

Meta Balanced Network for Fair Face Recognition

Mei Wang, Yaobin Zhang, Weihong Deng

Abstract—Although deep face recognition has achieved impressive progress in recent years, controversy has arisen regarding discrimination based on skin tone, questioning their deployment into real-world scenarios. In this paper, we aim to systematically and scientifically study this bias from both data and algorithm aspects. First, using the dermatologist approved Fitzpatrick Skin Type classification system and Individual Typology Angle, we contribute a benchmark called Identity Shades (IDS) database, which effectively quantifies the degree of the bias with respect to skin tone in existing face recognition algorithms and commercial APIs. Further, we provide two skin-tone aware training datasets, called BUPT-Globalface dataset and BUPT-Balancedface dataset, to remove bias in training data. Finally, to mitigate the algorithmic bias, we propose a novel meta-learning algorithm, called Meta Balanced Network (MBN), which learns adaptive margins in large margin loss such that the model optimized by this loss can perform fairly across people with different skin tones. To determine the margins, our method optimizes a meta skewness loss on a clean and unbiased meta set and utilizes backward-on-backward automatic differentiation to perform a second order gradient descent step on the current margins. Extensive experiments show that MBN successfully mitigates bias and learns more balanced performance for people with different skin tones in face recognition. The proposed datasets are available at <http://www.whdeng.cn/RFW/index.html>.

Index Terms—fairness with respect to skin tone, meta learning, adaptive margin, face recognition.

1 INTRODUCTION

RECENTLY, with the emergence of deep convolutional neural networks (CNN) [1], [2], [3], [4], [5], research focus of face recognition (FR) has shifted to deep-learning-based approaches [6], [7], [8] and the accuracy was dramatically boosted to above 99.80% on the Labeled Faces in the Wild (LFW) dataset [9]. However, the recognition accuracy is not only aspect to attend when designing learning algorithms. As a growing number of applications based on FR have integrated into our lives, its potential for unfairness is raising alarm. For example, Amazon’s Rekognition Tool incorrectly matched the photos of 28 U.S. congressmen with the faces of criminals, especially the error rate was up to 39% for Black faces. According to these reports [10], [11], FR system seems discriminative based on classes like race, demonstrating significantly different accuracy when applied to different groups. Such bias can result in mistreatment of certain demographic groups, by either exposing them to a higher risk of fraud, or by making access to services more difficult. Consequently, there is an increased need to guarantee fairness for automatic systems and prevent discriminatory decisions.

Although several studies [20], [21] have uncovered such discrimination in non-deep FR algorithms, there are still no sufficient research efforts in deep learning era. Without a dataset that has demographic labels for various people, one cannot

systematically examine the inappropriate biases in trained models. To facilitate the research towards this issue, in this paper, we have done the first step to overcome the major obstacle. Considering race labels are unstable, we decided to use skin tone as a more precise and scientific label. With the help of the Fitzpatrick Skin Type classification system [22] and Individual Typology Angle [23], [24], a new test dataset, called Identity Shades (IDS), is constructed which is phenotypically balanced on the basis of skin tone, as shown in Fig. 1. It can be used to fairly evaluate FR algorithms across faces with different skin tones. Based on experiments on IDS, we find that both commercial APIs and state-of-the-art (SOTA) algorithms indeed suffer from bias: the error rates on dark-skinned faces are about two times of the light-skinned ones, as shown in Table 1. Moreover, we demonstrated that this bias comes from both data and algorithm aspects. Hence, further research efforts on data and algorithms are requested to eliminate this bias.

A major driver of bias in FR is the training data. Large-scale datasets, such as CASIA-WebFace [25], VGGFace2 [19] and MS-Celeb-1M [26], are typically constructed by scraping websites like Google Images. Such data collecting methods can unintentionally produce data that encode biases with respect to skin tone. Thus, social awareness must be brought to the building of datasets for training. In this work, we take steps to ensure such datasets are diverse and do not under represent particular skin tone groups by constructing two new training datasets, i.e., BUPT-Globalface and BUPT-Balancedface dataset. One is built up according to the approximated distribution of average skin tones around the world, and the other strictly balances the number of samples in skin tone.

Another source of bias can be traced to the algorithms. The state-of-the-art (SOTA) face recognition methods, such as Cosface [27] and Arcface [18], apply a fixed margin between classes to maximize overall prediction accuracy for the training

- Mei Wang and Yaobin Zhang are with the Pattern Recognition and Intelligent System Lab., School of Artificial Intelligence, Beijing University of Posts and Telecommunications, Beijing, China. E-mail: {wangmei1,zhangyaobin}@bupt.edu.cn.
- Weihong Deng is with the Pattern Recognition and Intelligent System Lab., School of Artificial Intelligence, Beijing University of Posts and Telecommunications, Beijing, China, and also with Key Lab. of Trustworthy Distributed Computing and Service, Ministry of Education, Beijing University of Posts and Telecommunications, Beijing, China. E-mail: whdeng@bupt.edu.cn.

(Corresponding author: Weihong Deng)

TABLE 1

Bias with respect to skin tone in commercial APIs and SOTA FR algorithms. Verification accuracies (%) on our IDS-8 are given. Skin gradually darkens with the increase of tone value (from I to VIII). We make the best results bold, and make the worst in red.

| Model | The skin tone of IDS-8 | | | | | | | |
|----------------------------|------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | I | II | III | IV | V | VI | VII | VIII |
| Microsoft [12] | 88.67 | 87.15 | 80.20 | 78.45 | 83.13 | 82.09 | 76.10 | 74.90 |
| Face++ [13] | 93.68 | 93.89 | 92.45 | 92.49 | 88.58 | 89.36 | 88.23 | 87.40 |
| Baidu [14] | 90.52 | 88.04 | 90.52 | 89.68 | 86.72 | 87.19 | 78.18 | 78.62 |
| Amazon [15] | 91.22 | 90.12 | 84.96 | 85.13 | 86.85 | 88.40 | 85.79 | 86.36 |
| mean | 91.02 | 89.80 | 87.03 | 86.44 | 86.32 | 86.76 | 82.07 | 81.82 |
| Center-loss [16] | 87.59 | 87.14 | 79.21 | 77.87 | 81.90 | 83.68 | 78.55 | 77.51 |
| Sphereface [17] | 90.59 | 90.71 | 82.75 | 82.42 | 85.50 | 88.06 | 81.95 | 81.82 |
| Arcface ¹ [18] | 92.33 | 91.70 | 83.97 | 83.05 | 87.46 | 88.40 | 84.07 | 83.23 |
| VGGface2 ¹ [19] | 90.22 | 89.66 | 84.93 | 83.82 | 86.03 | 87.29 | 83.37 | 82.72 |
| mean | 90.18 | 89.80 | 82.71 | 81.79 | 85.22 | 86.86 | 81.98 | 81.32 |

¹ Arcface is a ResNet-34 model trained with CASIA-Webface. VGGFace2 is a SeNet model trained with VGGFace2.

data. However, maximizing overall accuracy might come at the expense of the under-represented populations and leads to poor performance for those subjects, which even amplifies the biases in data. To address this problem, the algorithms must trade-off the specific requirements of margins of various groups of people, and set adaptive margins for faces with different skin tones to produce more equitable recognition performance. But this demographic bias is a complex problem caused by many latent factors, including but not limited to the quantity, it's quite difficult to manually preset or tune these margins by heuristic strategies like cross-validation, naturally conducting efficiency issue and difficulty in practical implementations.

In our paper, we propose a novel meta-learning algorithm, named Meta Balanced Network (MBN), which automatically learns adaptive margin for each skin tone group based on their gradients from meta data. We additionally use a small but unbiased validation set, i.e., meta data, to guide training in our MBN. The meta data can help the model to distinguish fairness from bias and learn high-level causal relationships in bias. During training, we treat model optimization as the objective of inner algorithm trained by training data, and treat margin optimization as the objective of outer algorithm trained by meta data. The outer algorithm can evaluate the bias of the learned model on meta data and dynamically output margins in adaptive margin loss function; and the inner algorithm aims to optimize the model guided by adaptive margin loss function on training data such that the model performs well and fairly across faces with different skin tones. A specific backward-on-backward differentiation enables us to connect disjoint process between inner algorithm and outer algorithm and transmit meta gradient from meta loss to training loss to update the margins. This bilevel-optimization realizes a mutual amelioration between automatically tuning margin parameters involved in adaptive margin loss and learning suitable model parameters leading to balanced performance.

Our contributions can be summarized into three aspects.

1) We construct and release three large-scale in-the-wild datasets, i.e., IDS, BUPT-Globalface and BUPT-Balancedface, which are balanced by skin tone. They are the first series of databases for studying bias with respect to skin tone in both training and evaluating aspects. Based on the extensive experiments on them, we not only measure the bias in commercial systems and deep recognition models, but also validate the bias comes from both data and algorithm.

2) To mitigate the algorithmic bias, we propose to trade

off the specific requirements of margins of various people. Our novel MBN method leverages an additional small meta set to automatically learn the optimal margins by utilizing backward-on-backward automatic differentiation to take a second order gradient pass in meta-optimization. To the best of our knowledge, this is the first time that meta learning is used to solve the fairness of recognition problem.

3) Extensive experiments on the Globalface, Balancedface, and IDS datasets show that our meta balanced network (MBN) shows more balanced performance on different skin tone subjects than traditional sample re-weighting method, adversarial attribute removal method and our recently proposed reinforcement marginal learning method [28].

The remainder of this paper is structured as follows. In the next section, we discuss the studies of bias, and also review the related approaches of debiasing algorithms and meta learning. Then, we introduce the details of our databases in Section III. In Section IV, we propose the meta balanced network, i.e., MBN, to learn balanced performance for different skin tone groups. In Section V, experimental results on IDS are shown and validate that the existence and cause of such bias. Then we evaluate the effectiveness of our MBN method. Finally, we conclude and discuss future work.

2 RELATED WORK

2.1 Bias with respect to skin tone in face recognition

Some studies [29], [30], [31] have raised concerns about face analysis systems, e.g., gender classification, being biased based on some classes like race. The study of such bias in face recognition, likewise, has a nearly 30-year history that converges on the following aspects.

Algorithms. There are published studies [20], [21], [32], [33] analyzing the performances of face recognition algorithms over demographic groups and uncovering that these algorithms suffer from bias. The 2002 NIST Face Recognition Vendor Test (FRVT) [33] showed that non-deep FR algorithms have different recognition accuracies depending on demographic cohort, such as skin color, age and gender. Phillips et al. [20] suggested that training and testing on different races results in severe performance drop. Klare et al. [21] evaluated six different FR algorithms on three skin tone groups, and concluded that the Black cohorts are more difficult to recognize for all matchers. The FRVT 2019 [34] showed the demographic bias of over 100 face recognition algorithms. Krishnapriya et al. [35], [36] found

that darker-skinned subjects have a higher false match rate, and lighter-skinned ones have a higher false nonmatch rate through experiments on MORPH dataset [37].

Datasets. In non-deep learning era, some training and testing datasets are constructed for studying and measuring this demographic bias in FR algorithms. Klare et al. [21] collected mug shot face images of Whites, Blacks and Hispanics from the Pinellas County Sheriff’s Office (PCSO). Furl et al. [38] collected images of Caucasians, Asians, Africans, Indians and Hispanics from the FERET database [39]. Phillips et al. [20] utilized the images of FRVT 2006 [40] with a sufficient number of Caucasian and East Asian faces to conduct cross training and matching on White and Asian faces. However, few studies focus on bias in deep era because benchmark datasets are lacking to systematically examine the inappropriate biases in trained models. In this paper, we take a step in this direction by releasing three datasets. Considering race labels are unstable as we explained in Section 3.1, we decided to use skin tone as a more precise label.

2.2 Debiasing algorithms

There are some works that seek to introduce fairness into machine learning pipelines and mitigate data bias.

Unbalanced-training method. Sample reweighting methods [41], [42] have been believed to be an effective way in addressing class imbalance problem, which decreases weights for the over-sized classes and imposes larger weights on rarer data. For example, DB-VAE [31] used the learned latent distributions to re-weight the importance of data points. Other methods [28], [43] mitigate the bias via model regularization, taking into consideration of the fairness goal in the overall model objective function. For example, RL-RBN [28] formulated the process of finding the optimal margins for non-White people as a Markov decision process and employed deep Q-learning to learn policies based on large margin loss.

Attribute removal method. By confounding or removing demographic information of faces, attribute removal methods [44], [45], [46] make a great contribution to mitigating the hidden, and potentially unknown, biases within training data. Alvi et al. [44] applied a confusion loss to make a classifier fail to distinguish attributes of examples so that multiple spurious variations are removed from the feature representation. SensitiveNets [47] minimized sensitive information in triplet loss while retaining recognition ability. Gong et al. [48] debiased face recognition by disentangling features related to demographics and identity.

Domain adaptation. Some papers [49], [50], [51], [52] proposed to investigate bias problem from a domain adaptation point of view and attempted to learn domain-invariant feature representations to mitigate bias across domains. IMAN [49] simultaneously aligned global distribution to decrease gap at domain-level, and learned the discriminative target representations at cluster level. Kan [50] directly converted the White faces to non-White domain in the image space with the help of sparse reconstruction coefficients.

