 5, 6, 7
- [39] Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations. In *International Conference on Machine Learning (ICML)*, 2013. 2
- [40] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating unwanted biases with adversarial learning. In *AAAI/ACM Conference on AI, Ethics, and Society (AI/ES)*, 2018. 2, 5, 6, 7
- [41] X Zhang, X Zhou, M Lin, and J Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 5
- [42] Zhifei Zhang, Yang Song, and Hairong Qi. Age progression/regression by conditional adversarial autoencoder. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2, 3, 5, 7

# Supplementary Materials for Fair Feature Distillation for Visual Recognition

Sangwon Jung<sup>1\*</sup>, Donggyu Lee<sup>1\*</sup>, Taeon Park<sup>1\*</sup> and Taesup Moon<sup>2†</sup>

<sup>1</sup>Department of Electrical and Computer Engineering, Sungkyunkwan University, Suwon, Korea

<sup>2</sup>Department of Electrical and Computer Engineering, Seoul National University, Seoul, Korea

{s.jung, ldk308, pte1236}@skku.edu, tsmoon@snu.ac.kr

## 1. Implementation Details

For all datasets, we trained all methods for 50 epochs with a mini-batch size of 128 using Adam optimizer with an initial learning rate 0.001 and decaying it by a factor of 10 if no improvement in the test loss for 10 consecutive epochs. Also, all results were averaged over 4 different random runs.

### 1.1. Network Architecture for CIFAR-10S

We employed a simple convolutional neural network having six convolutional layers with the kernel size of  $3 \times 3$ , followed by two fully connected hidden layers with ReLU [1] activations. The number of channels was set to 32, 32, 64, 64, 128, and 128 for each convolutional layer, respectively. Dropout [4] and max-pooling were applied after every two convolutional layers.

### 1.2. Hyperparameters for Main Results

For fair comparison, we did the extensive search for one hyperparameter of each method including ours and baselines. We set one parameter to search for and fixed others using a suitable strategy for baselines having more than two hyperparameters. For HKD and FitNet, we focused on finding the optimal  $T$ , a temperature to soften the output, while we gradually annealed the strength of output distillation for both methods and fixed feature distillation strength for FitNet to 1 like in [3]. For AD, we tune the strength of the adversary loss while fixing the learning rate of it to 0.003, a commonly used value. For the variants of SS, we do the same search strategy as the knowledge distillation methods. For the variants of AD, we fixed the all hyperparameters of the knowledge distillation to the best values found in the experiment for single distillations and searched the strength of the adversary loss to control the balance between two methods. The values for hyperparameters used to report the results in the manuscript are in Table 1. In Table 1, we denote the strength for each method as  $\lambda$ .

### 1.3. Details on AD+FitNet

Three combined methods of the third class of baselines in the manuscript, except for AD+FitNet, are naturally implemented, but implementing AD+FitNet requires modification to FitNet. More specifically, FitNet originally has two stages of training, the hint training for feature distillation and the KD training for output distillation. However, since this stage-wise training of FitNet has difficulty to being incorporated with the mini-max game with an adversary in AD, we modify the two stages training FitNet to one stage FitNet by minimizing the output and feature distillation loss simultaneously, as in [5]. Then, we integrate the loss of an adversary of AD into the loss of one stage FitNet to implement AD+FitNet.

### 1.4. Hyperparameters for t-SNE

Hyperparameters of t-SNE feature visualization for (Figure 5, manuscript) are as follows : dimension of the embedded space (3), perplexity (200), early exaggeration(1.0), maximum number of iterations (250), metric (cosine), random state (5). For other factors, we remained default in scikit-learn [2].

Table 1. Hyperparameters for experiments.

Methods\Dataset	CIFAR-10S	UTKFace	CelebA
HKD	$T$ (1)	$T$ (3)	$T$ (5)
FitNet	$T$ (1)	$T$ (5)	$T$ (1)
AT	$\lambda$ (1)	$\lambda$ (30)	$\lambda$ (1)
NST	$\lambda$ (30)	$\lambda$ (3)	-
AD	$\lambda$ (0.001)	$\lambda$ (0.01)	$\lambda$ (10)
SS+HKD	$T$ (3)	$T$ (5)	$T$ (3)
SS+FitNet	$T$ (3)	$T$ (10)	$T$ (3)
AD+HKD	$T$ (1) $\lambda$ (1e-4)	$T$ (3) $\lambda$ (30)	$T$ (5) $\lambda$ (10)
AD+FitNet	$T$ (1) $\lambda$ (1e-3)	$T$ (5) $\lambda$ (1)	$T$ (1), $\lambda$ (1)
MFD	$\lambda$ (3)	$\lambda$ (3)	$\lambda$ (7)

## 2. Result Tables

Table 2, 3 and 4 show the detail results. The number in the parenthesis with  $\pm$  sign stands for the standard deviation of each metric obtained from 4 independent runs.

\*Equal contribution.

†Corresponding author (E-mail: tsmoon@snu.ac.kr)

Table 2. Average accuracy (%) and DEO (%) with standard deviation on CIFAR-10S.