2.3 Meta learning

The field of meta learning, or learning to learn, has seen a dramatic rise in interest in recent years. Recent meta-learning studies concentrate on: 1) learning a good weight initialization for fast adaptation on a new task [53], [54], [55], [56]. For example, Ravi and Larochelle [57] proposed to learn the few-shot optimization

algorithm with an LSTM-based meta-learner. MAML [53] and its variants Reptile [54] simplified the above meta-learner model and only learned the initial learner parameters. MFR [58] synthesized the source/target domain shift with a meta-optimization objective to improve generalization in face recognition. 2) learning an adaptive weighting scheme [59], [60]. For example, to address label noise and long-tail problem, L2RW [59] adopted a novel meta-learning algorithm that learns to assign weights to training data based on their gradient directions. Different from learning weights explicitly in L2RW [59], Shu et al. [60] proposed to utilize meta learning to learn a weight function, i.e., Meta-Weight-Net, which can learn the weights in a more stable way. Our approach is most related to L2RW [59]. However, 1) L2RW used meta learning to re-weight training samples, while our MBN learns adaptive margin. 2) Fairness is not taken into consideration in L2RW which only aims to improve total performance on long-tail data. Thus, L2RW utilized traditional classification loss as meta loss, while our MBN proposes a novel meta skewness loss to evaluate the bias of model during training such that model is required to learn fair representations on meta data.

3 DATABASES FOR UNBIASED TRAINING AND FAIRNESS EVALUATION

In deep face recognition, few studies focus on bias or fairness problem with respect to skin tone because so few balanced training and testing datasets are available. The deep networks are often designed to fit biased training data, and thus would naturally replicate the biases already existent in data; and test sets contained the same biases as the training set might fail to unveil the unfairness problem of trained models. To address this issue, we contribute two training and one testing databases, i.e., BUPT-Globalface, BUPT-Balancedface and IDS, to enable researchers to go deep into this issue.

3.1 Rationale for skin tone labeling

Although race or ethnicity based categories have been used in some fields [61], [62], they seem unstable in computer vision. First, the drawing boundaries between distinct races are complex and even confused. Thus, inconsistent definition and use have been chief problems with the race concept, particularly as the prevalence of admixture increases across populations. Second, the concept of race is a social construct without biological and physiological basis [29], [63]. Subjects’ phenotypic features, e.g., skin color, eye and nose shapes, can vary widely within a racial or ethnic category. As such, explicitly annotating and building models on top of race or ethnicity risks to perpetuate bias and potentially dangerous outcomes. To scientifically study this demographic bias, we follow the works of Buolamwini et al. [29] and Merle et al. [23], and use skin tone as a more visually precise label.

3.2 Data collection and annotation

We download face images of diverse regions around world from Google according to one-million FreeBase celebrity list [64], and clean them both automatically and manually in the similar way as other FR training datasets, e.g., VGGface2 [19] and Megaface [65]. Then, the skin tones of these downloaded images are estimated. We utilized Individual Typology Angle (ITA) to automatically measure skin tones of images following the work

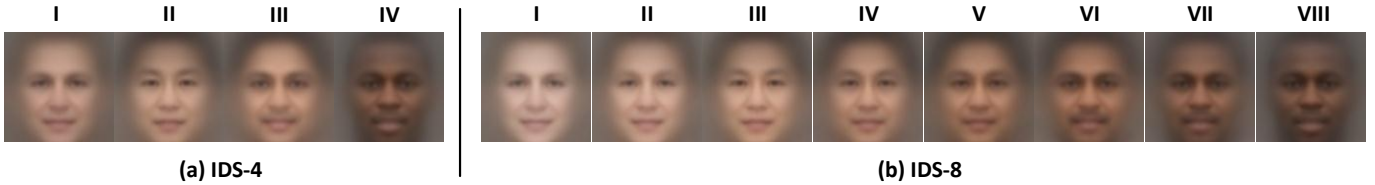


Fig. 1. The average faces of different skin tone bins of IDS database which consist of the average pixel values computed from aligned faces.

of Merle et al. [23]. Note that the skin tone annotations of testing images are further manually checked using Fitzpatrick Skin Type classification system [22], [29]. Thus, images can be divided into several bins based on a mapping to skin tone. We group people into 8 bins, named “Tone I-VIII”. Skin gradually darkens with the increase of tone value, that is, people with Tone-I have the lightest skin color and those with Tone-VIII have the darkest skin tone. After that, we construct our training and testing datasets, i.e., Globalface-8, Balancedface-8 and IDS-8, using these annotated images. Our recently published RFW testing dataset [49] is constructed by the similar way and is divided into 4 bins. For consistency, we call it IDS-4 throughout the paper.

3.3 BUPT-Globalface and BUPT-Balancedface

To remove bias from data aspect and represent people of different skin tones equally, we construct two training datasets, i.e., BUPT-Globalface and BUPT-Balancedface. Globalface contains 2M images from 38K celebrities in total and its distribution is approximately the same as real distribution of world’s population. Balancedface dataset contains 1.3M images from 28K celebrities and is approximately balanced with respect to skin tone. Specifically, Globalface-4 and Balancedface-4 are divided into 4 skin bins, named “Tone I-IV”, based on a mapping to skin tone; while the images in Globalface-8 and Balancedface-8 are divided into 8 bins, i.e., “Tone I-VIII”, as shown in Fig. 2.

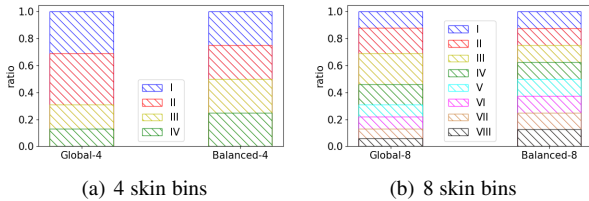


Fig. 2. The skin tone distributions of BUPT-Globalface and BUPT-Balancedface datasets.

3.4 Identity shades dataset: IDS

IDS-4 database is a testing set which contains four skin bins. Each bin contains about 10K images of 3K individuals for face verification. For easy comparison, 6K difficult pairs of images are selected for each bin. Similarly, IDS-8 consists of 8 skin bins and contains 3K difficult pairs of images per bin. All of these images have been carefully and manually cleaned. For the performance evaluation, we recommend to use both the biometric receiver operating characteristic (ROC) curve and LFW-like protocol. Specifically, ROC curve, which aims to report a comprehensive performance, evaluates algorithms on all pairs of identities (about 14K positive vs. 50M negative pairs for IDS-4). In contrast,

LFW-like protocol facilitates easy and fast comparison between algorithms with selected 6K (or 3K) pairs of images for IDS-4 (or IDS-8). Further, inspired by the ugly subset of GBU database [66], we have selected the “difficult” pairs (in term of cosine similarity) to avoid the saturated performance to be easily reported.

Some image examples of our IDS are shown in Fig. 4. In IDS, the images of each skin bin are randomly collected without any preference, there is no significant difference between different bins besides skin color, and thus they are suitable to fairly measure the bias with respect to skin tone. We have validated that, across varying bins, their distributions of pose, age, and gender are similar. As evidence, the detailed distributions measured by Face++ API are shown in Fig. 3.

4 META BALANCED NETWORK

In our meta balanced network, we introduce the idea of adaptive margin into fairness problem. The optimal margins are automatically learned by meta learning to trade-off the specific requirements of people with different skin tones. Compared with traditional sample re-weighting method, our MBN can balance the feature scatter of different skin tone groups in feature space instead of just imposing different weights on loss. In face recognition which is a fine-grained and open-set classification problem, MBN can obtain better generalization ability and more balanced performance across various skin tone subjects.

4.1 Adaptive margin loss

Although large margin losses, e.g., Cosface [27] and Arcface [18], successfully improve feature discrimination, and get better performance on FR benchmarks, they still fail to obtain balanced representations on different skin tone bins, as shown in Table 1. The decision boundary in these losses is assigned the same margin without considering the requirements of people leading to biased performance. In order to address this bias problem and prevent unintended discrimination in existing face recognition algorithms, we suggest that algorithms must trade-off the specific requirements of margins of various groups of people to control their fairness extents under different skin-tone distributions. Therefore, we introduce the idea of adaptive margin into fairness problem with respect to skin tone, and the details are as follows. Note that we take 4 skin bins as examples to describe our method and 8 skin bins can be processed by the similar way.

Since lighter-skinned (Tone-I) subjects are overwhelmingly dominant in numbers and perform best in existing FR datasets, we make lighter-skinned group as the benchmark (anchor) in our paper by remaining the margins of Tone-I subjects unchanged. Optimal margins are learned adaptively for other darker-skinned groups (II-IV) in order to minimize the skewness between lighter- and darker-skinned subjects. In our MBN, we replace the fixed margin in Arcface [18] by a skin tone related and training step

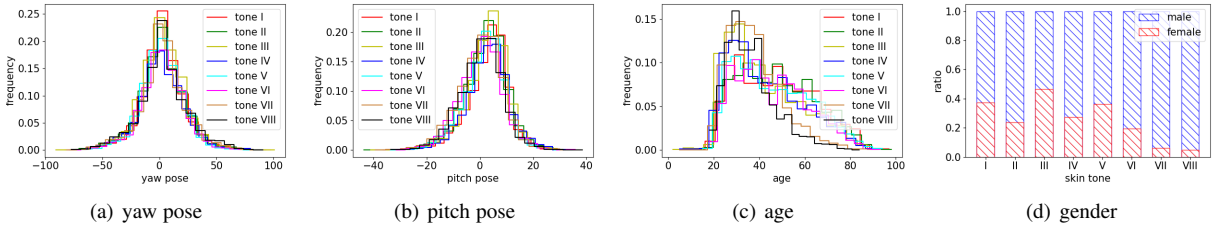


Fig. 3. IDS statistics. We show yaw pose, pitch pose, age and gender distribution of eight testing bins of IDS-8.



Fig. 4. Some images examples of our IDS. One can see from the figure that faces in our dataset are from diversity regions around world and exhibit variability in factors such as pose, age and expression.

related parameter $m_g(t)$, where $m_g \in \{m_{II}, m_{III}, m_{IV}\}$ is the margin corresponding to darker-skinned group g and t represents the stage of the training. The proposed adaptive margin loss function can be formulated as follows:

$$L_j^{T(arc)} = -\log \frac{e^{s(\cos(\theta_{y_j} + \lambda_{g_j}(t)))}}{e^{s(\cos(\theta_{y_j} + \lambda_{g_j}(t)))} + \sum_{i=1, i \neq y_j}^c e^{s \cdot \cos \theta_i}} \quad (1)$$

where, $\lambda_{g_j}(t) = \begin{cases} m, & \text{if } g_j = \text{Tone I} \\ m_g(t), & \text{otherwise} \end{cases}$

where g_j is the skin-tone label of j -th sample. θ_i is the angle between the weight W_i and the feature z_j . $z_j \in \mathbb{R}^d$ denotes the deep feature of the j -th sample, belonging to the y_j -th class, and $W_i \in \mathbb{R}^d$ denotes the i -th column of the weight $W \in \mathbb{R}^{d \times n}$. c is the number of classes and s is the scale factor. The similar modification can be made for Cosface [27] as follows:

$$L_j^{T(cos)} = -\log \frac{e^{s(\cos(\theta_{y_j}) - \lambda_{g_j}(t))}}{e^{s(\cos(\theta_{y_j}) - \lambda_{g_j}(t))} + \sum_{i=1, i \neq y_j}^c e^{s \cdot \cos \theta_i}} \quad (2)$$

where, $\lambda_{g_j}(t) = \begin{cases} m, & \text{if } g_j = \text{Tone I} \\ m_g(t), & \text{otherwise} \end{cases}$

So the key problem is to set optimal margin $m_g(t)$ for each darker-skinned bin g to minimize the skewness between lighter-skinned and darker-skinned subjects. However, bias is a complex problem caused by many latent factors, including but not limited to the quantity. Instead of manually presetting or tuning them by cross-validation, we provide the following algorithm to adaptively learn these margins, by borrowing the idea of recent meta-learning techniques [53], [59], [60].

4.2 Meta margin learning

Our meta margin learning consists of two parts: 1) margin optimization which is responsible for setting training loss (adaptive margin loss) for model optimization by outputting

appropriate margin m_g ; 2) model optimization which optimizes the network by training loss. We aim to learn the optimal training loss such that the network guided by this can perform fairly across different skin tone subjects. To accomplish this, first, we design a skewness loss to evaluate the fairness of learned model by which the margin parameters can be optimized. However, when the training set is biased, the skewness loss imposed on training data would have the wrong perception of bias. According to [59], without a proper definition of an unbiased set, solving the training set bias problem is inherently ill-defined. Therefore, we propose to utilize a clean and balanced meta set to evaluate bias of the learned model and use a large but biased training set to optimize network. Third, to connect disjoint process between training and meta-evaluating, we utilize backward-on-backward automatic differentiation in meta-optimization such that we can transmit meta gradient from meta skewness loss to training loss to update the margins. Benefiting from this bilevel-optimization, we realize a mutual amelioration between model and margin optimization. Then, we introduce the details of our meta learning algorithm.

Model optimization. Let $\mathcal{X} \in \mathbb{R}^d$ be the image space, $\mathcal{Y} = \{1, 2, \dots, c\}$ be the identity label space, and $\mathcal{G} = \{1, 2, \dots, k\}$ be the skin-tone label space. $D_{train} = \{x_i, y_i, g_i\}_{i=1}^N$ denotes a large biased training set, where N is the number of training samples. D_{train} is divided into 4 bins based on a mapping to skin tone, i.e., $D_{train}^g = \{x_i^g, y_i^g\}_{i=1}^{N_g}$, where $g \in \{I, II, III, IV\}$. Let $f(x; w)$ be our neural network model, and w be the model parameters. The training of model $f(x; w)$ is an optimization process that discovers a good model parameter w^* by minimizing our adaptive margin loss on the training data. Our adaptive margin loss can be formulated as $\frac{1}{N} \sum_{j=1}^N L_j^{T(arc)}(w; m_{g_j})$ or $\frac{1}{N} \sum_{j=1}^N L_j^{T(cos)}(w; m_{g_j})$. For notation convenience, we denote that $L_j^T(w; m_{g_j}) = L_j^T(f(x_j; w); y_j; m_{g_j})$. So the model optimization can be formulated as follows:

$$w^*(m_g) = \arg \min_w \frac{1}{N} \sum_{j=1}^N L_j^T(w; m_{g_j}) \quad (3)$$

where $m_g = \{m_{II}, m_{III}, m_{IV}\}$ denotes the margins of darker-skinned groups and g_j is the skin tone label of j -th sample. Note that m_g can be understood as training hyper-parameters and learned automatically during training.

Margin optimization. To cater for different requirements of different people and different status of model training, we should learn different margins for different skin tone bins at each training step t to balance the performance. Assume that we additionally have a small meta-data set $D_{meta} = \{x_i^v, y_i^v, g_i^v\}_{i=1}^M$ with clean labels and balanced data distribution, where M is the number of meta-samples, and $M \ll N$. D_{meta} is also divided into 4 bins,

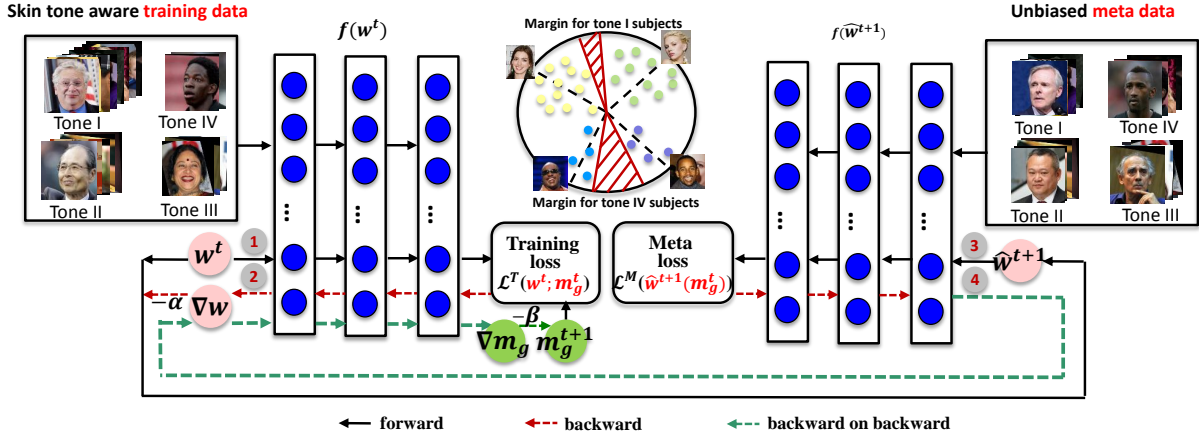


Fig. 5. An illustration of our method. At iteration step t , **Margin parameters learning**: first, given a mini-batch of training samples, we use SGD to update the model parameters \hat{w}^{t+1} on the basis of the current parameters w^t and margins m_g^t (step 1 and 2). Then, a mini-batch of meta samples is sent to the updated model (step 3). Guided by our meta skewness loss, a meta-gradient (high-order gradient) is back-propagated from meta data to training data, respectively, to update the margins m_g^{t+1} (step 4). **Model parameters learning**: with a fixed margin for Tone-I subjects and the updated margins m_g^{t+1} for other faces, we use adaptive margin loss and training data to update the model parameters w^{t+1} such that it performs fairly across people with different skin tones.

i.e., $D_{meta}^g = \{x_i^{v,g}, y_i^{v,g}\}_{i=1}^{M/4}$, where $g \in \{I, II, III, IV\}$. The optimal selection of m_g is based on the performance on meta set.

We aim to learn the optimal margins m_g^* based on this basic assumption: *when testing the model performance on meta set, the model parameters $w^*(m_g^*)$ trained by adaptive margin loss $L^T(w; m_g^*)$ should reduce the bias (skewness) between lighter-skinned (I) and darker-skinned (II-IV) subjects in meta set.* Thus, we can then formulate a meta skewness loss minimization problem with respect to m_g as:

$$m_g^* = \arg \min_{m_g} L^M(w^*(m_g)) \quad (4)$$

where $m_g = \{m_{II}, m_{III}, m_{IV}\}$, and L^M represents the meta skewness loss imposed on meta data and can reflect the skewness between lighter- and darker-skinned subjects in meta set. It consists of three parts: B^{II} , B^{III} and B^{IV} :

$$L^M = B^{II} + B^{III} + B^{IV} \quad (5)$$

B^{II} represents the bias (skewness) of feature scatter between Tone-I and Tone-II subjects, which is formulated as follows. B^{III} and B^{IV} can be computed by the same way as B^{II} .

$$B^{II} = \left| \frac{1}{M/4} \sum_{i=1}^{M/4} l_i^{II} - \frac{1}{M/4} \sum_{k=1}^{M/4} l_k^I \right| \quad (6)$$

where l_i^{II} can reflect intra-class compactness and inter-class discrepancy of meta samples with skin tone II, and we use it to measure the model generalization on Tone-II subjects. l_k^I can be computed by the same way as l_i^{II} to measure the model generalization on subjects with skin tone I.

$$l_i^{II} = \exp \left(\left\| f(x_i^{v,II}) - f(x_{i,p}^{v,II}) \right\|_2^2 - \gamma \left\| f(x_i^{v,II}) - f(x_{i,n}^{v,II}) \right\|_2^2 \right) \quad (7)$$

where $f(x_i^{v,II})$ represents the embedding of meta sample with skin tone II. For notation convenience, we denote that $f(x_i^{v,II}) = f(x_i^{v,II}; w)$. The subscripts i,p and i,n denote meta sample $x_i^{v,II}$'s hardest positive and negative samples,

$\left\| f(x_i^{v,II}) - f(x_{i,p}^{v,II}) \right\|_2^2$ is the L2-norm distance between sample $x_i^{v,II}$ and its positive sample $x_{i,p}^{v,II}$ to measure their similarity. γ is the parameter for the trade-off between positive-pair distance and negative-pair distance. Therefore, B^{II} represents the skewness of model performance between Tone-I and Tone-II subjects. Through minimizing L^M , we can learn the optimal margin m_g^* in large margin loss $L^T(w; m_g^*)$ which makes the learned model mitigate bias between lighter-skinned (I) and darker-skinned (II-IV) groups on meta data.

4.3 Iterative meta-learning training strategy

Calculating the optimal w^* and m_g^* requires two nested loops of optimization, which is expensive to obtain the exact solution [67]. Here we adopt an online approximation strategy [53], [59] to jointly update both w^* and m_g^* in an iterative manner to guarantee the efficiency of the algorithm.

For most training of deep networks, SGD or its variants are used to optimize such loss functions. At every step t of training, a mini-batch of training samples $\{(x_j, y_j), 1 \leq j \leq n\}$ is sampled, where n is the mini-batch size and $n \ll N$. Then, the optimization of the model parameters can be formulated by moving the current w^t along the descent direction of training loss on a mini-batch training data:

$$\hat{w}^{t+1}(m_g^t) = w^t - \alpha \frac{1}{n} \sum_{j=1}^n \nabla_w L_j^T(w; m_{g_j}^t) \Big|_{w=w^t} \quad (8)$$

where α is the descent step size on w .

Then, a mini-batch of meta samples $\{(x_i^v, y_i^v), 1 \leq i \leq n\}$ is sampled with $n/4$ samples per skin tone bin. We extract the features of these meta samples by the updated model parameters $\hat{w}^{t+1}(m_g^t)$, and calculate meta skewness loss L^M . The margin m_g can then be readily updated guided by Eq. 4, i.e., moving the current parameter m_g^t along the objective gradient of L^M calculated on the meta-data:

$$m_g^{t+1} = m_g^t - \beta \nabla_{m_g} L^M(\hat{w}^{t+1}(m_g)) \Big|_{m_g=m_g^t} \quad (9)$$

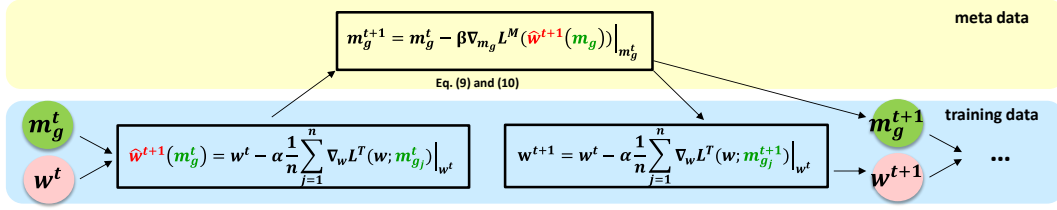


Fig. 6. Main flowchart of the proposed algorithm. Notice that \hat{w}^{t+1} here is a variable instead of a quantity, which makes $\hat{w}^{t+1}(m_g^t)$ a function of m_g^t and the gradient in Eq. (9) and (10) be able to be computed by backward-on-backward differentiation.

where β is the descent step size on margin m_g . Note that we need to use meta loss to optimize the margins of training loss by computing the gradient of L^M with respect to margin m_g . However, the skewness (bias) evaluated by meta skewness loss L^M is incurred on meta data, while the margin m_g only plays effect in adaptive margin loss L^T in training phase. We connect the disjoint process between training and meta-evaluating via the updated model parameter $\hat{w}^{t+1}(m_g^t)$. Notice that \hat{w}^{t+1} here is a variable instead of a quantity, which makes $\hat{w}^{t+1}(m_g^t)$ a function of m_g^t and the gradient in Eq. 9 be able to be computed. Thus, the meta-optimization (Eq. 9) can be performed over the margin parameter m_g as follows:

$$\begin{aligned}
 m_g^{t+1} &= m_g^t - \beta \nabla_{m_g} L^M(\hat{w}^{t+1}(m_g)) \Big|_{m_g=m_g^t} \\
 &= m_g^t - \beta \frac{\partial L^M(\hat{w}^{t+1}(m_g))}{\partial m_g} \Big|_{m_g=m_g^t} \\
 &= m_g^t - \beta \frac{\partial L^M(\hat{w})}{\partial \hat{w}} \Big|_{\hat{w}=\hat{w}^{t+1}} \frac{\partial \hat{w}^{t+1}(m_g)}{\partial m_g} \Big|_{m_g=m_g^t} \\
 &= m_g^t + \frac{\alpha\beta}{n} \sum_{j=1}^n \frac{\partial L^M(\hat{w})}{\partial \hat{w}} \Big|_{\hat{w}=\hat{w}^{t+1}} \frac{\partial^2 L_j^T(w; m_{g_j})}{\partial w \partial m_g} \Big|_{w=w^t, m_g=m_g^t}
 \end{aligned} \tag{10}$$

Therefore, the meta-gradient update involves a gradient through a gradient, i.e., $\frac{\partial^2 L_j^T(w; m_{g_j})}{\partial w \partial m_g}$, which can transmit meta gradient from meta loss to training loss to update the margins. When implementing, we can leverage automatic differentiation techniques to compute it. We can first unroll the gradient graph of the training batch, and then use backward-on-backward automatic differentiation to take a second order gradient pass.

Then, the adaptive margin loss with updated margin m_g^{t+1} is employed to ameliorate the model parameters w :

$$w^{t+1} = w^t - \alpha \frac{1}{n} \sum_{j=1}^n \nabla_w L_j^T(w; m_{g_j}^{t+1}) \Big|_{w=w^t} \tag{11}$$

The meta-learning algorithm can then be summarized in Algorithm 1 and Fig. 6, and Fig. 5 illustrates its main implementation process. In our algorithm, both the model and margin gradually ameliorate their parameters during the learning process based on their values calculated in the last step. Finally, the optimal model parameters are learned which perform well and fairly across different skin tone bins.