	Accuracy	DEO <sub>A</sub>	DEO <sub>M</sub>
Teacher	79.62 ( $\pm 0.14$ )	15.63 ( $\pm 0.44$ )	31.32 ( $\pm 1.47$ )
HKD	80.34 ( $\pm 0.35$ )	15.54 ( $\pm 0.67$ )	34.12 ( $\pm 2.21$ )
FitNet	81.66 ( $\pm 0.20$ )	14.83 ( $\pm 0.26$ )	32.28 ( $\pm 1.59$ )
AT	79.00 ( $\pm 0.99$ )	15.57 ( $\pm 0.71$ )	31.25 ( $\pm 1.20$ )
NST	79.70 ( $\pm 0.99$ )	15.11 ( $\pm 0.75$ )	30.87 ( $\pm 2.38$ )
SS	82.69 ( $\pm 0.22$ )	3.29 ( $\pm 0.30$ )	7.13 ( $\pm 1.36$ )
AD	62.49 ( $\pm 30.32$ )	11.59 ( $\pm 6.75$ )	23.07 ( $\pm 13.36$ )
SS+HKD	82.27 ( $\pm 0.33$ )	10.15 ( $\pm 0.20$ )	20.37 ( $\pm 1.14$ )
SS+FitNet	81.73 ( $\pm 0.39$ )	10.35 ( $\pm 0.47$ )	20.92 ( $\pm 0.54$ )
AD+HKD	79.27 ( $\pm 0.33$ )	16.19 ( $\pm 0.50$ )	33.25 ( $\pm 0.72$ )
AD+FitNet	79.59 ( $\pm 0.37$ )	15.90 ( $\pm 0.51$ )	32.47 ( $\pm 1.66$ )
MFD	<b>82.77 (<math>\pm 0.14</math>)</b>	<b>2.73 (<math>\pm 0.41</math>)</b>	<b>6.08 (<math>\pm 0.91</math>)</b>

Table 3. Average accuracy (%) and DEO (%) with standard deviation on UTKFace.

	Accuracy	DEO <sub>A</sub>	DEO <sub>M</sub>
Teacher	74.54 ( $\pm 1.07$ )	21.92 ( $\pm 1.36$ )	39.25 ( $\pm 2.86$ )
HKD	<b>76.17 (<math>\pm 0.58</math>)</b>	22.5 ( $\pm 0.76$ )	41.25 ( $\pm 3.49$ )
FitNet	75.23 ( $\pm 0.52$ )	21.50 ( $\pm 1.59$ )	40.00 ( $\pm 4.64$ )
AT	75.17 ( $\pm 0.82$ )	22.67 ( $\pm 3.41$ )	40.50 ( $\pm 6.87$ )
NST	75.10 ( $\pm 0.39$ )	22.75 ( $\pm 0.49$ )	42.00 ( $\pm 4.18$ )
SS	75.23 ( $\pm 0.87$ )	24.33 ( $\pm 1.75$ )	38.50 ( $\pm 2.29$ )
AD	74.67 ( $\pm 1.01$ )	20.42 ( $\pm 1.55$ )	36.00 ( $\pm 2.55$ )
SS+HKD	76.08 ( $\pm 0.42$ )	21.92 ( $\pm 1.07$ )	37.50 ( $\pm 2.05$ )
SS+FitNet	75.5 ( $\pm 0.99$ )	21.92 ( $\pm 1.75$ )	38.00 ( $\pm 2.06$ )
AD+HKD	69.48 ( $\pm 3.21$ )	18.75 ( $\pm 1.93$ )	32.50 ( $\pm 4.15$ )
AD+FitNet	70.23 ( $\pm 6.64$ )	21.17 ( $\pm 6.03$ )	33.75 ( $\pm 6.06$ )
MFD	74.69 ( $\pm 0.69$ )	<b>17.75 (<math>\pm 1.38</math>)</b>	<b>28.50 (<math>\pm 1.80</math>)</b>

Table 4. Average accuracy (%) and DEO (%) with standard deviation on CelebA.

	Accuracy	DEO <sub>A</sub>	DEO <sub>M</sub>
Teacher	78.33 ( $\pm 0.08$ )	21.04 ( $\pm 0.48$ )	21.81 ( $\pm 0.13$ )
HKD	78.64 ( $\pm 0.37$ )	21.56 ( $\pm 0.92$ )	22.54 ( $\pm 0.60$ )
FitNet	78.62 ( $\pm 0.20$ )	20.66 ( $\pm 0.81$ )	21.70 ( $\pm 0.56$ )
AT	78.63 ( $\pm 0.22$ )	21.28 ( $\pm 0.28$ )	22.24 ( $\pm 0.51$ )
SS	79.67 ( $\pm 0.36$ )	4.87 ( $\pm 0.69$ )	5.22 ( $\pm 0.81$ )
AD	76.10 ( $\pm 1.12$ )	<b>2.51 (<math>\pm 2.12</math>)</b>	<b>3.34 (<math>\pm 3.09</math>)</b>
SS+HKD	79.95 ( $\pm 0.42$ )	8.41 ( $\pm 1.78$ )	8.27 ( $\pm 1.83$ )
SS+FitNet	79.77 ( $\pm 0.28$ )	9.31 ( $\pm 1.77$ )	8.61 ( $\pm 2.23$ )
AD+HKD	80.31 ( $\pm 0.30$ )	3.40 ( $\pm 2.46$ )	4.05 ( $\pm 2.86$ )
AD+FitNet	<b>80.60 (<math>\pm 0.14</math>)</b>	5.12 ( $\pm 1.67$ )	5.51 ( $\pm 1.64$ )
MFD	80.15 ( $\pm 0.29$ )	5.46 ( $\pm 0.95$ )	5.86 ( $\pm 0.83$ )

## References

- [1] Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. In *International Conference on Machine Learning (ICML)*, 2010. [1](#)
- [2] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,

12:2825–2830, 2011. [1](#)

- [3] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fit-nets: Hints for thin deep nets. In *International Conference on Learning Representations (ICLR)*, 2015. [1](#)
- [4] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958, 2014. [1](#)
- [5] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. In *International Conference on Learning Representations (ICLR)*, 2020. [1](#)