5 EXPERIMENTS

5.1 Experimental study on bias

Experimental Settings. We use the similar ResNet-34 architecture described in [18]. It is trained with the guidance

Algorithm 1 Meta balanced network.

Input:

Training data D_{train} , meta data D_{meta} , batch size n and max iterations T .

Output:

Model parameter w and margin parameter m_g .

- 1: Initialize model parameter w^0 and margin parameter m_g^0 .
 - 2: **for** $t = 0$ **to** $T - 1$ **do**
 - 3: $\{X, Y\} \leftarrow \text{SampleMiniBatch}(D_{train}, n)$
 - 4: $\{X^v, Y^v\} \leftarrow \text{SampleMiniBatch}(D_{meta}, n)$
 - 5: $f(X, w^t) \leftarrow \text{Forward}(X, w^t)$.
 - 6: $L^T \leftarrow \frac{1}{n} \sum_{j=1}^n L_j^T(f(x_j; w^t); y_j; m_{g_j}^t)$ by Eq. 1 or 2.
 - 7: $\nabla w^t \leftarrow \text{BackwardAD}(L^T, w^t)$.
 - 8: $\hat{w}^{t+1} \leftarrow w^t - \alpha \nabla w^t$ by Eq. 8.
 - 9: $f(X^v, \hat{w}^{t+1}) \leftarrow \text{Forward}(X^v, \hat{w}^{t+1})$.
 - 10: $L^M \leftarrow L^M(f(X^v; \hat{w}^{t+1}); Y^v)$ by Eq. 5, 6 and 7.
 - 11: $\nabla m_g^t \leftarrow \text{BackwardAD}(L^M, m_g^t)$.
 - 12: $m_g^{t+1} \leftarrow m_g^t - \beta \nabla m_g^t$ by Eq. 10.
 - 13: $\hat{L}^T \leftarrow \frac{1}{n} \sum_{j=1}^n L_j^T(f(x_j; w^t); y_j; m_{g_j}^{t+1})$ by Eq. 1 or 2.
 - 14: $\nabla w^t \leftarrow \text{BackwardAD}(\hat{L}^T, w^t)$.
 - 15: $w^{t+1} \leftarrow w^t - \alpha \nabla w^t$ by Eq. 11.
 - 16: **end for**
-

of Arcface loss [18] on the CASIA-Webface [25], and is called Arcface(CASIA) model. CASIA-Webface consists of 0.5M images of 10K celebrities in which 85% of the photos are lighter-skinned subjects. For preprocessing, we use five facial landmarks for similarity transformation, then crop and resize the faces to 112×112 . Each pixel ($[0, 255]$) in RGB images is normalized by subtracting 127.5 and then being divided by 128. We set the batch size, momentum, and weight decay as 200, 0.9 and $5e - 4$, respectively. LFW [9], CFP-FP [68] and AgeDB-30 [69] are utilized as validations to determine when to decrease the learning rate or stop training. The learning rate is started from 0.1 and decreased twice with a factor of 10 when errors plateau on LFW, CFP-FP and AgeDB-30.

Existence of bias. We examine some SOTA algorithms, i.e., Center-loss [16], Sphereface [17], VGGFace2 [19] and ArcFace [18], as well as four commercial recognition APIs, i.e., Face++, Baidu, Amazon, Microsoft on our IDS-4 and IDS-8, respectively. The results on IDS-4 are presented in Table 2, Fig. 7 and Fig. 10. All SOTA algorithms perform best on Tone-I bin and worst on subjects with skin tone II and IV, which proves the existence of bias with respect to skin tone. The results on IDS-8 are presented in Table 1 and Fig. 8. The similar bias can be observed on IDS-8. For example, Arcface has an accuracy of 92.33% for people with Tone-I but drops to 83.23% for subjects with Tone-VIII.

Existence of domain gap. We extract the features of

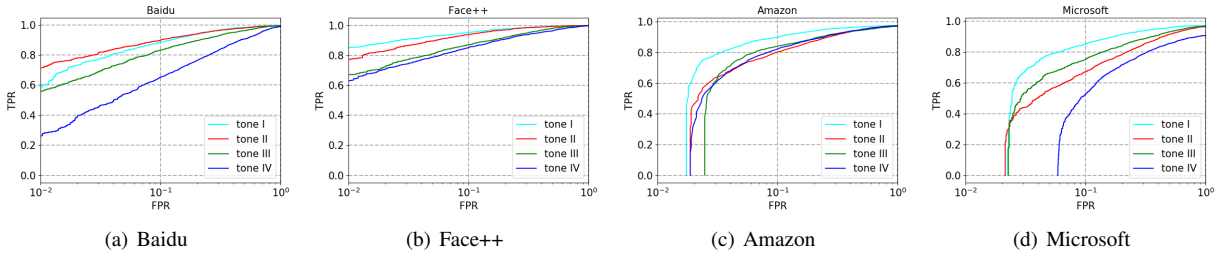


Fig. 7. The ROC curves of (a) Baidu, (b) Face++, (c) Amazon and (d) Microsoft evaluated on 6K pairs of IDS-4. Due to limited number of negative pairs, the performances cannot be reliably estimated at lower FPR values. Besides, once API fails to detect faces, we assume that it will give an incorrect verification result whatever decision thresholds are.

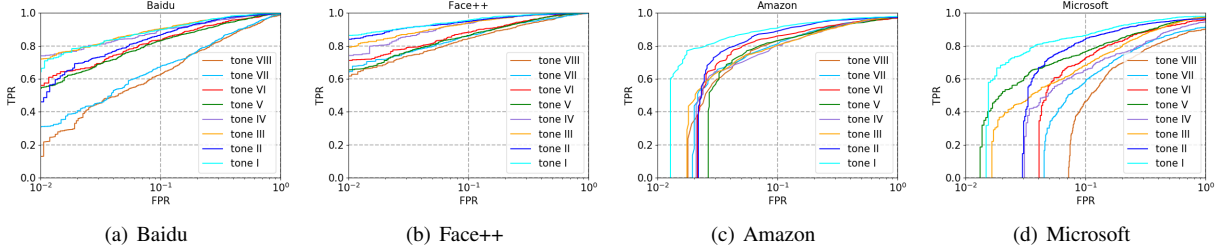


Fig. 8. The ROC curves of (a) Baidu, (b) Face++, (c) Amazon and (d) Microsoft evaluated on 3K pairs of IDS-8. Due to limited number of negative pairs, the performances cannot be reliably estimated at lower FPR values.

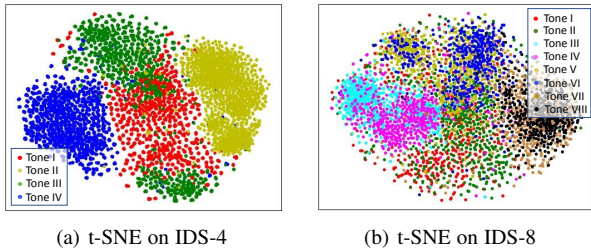


Fig. 9. The feature space of Arcface(CASIA) model.

TABLE 2
Bias with respect to skin tone in commercial APIs and SOTA FR algorithms. Accuracies (%) on our IDS-4 are given.

| Model | LFW [9] | IDS-4 | | | |
|----------------------------|---------|-------|-------|-------|-------|
| | | I | II | III | IV |
| Microsoft [12] | 98.22 | 87.60 | 79.67 | 82.83 | 75.83 |
| Face++ [13] | 97.03 | 93.90 | 92.47 | 88.55 | 87.50 |
| Baidu [14] | 98.67 | 89.13 | 90.27 | 86.53 | 77.97 |
| Amazon [15] | 98.50 | 90.45 | 84.87 | 87.20 | 86.27 |
| mean | 98.11 | 90.27 | 86.82 | 86.28 | 81.89 |
| Center-loss [16] | 98.75 | 87.18 | 79.32 | 81.92 | 78.00 |
| Sphereface [17] | 99.27 | 90.80 | 82.95 | 87.02 | 82.28 |
| Arcface ¹ [18] | 99.40 | 92.15 | 83.98 | 88.00 | 84.93 |
| VGGface2 ² [19] | 99.30 | 89.90 | 84.93 | 86.13 | 83.38 |
| mean | 99.18 | 90.01 | 82.80 | 85.77 | 82.15 |

¹ Arcface is a ResNet-34 model trained with CASIA-Webface.
² VGGFace2 here is a SeNet model trained with VGGFace2.

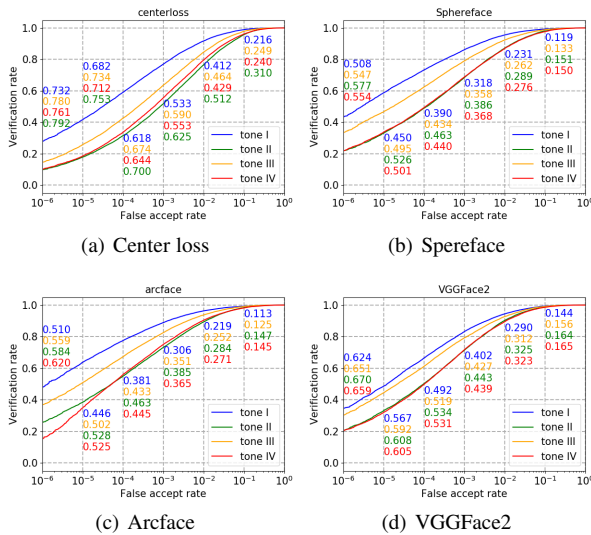


Fig. 10. The ROC curves of (a) Center loss, (b) Sphereface (c) Arcface, (d) VGGFace2 evaluated on all pairs of IDS-4. The cosine similarity thresholds of different skin tone bins are showed at each axis point (FAR={ $10e-6, 10e-5, 10e-4, 10e-3, 10e-2, 10e-1$ }).

1.2K images by our Arcface(CASIA) model and visualize them using t-SNE embeddings [70] in Fig. 9. In Fig. 9(a), although Arcface(CASIA) model is a face recognition model, not skin detection, the skin tone information is highly embedded in the feature space of IDS-4, resulting in four clusters highly correlated with skin tone. On IDS-8, because more fine-grained division according to skin tone increases the similarity of faces of different bins, there is not a clear boundary between these eight skin tone bins. However, they are also separated from each other as shown in Fig. 9(b). From these figures, we make the conclusions: the distribution discrepancies across skin tone bins are large, which conforms that there is domain gap between people with different skin tones.

Cause of bias. We further select a subset from BUPT-Balancedface, called Balancedface*. It contains 590K images from 14K celebrities which has the similar scale with CASIA-Webface database but is approximately balanced with

TABLE 3
Verification accuracy (%) of ResNet-34 models trained with CASIA-Webface [25] and our Balancedface*.

| Training Databases | LFW [9] | CFP-FP [68] | AgeDB-30 [69] | IDS-4 | | | | IDS-8 | | | | | | | |
|----------------------|--------------|--------------|---------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | | | | I | II | III | IV | I | II | III | IV | V | VI | VII | VIII |
| CASIA-WebFace [25] | 99.40 | 93.91 | 93.35 | 92.15 | 83.98 | 88.00 | 84.93 | 92.33 | 91.70 | 83.97 | 83.05 | 87.46 | 88.40 | 84.07 | 83.23 |
| Balancedface* (ours) | 99.55 | 92.74 | 95.15 | 93.92 | 90.60 | 92.98 | 90.98 | 93.95 | 93.85 | 90.14 | 91.14 | 92.97 | 93.55 | 91.70 | 90.32 |

respect to skin tone. Using Balancedface* as training data, we train an Arcface(Balanced) model in the same way as Arcface(CASIA) model and compare their performances on IDS-4 and IDS-8, as shown in Table 3. Compared with Arcface(CASIA) model, Arcface(Balanced) model trained equally on all skin tones performs much better on darker-skinned subjects which proves that bias in databases will reflect in FR algorithm. However, even with balanced training, we see that darker-skinned subjects still perform poorly than subjects with skin tone I. The reason may be that faces of dark skin are inherently difficult to recognize.

Decision thresholds of different skin tone groups. FR algorithms use a threshold similarity score to determine whether two images are of the same subject. Different thresholds yield a different number of true/false positives and true/false negatives, and consequently different accuracy for a given dataset. Here, we look at the influences of the choice of threshold cutoff on performances of different skin tone bins. In Fig. 11, we plot the false positive rate (FPR) as a function of threshold when evaluating commercial APIs on IDS-4. For ROC curves of SOTA algorithms, thresholds are also showed at each axis point in Fig. 10. In all cases, the thresholds that produce the same FPRs are shifted for different bins showing that threshold is skin-tone-specific, i.e., different for different skin tone groups. Therefore, to operate at a particular FPR, threshold should be carefully selected taking skin tone into consideration.

Correlation of bias with other variations. First, we degrade images of IDS-4 by occlusion, low illumination and Gaussian noise, and observe the influence of these variations on accuracy gap between Tone-I and Tone-IV subjects. We train ResNet-34 with guidance of Arcface loss [18] using Globalface-4 dataset, and show the results using red lines in Fig. 12. We can see that subjects with skin tone I and IV are both found to be sensitive to illumination, occlusion and Gaussian noise. Increasing variation level decreases the accuracy and fairness dramatically. Second, the correlations of bias with pose, age and gender are also studied by cohort analysis within each skin tone bin. For pose, we partition images of each skin tone bin into cohorts of (1) small-pose (0°-20°), (2) middle-pose (20°-45°) and (3) large-pose (45°-90°). For age, we partition the images into three cohorts: (1) young (0 to 25 years old), (2) middle-age (25 to 50 years old), and (3) old (50 to 100 years old). For gender, we partition the images into cohorts of (1) male and (2) female. We evaluate the trained Arcface model on different pose and age cohorts of Tone-I and Tone-IV bins. Because females with skin tone IV are too few to be evaluated, we evaluate the model on gender cohorts of our Tone-I and Tone-II bins instead. As we can see from red lines in Fig. 12(g)-(l), for Arcface model, fairness with respect to skin tone decreases on large pose, young and female cohorts. Therefore, we conclude that faces with darker skin are more susceptible to different variations than the lighter-skinned ones, which results in a higher bias across groups when more difficult faces with large variations are recognized.

5.2 Meta margin learning experiment

In this section, we verify the effectiveness of our MBN.

Datasets. We use our Globalface and Balancedface datasets to train our models, and use IDS to fairly measure performances of different skin tone groups. Moreover, in order to adopt meta learning to learn adaptive margins, we additionally collect a meta set with clean labels and balanced distribution. It contains 500 identities per skin tone bin and has been carefully and manually cleaned. We further remove its overlapping subjects with our training and testing datasets.

Experimental Settings. For preprocessing, we share the uniform alignment methods and backbone as Arcface(CASIA) model as mentioned above. The batch size is set to be 240. LFW [9], CFP-FP [68] and AgeDB-30 [69] are utilized as validations to determine when to decrease the learning rate or stop training. For model optimization, the learning rate α starts from 0.1 and is divided by 10 at 80K, 120K, 155K iterations on Globalface. The training process is finished at 180K iterations when errors reach a plateau on LFW, CFP-FP and AgeDB-30. On Balancedface, we divide the learning rate at 70K, 100K, 120K iterations and finish at 145K iterations.

We train the models with SGD and set momentum as 0.9 and weight decay as $5e - 4$. For meta-optimization, we use SGD to optimize the margin parameters with the momentum of 0.9. The learning rate β begins with $1e - 3$ and is decreased twice with a factor of 10. The hyper-parameter γ in meta skewness loss L^M is 0.5. When computing L_M , the triplets (anchor, positive and negative samples) are selected online during the training process for efficiency. Among all the triplets in the generated batches, the online selection chooses those for which: $\left\|f(x_i^v) - f(x_{i,n}^v)\right\|_2^2 - \left\|f(x_i^v) - f(x_{i,p}^v)\right\|_2^2 \leq 0.2$.

When training models on datasets with 4 skin bins, we utilize Tone-I subjects as anchors, and keep their margins unchanged. Optimal margins are learned for other skin tone groups to mitigate their performance biases compared with anchors. According to the loss function and margin of anchors, three variations of our MBN are defined. In MBN(soft), Norm-Softmax [71] is adopted by Tone-I subjects; while MBN(arc) utilizes Arcface loss [18] to optimize anchors and their margins are set to be 0.3. In these two above cases, other skin tone groups are optimized by Arcface loss and their margins are initialized by 0.3. In MBN(cos), Cosface loss is used. Margins are set to be 0.15 for people with skin tone I and are initialized by 0.15 for other groups. When training models on datasets with 8 skin bins, we utilize faces with skin tone I and II as anchors and similar experimental settings are taken.

For evaluating, we test deep models on our IDS dataset and report their accuracies. Following the metric proposed in [28], we utilize average accuracy of different skin tone bins to evaluate total recognition performance, and use the standard deviation (STD) and the skewed error ratio (SER) to evaluate the fairness

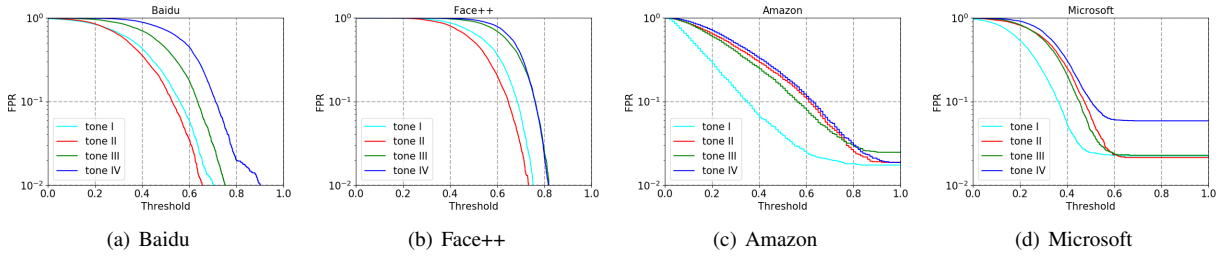


Fig. 11. False positive rate is plotted as a function of threshold when (a) Baidu, (b) Face++, (c) Amazon and (d) Microsoft are evaluated on 6K pairs of IDS-4. Due to limited number of negative pairs, the performances cannot be reliably estimated at lower FPR values.

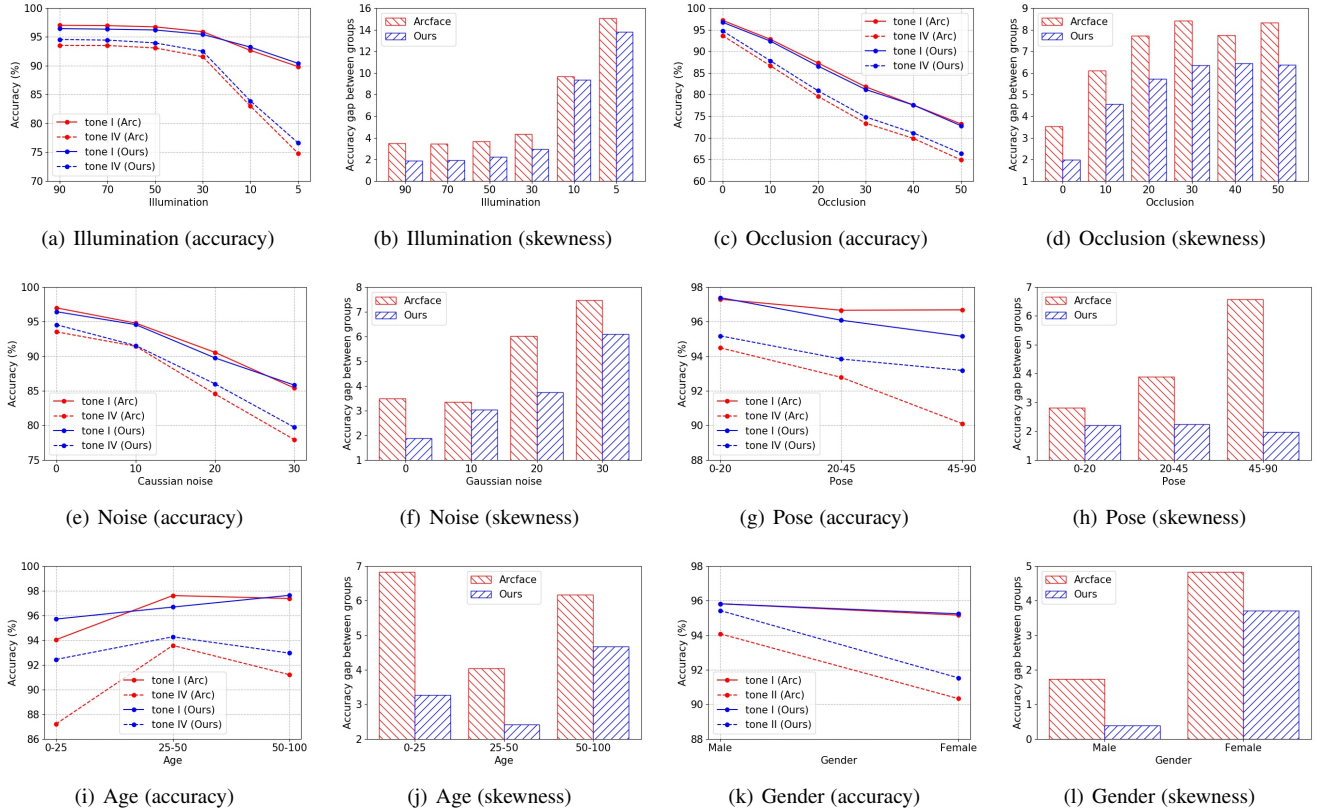


Fig. 12. The influence of different variations on accuracy and skewness (accuracy gap between different skin tone groups) of Arcface and our MBN model. Lower accuracy gap is better and means fairer performance across groups.

performance. STD reflects the amount of dispersion of accuracies of different bins. SER is calculated as the highest error rate divided by the lowest one across skin tone bins, which is formulated as
$$SER = \frac{\max Error_g}{\min Error_g}$$
 and g means skin tone bin.

Results on simulated dataset. We follow the same setting of [28] and evaluate the effectiveness of MBN on some simulated training sets with different skin tone distributions. The images of each training set are selected from BUPT-Globalface-4 and have the similar scale with CASIA-Webface [25]. The ratio of Tone-I subjects varies from $\frac{2}{5}$ to $\frac{7}{10}$. Norm-Softmax [71], which normalizes weights and features based on Softmax, is compared with our MBN.

The results of IDS-4 reported in Table 4 exhibit large gaps between the performances of different skin tone groups, suggesting that training on biased datasets results in algorithmic bias. Second, the skin tone distribution of the training set generally

has a clear impact on the performances of different groups. With the change of distribution (from 4:2:2:2 to 7:1:1:1), we observe a decrease in fairness between different bins, indicative of more uneven distribution with greater bias. Third, our MBN(soft) significantly improves performance on darker-skinned (II-IV) subjects and obtains more balanced performance than Norm-Softmax on different skin tone groups. It verifies the effectiveness of the idea of adaptive margin and the superiority of MBN(soft) under different skin tone distributions.

Results on BUPT-Globalface dataset. We utilize our BUPT-Globalface dataset to validate the effectiveness of our method on large-scale training data. We compare our MBN with Softmax, Cosface [27] and Arcface [18]. Cosface [27] and Arcface [18] assign the same margin for all images and report SOTA performance on the LFW [9] and MegaFace [65]. The scaling parameter is set as 60 for Cosface [27] and Arcface [18]; the margin parameters are set as 0.2 and 0.3, respectively.

TABLE 4

Results on the verification experiments by varying distribution in the training set. Fairness is measured by the standard deviation (STD) (lower is better) and the skewed error ratio (SER) (1 is the best).

| Train ratio ↓ | Method ↓ | Test → | | | | Avg | Fairness | |
|---|----------------|--------|-------|-------|-------|-------|-------------|-------------|
| | | I | II | III | IV | | STD | SER |
| 4 : 2 : 2 : 2 | N-Softmax [71] | 89.67 | 84.68 | 87.97 | 84.17 | 86.62 | 2.64 | 1.53 |
| | MBN(soft) | 91.05 | 88.10 | 90.17 | 88.80 | 89.53 | 1.33 | 1.32 |
| 5 : $\frac{5}{3}$: $\frac{5}{3}$: $\frac{5}{3}$ | N-Softmax [71] | 89.88 | 85.13 | 88.52 | 83.42 | 86.74 | 2.98 | 1.64 |
| | MBN(soft) | 91.28 | 87.73 | 90.68 | 88.15 | 89.46 | 1.78 | 1.41 |
| 6 : $\frac{4}{3}$: $\frac{4}{3}$: $\frac{4}{3}$ | N-Softmax [71] | 90.43 | 84.75 | 88.32 | 83.32 | 86.70 | 3.26 | 1.74 |
| | MBN(soft) | 91.28 | 87.57 | 90.35 | 87.22 | 89.10 | 2.02 | 1.47 |
| 7 : 1 : 1 : 1 | N-Softmax [71] | 90.67 | 84.37 | 87.77 | 82.97 | 86.44 | 3.46 | 1.83 |
| | MBN(soft) | 90.85 | 86.82 | 89.20 | 86.08 | 88.24 | 2.19 | 1.52 |

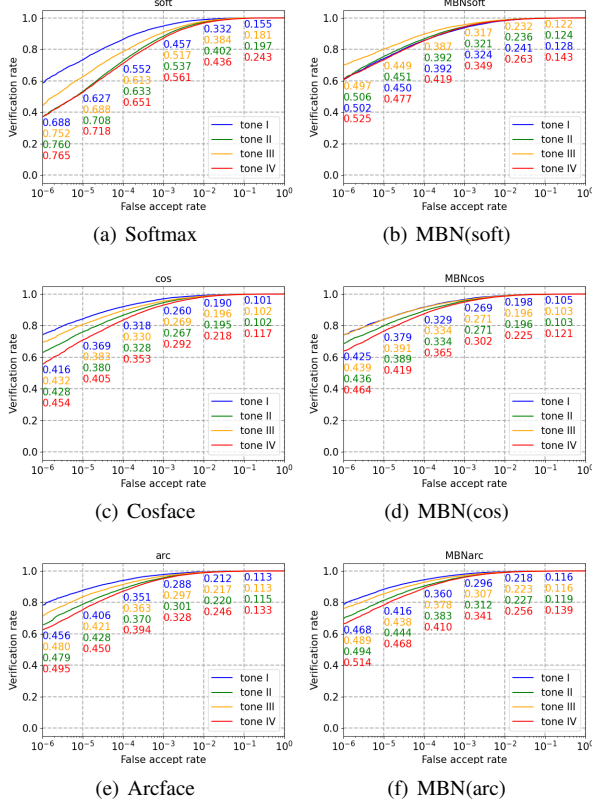


Fig. 13. The ROC curves of (a) Softmax, (b) MBN(soft) (c) Cosface [27], (d) MBN(cos), (e) Arcface [18] and (f) MBN(arc) evaluated on all pairs of IDS-4. The cosine similarity thresholds of different bins are showed at each axis point (FAR={10e−6, 10e−5, 10e−4, 10e−3, 10e−2, 10e−1}).

We train the models on Globalface-4 and show the results tested on IDS-4 in Table 5 and Fig. 13. First, with the help of more separate inter-class data, Cosface [27] and Arcface [18] outperform Softmax by improving the fairness metrics. However, bias cannot be eliminated completely by a uniform margin for different skin tone groups. The performance of darker-skinned subjects is still inferior to that of people with skin tone I. Second, MBN(soft), MBN(cos) and MBN(arc) obtain fairer performances than Softmax, Cosface and Arcface. For example, MBN(arc) is superior to Arcface by reducing the SER from 2.33 to 1.59. It shows the effectiveness of our algorithm on learning balanced features from a biased dataset. Moreover, we additionally train the models on Globalface-8 and show the results tested on IDS-8 in Table 6. The same conclusion can be observed when the images

TABLE 5

Verification accuracy (%) of methods trained with different loss function ([BUPT-Globalface-4, ResNet34, loss*]).

| Methods | IDS-4 | | | | Avg | Fairness | |
|--------------------|-------|-------|-------|-------|--------------|-------------|-------------|
| | I | II | III | IV | | STD | SER |
| Fisher vector [72] | 70.23 | 66.88 | 69.55 | 63.22 | 67.47 | 3.18 | 1.24 |
| Triplet [8] | 95.80 | 91.03 | 92.77 | 90.47 | 92.52 | 2.40 | 2.27 |
| Softmax | 95.62 | 90.85 | 91.97 | 89.98 | 92.10 | 2.48 | 2.29 |
| M-RBN(soft) [28] | 93.50 | 90.06 | 94.50 | 93.43 | 92.87 | 1.94 | 1.81 |
| RL-RBN(soft) [28] | 94.53 | 94.20 | 95.03 | 94.05 | 94.45 | 0.44 | 1.20 |
| MBN(soft) | 94.62 | 94.18 | 94.72 | 93.87 | 94.35 | 0.39 | 1.16 |
| Cosface [27] | 96.63 | 93.50 | 94.68 | 92.17 | 94.25 | 1.90 | 2.33 |
| M-RBN(cos) [28] | 96.15 | 93.43 | 95.73 | 94.76 | 95.02 | 1.21 | 1.70 |
| RL-RBN(cos) [28] | 96.03 | 94.58 | 95.15 | 94.27 | 95.01 | 0.77 | 1.45 |
| MBN(cos) | 96.07 | 94.87 | 95.52 | 94.43 | 95.22 | 0.72 | 1.42 |
| Arcface [18] | 97.37 | 94.55 | 95.68 | 93.87 | 95.37 | 1.53 | 2.33 |
| M-RBN(arc) [28] | 97.03 | 94.40 | 95.58 | 95.18 | 95.55 | 1.10 | 1.89 |
| RL-RBN(arc) [28] | 97.08 | 95.57 | 95.63 | 94.87 | 95.79 | 0.93 | 1.76 |
| MBN(arc) | 96.87 | 95.63 | 96.20 | 95.00 | 95.93 | 0.80 | 1.59 |

are divided into 8 skin bins. Cosface and Arcface still show biased performance on faces with different skin tones, and our method performs more fairly benefitting from adaptive margin. For example, MBN(cos) is superior to Cosface by reducing the SER from 2.28 to 1.55 on IDS-8.

Results on BUPT-Balancedface dataset. We also compare our MBN with Softmax, Cosface [27] and Arcface [18] on BUPT-Balancedface. We train the models on Balancedface-4 and Balancedface-8, and then evaluate them on IDS-4 and IDS-8 respectively, as shown in Table 7 and Table 8. First, with balanced training, Softmax, Cosface and Arcface indeed obtain more balanced performances compared with trained on biased data. So training models on datasets well distributed across all skin tone groups can help to reduce face matcher vulnerabilities on specific cohorts to some extent. For instance, Arcface trained on Balancedface-8 decreases STD from 1.52 to 1.16 compared with trained on Globalface-8. Second, the results obtained by our MBN outperform all compared approaches by improving the fairness metrics in performances across skin tone bins. For example, compared with Arcface [18], MBN(arc) reduces the SER from 1.65 to 1.37 and reduces the standard deviation by 47% on IDS-4. Combining our debiasing algorithm and balanced data, we can obtain the fairest performance.

Adaptive margin mechanism visualization. To better understand our MBN, we plot the learned margins, skewness and accuracies of different skin tone bins during the learning process in Fig. 14 and Fig. 15. First, we can see from Fig. 14(a-b) and Fig. 15(a-b) that our method automatically learns

TABLE 6
Verification accuracy (%) of methods trained with different loss function ([BUPT-Globalface-8, ResNet34, loss*]).

| Methods | The skin tone of IDS-8 | | | | | | | | Avg | Fairness | |
|--------------|------------------------|-------|-------|-------|-------|-------|-------|-------|--------------|-------------|-------------|
| | I | II | III | IV | V | VI | VII | VIII | | STD | SER |
| Softmax | 95.06 | 94.98 | 90.29 | 90.02 | 90.64 | 93.02 | 89.88 | 87.97 | 91.48 | 2.58 | 2.43 |
| MBN(soft) | 94.05 | 94.18 | 93.47 | 93.75 | 93.37 | 94.94 | 92.96 | 92.34 | 93.63 | 0.80 | 1.51 |
| Cosface [27] | 96.30 | 95.97 | 93.19 | 93.08 | 94.16 | 94.56 | 92.86 | 91.57 | 93.96 | 1.61 | 2.28 |
| MBN(cos) | 95.87 | 95.84 | 94.64 | 94.82 | 94.70 | 96.20 | 94.58 | 94.12 | 95.09 | 0.76 | 1.55 |
| Arcface [18] | 97.34 | 97.09 | 94.43 | 94.72 | 95.36 | 95.86 | 94.25 | 92.78 | 95.23 | 1.52 | 2.72 |
| MBN(arc) | 97.11 | 96.83 | 94.94 | 95.45 | 95.54 | 96.73 | 94.81 | 94.83 | 95.78 | 0.96 | 1.79 |

TABLE 7
Verification accuracy (%) of methods trained with different loss function ([BUPT-Balancedface-4, ResNet34, loss*]).

| Methods | IDS-4 | | | | Avg | Fairness | |
|--------------------|-------|-------|-------|-------|--------------|-------------|-------------|
| | I | II | III | IV | | STD | SER |
| Fisher vector [72] | 70.23 | 66.88 | 69.55 | 63.22 | 67.47 | 3.18 | 1.24 |
| Triplet [8] | 94.58 | 91.48 | 93.17 | 91.60 | 92.71 | 1.47 | 1.57 |
| Softmax | 94.18 | 91.23 | 92.82 | 91.42 | 92.37 | 1.42 | 1.51 |
| RL-RBN(soft) [28] | 94.30 | 93.87 | 94.13 | 94.45 | 94.19 | 0.25 | 1.10 |
| MBN(soft) | 93.45 | 93.82 | 93.90 | 93.83 | 93.75 | 0.20 | 1.07 |
| Cosface [27] | 95.12 | 92.98 | 93.93 | 92.93 | 93.74 | 1.03 | 1.45 |
| RL-RBN(cos) [28] | 95.47 | 94.52 | 95.15 | 95.27 | 95.10 | 0.41 | 1.21 |
| MBN(cos) | 95.37 | 94.43 | 95.17 | 95.05 | 94.99 | 0.39 | 1.19 |
| Arcface [18] | 96.18 | 93.72 | 94.67 | 93.98 | 94.64 | 1.11 | 1.65 |
| RL-RBN(arc) [28] | 96.27 | 94.82 | 94.68 | 95.00 | 95.19 | 0.73 | 1.43 |
| MBN(arc) | 96.25 | 94.85 | 95.32 | 95.38 | 95.45 | 0.58 | 1.37 |

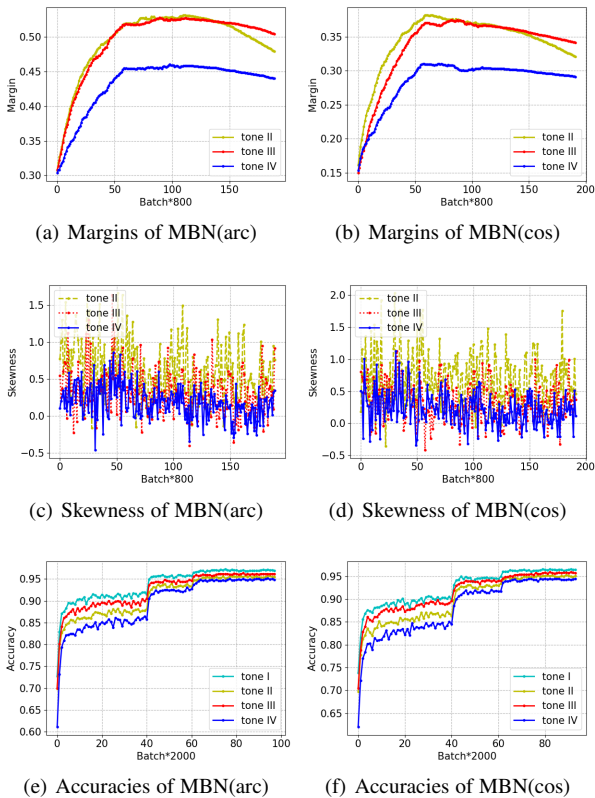


Fig. 14. The margins learned for people with skin tone II-IV in (a) MBN(arc) and (b) MBN(cos) trained on Globalface-4. The skewness between lighter-skinned (I) and darker-skinned (II-IV) subjects in (c) MBN(arc) and (d) MBN(cos). The accuracies of different bins in (e) MBN(arc) and (f) MBN(cos).

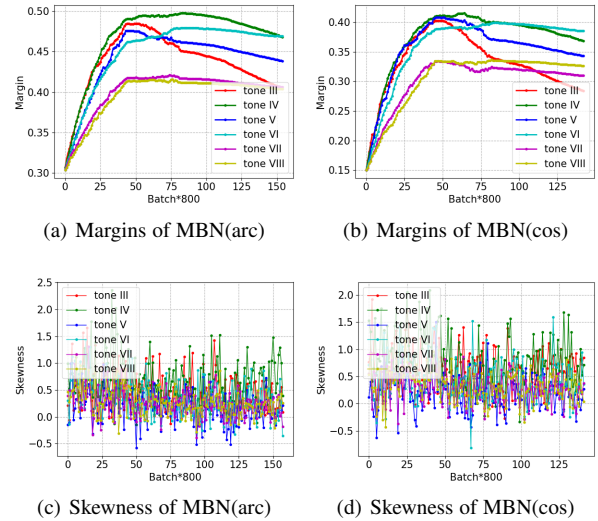


Fig. 15. The margins learned for people with skin tone III-VIII in (a) MBN(arc) and (b) MBN(cos) trained on Globalface-8. The skewness between lighter-skinned (I-II) and darker-skinned (III-VIII) subjects in (c) MBN(arc) and (d) MBN(cos).

larger margins for darker-skinned subjects compared with people with skin tone I. This is consistent with our theoretical analysis that we prefer stricter constraints for dark faces since they are difficult to recognize. Second, the skewness between lighter- and darker-skinned people calculated on meta data by $B^g = \frac{1}{n/4} \sum_{i=1}^{n/4} l_i^g - \frac{1}{n/4} \sum_{k=1}^{n/4} l_k^I$ is shown in Fig. 14(c-d) and Fig. 15(c-d), where $g \in \{II, III, IV\}$ and l^g and l^I can be computed by Eq. 7. On meta data, the skewness is always larger than zero which shows that darker-skinned subjects can not obtain as good intra-class compactness and inter-class discrepancy as Tone-I people do. However, we can see from the results on 4 skin bins shown in Fig. 14(c-d), the skewness of Tone-II group is relatively high on meta data even if people with skin tone IV seem more difficult to recognize in test set. Our hypothesis to explain this phenomenon is that there is a little domain shift between meta and test data which is hard to remove entirely. Despite this little shift, our MBN can still learn optimal margin for each bin, leading to more balanced performance on IDS as shown in Fig. 14(e) and Fig. 14(f).

Feature visualization. Similar to Fig. 9(a), we extract the features by our MBN(arc) and visualize them in Fig. 16(a). We can find that there are still domain gaps between different skin tone bins. It is reasonable because our MBN learns fairer representations through balancing the feature scatter across groups instead of domain adaptation or feature disentanglement. The skin

TABLE 8
Verification accuracy (%) of methods trained with different loss function ([BUPT-Balancedface-8, ResNet34, loss*]).

| Methods | The skin tone of IDS-8 | | | | | | | | Avg | Fairness | |
|--------------|------------------------|-------|-------|-------|-------|-------|-------|-------|--------------|-------------|-------------|
| | I | II | III | IV | V | VI | VII | VIII | | STD | SER |
| Softmax | 94.28 | 93.09 | 90.70 | 90.75 | 91.49 | 93.31 | 91.27 | 91.57 | 92.06 | 1.33 | 1.63 |
| MBN(soft) | 93.61 | 93.23 | 93.21 | 92.78 | 92.89 | 94.55 | 93.82 | 93.11 | 93.40 | 0.58 | 1.33 |
| Cosface [27] | 94.38 | 94.58 | 92.07 | 92.30 | 93.68 | 94.70 | 92.99 | 91.77 | 93.31 | 1.18 | 1.55 |
| MBN(cos) | 94.92 | 95.37 | 94.28 | 94.53 | 94.06 | 95.09 | 94.84 | 93.58 | 94.58 | 0.59 | 1.39 |
| Arcface [18] | 96.03 | 95.90 | 93.57 | 93.36 | 93.48 | 95.23 | 94.28 | 93.28 | 94.39 | 1.16 | 1.69 |
| MBN(arc) | 95.93 | 95.57 | 94.28 | 94.77 | 94.57 | 95.57 | 94.21 | 94.46 | 94.92 | 0.67 | 1.42 |

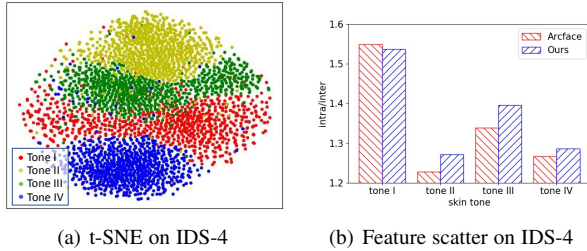


Fig. 16. (a) The feature space and (b) feature scatter of our MBN(arc) model trained on Globalface-4.

tone information is still embedded in the feature space resulting in four separated clusters. Moreover, we also show the feature scatter of Arcface and MBN(arc) in Fig. 16(b). The feature scatter is computed by intra-class scatter divided by inter-class scatter, where intra-class scatter refers to the mean of cosine similarities between features and their corresponding feature centres and inter-class scatter refers to the mean of cosine similarities between embedding feature centres. According to the definition, larger feature scatter means better performance. From figure, we can see that feature scatter of Arcface is biased drastically towards subjects with skin tone I while our MBN improves the feature scatter of darker-skinned subjects leading to balanced performance.

Tolerance to different variations. As we proved in Section 5.1, large variations will result in a higher bias across different skin tone groups. Therefore, we measure the influence of illumination, occlusion, Gaussian noise, pose, age and gender on our debiasing method to validate its robustness. As shown in Fig. 12, compared with Arcface model, our MBN can consistently improve the performances of darker-skinned subjects and present less bias with respect to skin tone under different conditions, even if the variations are extremely large. Moreover, our MBN shows more excellent debiasing ability under pose and age variations, especially on larger-pose and younger cohorts which are more difficult to recognize. Compared with Arcface model, MBN reduces the accuracy gap between Tone-I and Tone-IV subjects from 6.58% to 1.97% on large-pose group (45° - 90°) and reduces the gap from 6.83% to 3.27% on cohorts with age between 0 to 25. Besides, the results in Fig. 12(k) and 12(l) show that our debiasing method also performs fairly on different gender cohorts. The skewness with respect to skin tone is reduced from 1.74% to 0.39% on males, and from 4.83% to 3.70% on females.

Parameter sensitivity. The hyper-parameter γ in Eqn. 7 is a tradeoff parameter that balances positive-pair distance and negative-pair distance. It affects the measurement of fairness in meta skewness loss and so the debiasing performance. We studied this parameter by setting it to different values and checking the debiasing performance. A ResNet-34 is trained using Globalface-4

dataset and the results tested on IDS-4 are shown in Fig. 17. We observe that bias first decreases and then increases as γ varies from 0.125 to 2 and our MBN performs best when $\gamma = 0.5$. This suggests that positive-pair distance contributes more in measuring fairness in meta skewness loss. Actually, we find, through experiments, that intra-class compactness in Arcface model is indeed more biased across skin tone bins compared with inter-class discrepancy, so we should focus more on the positive-pair distance and less on the negative-pair distance in debiasing process.

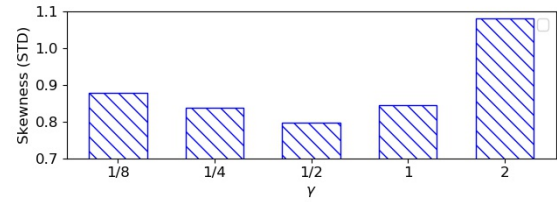


Fig. 17. The effect of hyper-parameter γ on debiasing ability of our method (STDs of accuracies of different bins are showed).

Comparison with other adaptive margin methods. RL-RBN [28] and M-RBN [28] published in CVPR'20 are closely related methods with ours, which adopt a similar adaptive margin mechanism. M-RBN utilizes different fixed margins for different skin tone bins inversely proportional to their quantity; and RL-RBN formulates the process of finding the optimal margins for darker-skinned subjects as a Markov decision process and employs deep reinforcement learning to learn margin policies. We show the compared results in Table 5 and Table. 7. Although M-RBN [28] can improve the fairness of model compared with Cosface [27] and Arcface [18] benefitting from different margins, the performance of people with skin tone II is always a drag on fairness. This is because bias is a complex problem in which the quantity is not only fact affecting out-of-balance accuracy. Although the number of Tone-II people is much larger than that of people with skin tone III and IV in Globalface-4, Tone-II group still needs a larger margin. Moreover, our MBN obtains more uniform performance than M-RBN and RL-RBN which shows the superiority of our algorithm. RL-RBN utilizes deep Q-learning and makes the search space of margin \mathcal{M} to be discrete by assuming $\mathcal{M} = \{m_1, m_2, m_3, m_4\}$. For example, the optimal margin can only be searched from a subspace $\mathcal{M}_s = \{0.3, 0.4, 0.5, 0.6\}$ in RL-RBN(arc). Our MBN uses gradient based optimization and can conduct a continuous search in the whole space, which enables more suitable margin parameters for darker-skinned people.

Comparison with other debiasing methods. To address the class-imbalance and bias issue, sample re-weighting [59], [60]

TABLE 9

Verification accuracy (%) of other debiasing methods trained with Arcface loss ([BUPT-Globalface-4, ResNet34, Arcface]).

| Methods | IDS-4 | | | | Avg | Fairness | |
|------------------|-------|-------|-------|-------|-------|----------|------|
| | I | II | III | IV | | STD | SER |
| Arcface [18] | 97.37 | 94.55 | 95.68 | 93.87 | 95.37 | 1.53 | 2.33 |
| Re-weight [59] | 96.35 | 94.25 | 95.32 | 93.48 | 94.85 | 1.25 | 1.79 |
| Adversarial [44] | 96.63 | 94.17 | 95.27 | 93.70 | 94.94 | 1.30 | 1.87 |
| MBN(arc) | 96.87 | 95.63 | 96.20 | 95.00 | 95.93 | 0.80 | 1.59 |

and adversarial attribute removal methods [44], [73] have been exploited and achieved improved performance in some tasks, e.g., object and gender classification. Here, we also compare our MBN with these methods in Table 9. All comparison methods are trained with Arcface loss [18] on Globalface-4 dataset.

From the results, we can find that our MBN is superior to all compared methods. The reasons may be as follows. First, there is a tradeoff between model identification performance and model fairness in adversarial attribute removal method. Such identity feature contained little skin-tone information could undermine the recognition competence since skin tone information is a part of identity-related facial appearance. Therefore, in order to maintain satisfactory model identification performance, fairness cannot be improved significantly. Second, sample re-weighting method just makes the network pay more attention to darker-skinned people during training by imposing larger weights on their losses, but actually, this improvement of network generalization on darker-skinned people is very limited since face recognition is a fine-grained and open-set classification problem. However, our MBN can adjust the feature scatter in feature space by learning optimal margins for darker-skinned groups such that intra-class compactness and inter-class discrepancy of darker-skinned subjects can be improved significantly leading to balanced performance.

Effectiveness on age and gender bias. For age bias, we can see from Fig. 12(i) that Arcface performs best on 25-50 years old cohort, followed by 50-100 age cohort, and worst on younger cohort, which proves that age bias also exists in face recognition algorithms. Compared with Arcface, our MBN can successfully reduce age bias. For example, on Tone-VI bin of IDS-4, Arcface model obtains the accuracies of 87.22%, 93.58% and 91.21% on young, middle-age and old cohorts, respectively, presenting a standard deviation of 3.21; while our MBN decreases the standard deviation from 3.21 to 0.94 and achieves the accuracies of 92.44%, 94.28% and 92.96% on different age cohorts. For gender bias, Arcface always performs worse on females, especially darker-skinned females, as shown in Fig. 12(k), which proves the existence of gender bias. Compared with Arcface, our MBN has a little effect on mitigating gender bias. For example, on Tone-I bin, Arcface model obtains the accuracies of 95.82% on males and 95.16% on females, presenting an accuracy gap of 0.66; while our MBN achieves the accuracies of 95.82% and 95.24%, decreasing the accuracy gap from 0.66 to 0.58. However, the improvement of fairness with respect to gender is limited. It is reasonable because MBN is designed for bias with respect to skin tone by learning adaptive margins for different skin tone bins without considering gender.

6 CONCLUSION

Considering that the problem of bias with respect to skin tone is yet to be comprehensively studied, we have done the first step for this bias of face recognition and create a benchmark for it. Our IDS database encourages FR algorithms to be fairly evaluated and compared on different skin tone groups. Through experiments on our IDS, we first prove that the deep models trained on the current benchmarks do not perform well on darker-skinned faces and that this bias comes from both data and algorithm aspects. Then, we provide two large-scale training datasets, i.e., BUPT-Globalface and BUPT-Balancedface, to remove bias from data aspect. Finally, a meta balanced network is proposed to alleviate bias and learn more balanced features from algorithm aspect. The comprehensive experiments prove the potential and effectiveness of our MBN on reducing bias with respect to skin tone. However, our MBN should be trained on skin tone aware training datasets and requires access to the sample-level protected attribute, i.e., skin tone label, during training. Hence, one future trend is to investigate more elegant debiasing algorithm such that a balanced model can be learned without the explicit usage of any demographic attributes.

ACKNOWLEDGMENTS

This work was supported in part by the National Natural Science Foundation of China under Grants No. 61871052 and BUPT Excellent Ph.D. Students Foundation CX2020207.

REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, pp. 1097–1105, 2012.
- [2] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [3] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [4] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [5] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [6] M. Wang and W. Deng, "Deep face recognition: A survey," *Neurocomputing*, vol. 429, pp. 215–244, 2021.
- [7] Y. Sun, Y. Chen, X. Wang, and X. Tang, "Deep learning face representation by joint identification-verification," in *Proceedings of the 27th International Conference on Neural Information Processing Systems-Volume 2*, 2014, pp. 1988–1996.
- [8] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 815–823.
- [9] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller, "Labeled faces in the wild: A database for studying face recognition in unconstrained environments," Technical Report 07-49, University of Massachusetts, Amherst, Tech. Rep., 2007.
- [10] C. Garvie, *The perpetual line-up: Unregulated police face recognition in america*. Georgetown Law, Center on Privacy & Technology, 2016.
- [11] "Are face recognition systems accurate? depends on your race." <https://www.technologyreview.com/s/601786>.
- [12] "Microsoft azure," <https://www.azure.cn>.
- [13] "Face++ research toolkit," www.faceplusplus.com.
- [14] "Baidu cloud vision api," <http://ai.baidu.com>.
- [15] "Amazon's rekognition tool," <https://aws.amazon.com/rekognition/>.
- [16] Y. Wen, K. Zhang, Z. Li, and Y. Qiao, "A discriminative feature learning approach for deep face recognition," in *European conference on computer vision*. Springer, 2016, pp. 499–515.

- [17] W. Liu, Y. Wen, Z. Yu, M. Li, B. Raj, and L. Song, "Sphereface: Deep hypersphere embedding for face recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, vol. 1, 2017.
- [18] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "Arcface: Additive angular margin loss for deep face recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4690–4699.
- [19] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, "Vggface2: A dataset for recognising faces across pose and age," in *2018 13th IEEE international conference on automatic face & gesture recognition (FG 2018)*. IEEE, 2018, pp. 67–74.
- [20] P. J. Phillips, F. Jiang, A. Narvekar, J. Ayyad, and A. J. O'Toole, "An other-race effect for face recognition algorithms," *ACM Transactions on Applied Perception (TAP)*, vol. 8, no. 2, p. 14, 2011.
- [21] B. F. Klare, M. J. Burge, J. C. Klontz, R. W. V. Bruegge, and A. K. Jain, "Face recognition performance: Role of demographic information," *IEEE Transactions on Information Forensics and Security*, vol. 7, no. 6, pp. 1789–1801, 2012.
- [22] T. B. Fitzpatrick, "The validity and practicality of sun-reactive skin types i through vi," *Archives of dermatology*, vol. 124, no. 6, pp. 869–871, 1988.
- [23] M. Merler, N. Ratha, R. S. Feris, and J. R. Smith, "Diversity in faces," *arXiv preprint arXiv:1901.10436*, 2019.
- [24] A. Chardon, I. Cretois, and C. Hourseau, "Skin colour typology and suntanning pathways," *International journal of cosmetic science*, vol. 13, no. 4, pp. 191–208, 1991.
- [25] D. Yi, Z. Lei, S. Liao, and S. Z. Li, "Learning face representation from scratch," *arXiv preprint arXiv:1411.7923*, 2014.
- [26] Y. Guo, L. Zhang, Y. Hu, X. He, and J. Gao, "Ms-celeb-1m: A dataset and benchmark for large-scale face recognition," in *European conference on computer vision*. Springer, 2016, pp. 87–102.
- [27] H. Wang, Y. Wang, Z. Zhou, X. Ji, D. Gong, J. Zhou, Z. Li, and W. Liu, "Cosface: Large margin cosine loss for deep face recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5265–5274.
- [28] M. Wang and W. Deng, "Mitigating bias in face recognition using skewness-aware reinforcement learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9322–9331.
- [29] J. Buolamwini and T. Gebru, "Gender shades: Intersectional accuracy disparities in commercial gender classification," in *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, vol. 81, 2018, pp. 77–91.
- [30] I. D. Raji and J. Buolamwini, "Actionable auditing: Investigating the impact of publicly naming biased performance results of commercial ai products," in *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, 2019, pp. 429–435.
- [31] A. Amini, A. P. Soleimany, W. Schwarting, S. N. Bhatia, and D. Rus, "Uncovering and mitigating algorithmic bias through learned latent structure," in *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, 2019, pp. 289–295.
- [32] P. J. Grother, G. W. Quinn, and P. J. Phillips, "Report on the evaluation of 2d still-image face recognition algorithms," *NIST interagency report*, vol. 7709, p. 106, 2010.
- [33] P. J. Phillips, P. Grother, R. Micheals, D. M. Blackburn, E. Tabassi, and M. Bone, "Face recognition vendor test 2002," in *Analysis and Modeling of Faces and Gestures, 2003. AMFG 2003. IEEE International Workshop on*. IEEE, 2003, p. 44.
- [34] P. Grother, M. Ngan, and K. Hanaoka, *Face Recognition Vendor Test (FVRT): Part 3, Demographic Effects*. National Institute of Standards and Technology, 2019.
- [35] K. Vangara, M. C. King, V. Albiero, K. Bowyer *et al.*, "Characterizing the variability in face recognition accuracy relative to race," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.
- [36] K. Krishnapriya, V. Albiero, K. Vangara, M. C. King, and K. W. Bowyer, "Issues related to face recognition accuracy varying based on race and skin tone," *IEEE Transactions on Technology and Society*, vol. 1, no. 1, pp. 8–20, 2020.
- [37] K. Ricanek and T. Tesafaye, "Morph: A longitudinal image database of normal adult age-progression," in *7th International Conference on Automatic Face and Gesture Recognition (FG'06)*. IEEE, 2006, pp. 341–345.
- [38] N. Furl, P. J. Phillips, and A. J. O'Toole, "Face recognition algorithms and the other-race effect: computational mechanisms for a developmental contact hypothesis," *Cognitive Science*, vol. 26, no. 6, pp. 797–815, 2002.
- [39] P. J. Phillips, H. Wechsler, J. Huang, and P. J. Rauss, "The feret database and evaluation procedure for face-recognition algorithms," *Image & Vision Computing*, vol. 16, no. 5, pp. 295–306, 1998.
- [40] J. R. Beveridge, G. H. Givens, P. J. Phillips, B. A. Draper, and Y. M. Lui, "Focus on quality, predicting frvt 2006 performance," in *Automatic Face & Gesture Recognition, 2008. FG'08. 8th IEEE International Conference on*. IEEE, 2008, pp. 1–8.
- [41] P. Sattigeri, S. C. Hoffman, V. Chenthamarakshan, and K. R. Varshney, "Fairness gan," *arXiv preprint arXiv:1805.09910*, 2018.
- [42] F. Calmon, D. Wei, B. Vinzamuri, K. N. Ramamurthy, and K. R. Varshney, "Optimized pre-processing for discrimination prevention," in *Advances in Neural Information Processing Systems*, 2017, pp. 3992–4001.
- [43] S. Gong, X. Liu, and A. K. Jain, "Mitigating face recognition bias via group adaptive classifier," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3414–3424.
- [44] M. Alvi, A. Zisserman, and C. Nellåker, "Turning a blind eye: Explicit removal of biases and variation from deep neural network embeddings," in *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, 2018, pp. 0–0.
- [45] V. Mirjalili, S. Raschka, and A. Ross, "Gender privacy: An ensemble of semi adversarial networks for confounding arbitrary gender classifiers," in *2018 IEEE 9th International Conference on Biometrics Theory, Applications and Systems (BTAS)*. IEEE, 2018, pp. 1–10.
- [46] A. Othman and A. Ross, "Privacy of facial soft biometrics: Suppressing gender but retaining identity," in *European Conference on Computer Vision*. Springer, 2014, pp. 682–696.
- [47] A. Morales, J. Fierrez, R. Vera-Rodriguez, and R. Tolosana, "Sensitiveness: Learning agnostic representations with application to face images," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [48] S. Gong, X. Liu, and A. K. Jain, "Jointly de-biasing face recognition and demographic attribute estimation," in *European Conference on Computer Vision*. Springer, 2020, pp. 330–347.
- [49] M. Wang, W. Deng, J. Hu, X. Tao, and Y. Huang, "Racial faces in the wild: Reducing racial bias by information maximization adaptation network," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 692–702.
- [50] M. Kan, S. Shan, and X. Chen, "Bi-shifting auto-encoder for unsupervised domain adaptation," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 3846–3854.
- [51] M. Kan, J. Wu, S. Shan, and X. Chen, "Domain adaptation for face recognition: Targetize source domain bridged by common subspace," *International Journal of Computer Vision*, vol. 109, no. 1-2, pp. 94–109, 2014.
- [52] Z. Zhang, Y. Song, and H. Qi, "Age progression/regression by conditional adversarial autoencoder," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [53] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 1126–1135.
- [54] A. Nichol, J. Achiam, and J. Schulman, "On first-order meta-learning algorithms," *arXiv preprint arXiv:1803.02999*, 2018.
- [55] Q. Sun, Y. Liu, T.-S. Chua, and B. Schiele, "Meta-transfer learning for few-shot learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 403–412.
- [56] A. Rajeswaran, C. Finn, S. M. Kakade, and S. Levine, "Meta-learning with implicit gradients," in *Advances in Neural Information Processing Systems*, 2019, pp. 113–124.
- [57] S. Ravi and H. Larochelle, "Optimization as a model for few-shot learning," 2016.
- [58] J. Guo, X. Zhu, C. Zhao, D. Cao, Z. Lei, and S. Z. Li, "Learning meta face recognition in unseen domains," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [59] M. Ren, W. Zeng, B. Yang, and R. Urtasun, "Learning to reweight examples for robust deep learning," in *International Conference on Machine Learning*. PMLR, 2018, pp. 4334–4343.
- [60] J. Shu, Q. Xie, L. Yi, Q. Zhao, S. Zhou, Z. Xu, and D. Meng, "Meta-weight-net: Learning an explicit mapping for sample weighting," in *Advances in Neural Information Processing Systems*, 2019, pp. 1917–1928.
- [61] S. A. Friedler, C. Scheidegger, and S. Venkatasubramanian, "On the (im) possibility of fairness," *arXiv preprint arXiv:1609.07236*, 2016.

- [62] A. Julia, L. Jeff, M. Surya, and K. Lauren, "Machine bias: Theres software used across the country to predict future criminals. and its biased against blacks," *ProPublica*, 2016.
- [63] M. Yudell, D. Roberts, R. DeSalle, and S. Tishkoff, "Taking race out of human genetics," *Science*, vol. 351, no. 6273, pp. 564–565, 2016.
- [64] Google, "Freebase data dumps." <https://developers.google.com/freebase/data>, 2015.
- [65] I. Kemelmacher-Shlizerman, S. M. Seitz, D. Miller, and E. Brossard, "The megaface benchmark: 1 million faces for recognition at scale," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4873–4882.
- [66] P. J. Phillips, J. R. Beveridge, B. A. Draper, G. Givens, A. J. O'Toole, D. Bolme, J. Dunlop, Y. M. Lui, H. Sahibzada, and S. Weimer, "The good, the bad, and the ugly face challenge problem," *Image & Vision Computing*, vol. 30, no. 3, pp. 177–185, 2012.
- [67] L. Franceschi, P. Frasconi, S. Salzo, R. Grazzi, and M. Pontil, "Bilevel programming for hyperparameter optimization and meta-learning," in *International Conference on Machine Learning*. PMLR, 2018, pp. 1568–1577.
- [68] S. Sengupta, J.-C. Chen, C. Castillo, V. M. Patel, R. Chellappa, and D. W. Jacobs, "Frontal to profile face verification in the wild," in *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2016, pp. 1–9.
- [69] S. Moschoglou, A. Papaioannou, C. Sagonas, J. Deng, I. Kotsia, and S. Zafeiriou, "Agedb: the first manually collected, in-the-wild age database," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2017, pp. 51–59.
- [70] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell, "Decaf: A deep convolutional activation feature for generic visual recognition," in *International conference on machine learning*, 2014, pp. 647–655.
- [71] F. Wang, X. Xiang, J. Cheng, and A. L. Yuille, "Normface: 1 2 hypersphere embedding for face verification," in *Proceedings of the 25th ACM international conference on Multimedia*. ACM, 2017, pp. 1041–1049.
- [72] K. Simonyan, O. M. Parkhi, A. Vedaldi, and A. Zisserman, "Fisher vector faces in the wild," in *BMVC*, vol. 2, no. 3, 2013, p. 4.
- [73] Y. Ganin and V. Lempitsky, "Unsupervised domain adaptation by backpropagation," in *International Conference on International Conference on Machine Learning*, 2015, pp. 1180–1189.



Weihong Deng received the B.E. degree in information engineering and the Ph.D. degree in signal and information processing from the Beijing University of Posts and Telecommunications (BUPT), Beijing, China, in 2004 and 2009, respectively. From Oct. 2007 to Dec. 2008, he was a postgraduate exchange student in the School of Information Technologies, University of Sydney, Australia. He is currently a professor in School of Artificial Intelligence, BUPT. His research interests include trustworthy biometrics and affective computing, with a particular emphasis in face recognition and expression analysis. He has published over 100 papers in international journals and conferences, such as IEEE TPAMI, TIP, IJCV, CVPR and ICCV. He serves as area chair for major international conferences such as IJCB, FG, IJCAI, ACMMM, and ICME, guest editor for IEEE Transactions on Biometrics, Behavior, and Identity Science, and Image and Vision Computing Journal, and the reviewer for dozens of international journals and conferences. His Dissertation was awarded the outstanding doctoral dissertation award by Beijing Municipal Commission of Education. He has been supported by the programs for New Century Excellent Talents and Young Changjiang Scholar by Ministry of Education.



recognition, transfer learning and AI fairness.

Mei Wang received the B.E. degree in information and communication engineering from the Dalian University of Technology (DUT), Dalian, China, in 2013 and received M.E. degree in communication engineering from the Beijing University of Posts and Telecommunications (BUPT), Beijing, China, in 2016. From September 2018, she is a Ph.D. student in school of Artificial Intelligence of BUPT. Her research interests include computer vision, with a particular emphasis in face



Yaobin Zhang was born in Beijing, China, in 1996. He received his B.S. degree in Communication Engineering from Beijing University of Posts and Telecommunications, Beijing, China, in 2019. He is currently pursuing the M.S. degree in Information and Communication Engineering with Beijing University of Posts and Telecommunications. His research interests include deep learning and computer vision.