

DOMAIN ADAPTING ABILITY OF SELF-SUPERVISED LEARNING FOR FACE RECOGNITION

Chun-Hsien Lin, Student Member, IEEE, and Bing-Fei Wu, Fellow, IEEE

National Chiao Tung University
Institute of Electrical and Control Engineering
1001 University Road, Hsinchu, Taiwan 300, ROC

ABSTRACT

Although deep convolutional networks have achieved great performance in face recognition tasks, the challenge of domain discrepancy still exists in real world applications. Lack of domain coverage of training data (source domain) makes the learned models degenerate in a testing scenario (target domain). In face recognition tasks, classes in two domains are usually different, so classical domain adaptation approaches, assuming there are shared classes in domains, may not be reasonable solutions for this problem. In this paper, self-supervised learning is adopted to learn a better embedding space where the subjects in target domain are more distinguishable. The learning goal is maximizing the similarity between the embeddings of each image and its mirror in both domains. The experiments show its competitive results compared with prior works. To know the reason why it can achieve such performance, we further discuss how this approach affects the learning of embeddings.

Index Terms— Face recognition, domain adaptation, self-supervised learning

1. INTRODUCTION

With the growing of dataset and model capacity, the accuracy of face recognition is getting higher. Face recognition has been an efficient tool for authentication and has been widely deployed to many applications. Even if a learned model can perform nearly perfect in benchmark datasets, it may fail in some scenarios. The primary reason is the domain discrepancy between the training data (source domain) and testing data (target domain). The factors causing the domain discrepancy may be illumination, blur, pose, race, gender, or age. Most training images are collected on Internet, which makes the training set contain various domain information. However, only few samples come from same domain, so it is hard to train a model with good generalization. A straight forward solution is to fine-tune the learned model on target scenar-

ios with labels. This approach is not practice, however, because labeling enough data is too expensive. Therefore, how to leverage the unlabeled data in target domains becomes an important issue, and this is what domain adaptation does.

Current approaches of domain adaptation [1] can be divided into two main branches: adapting classifier and adapting embedding space. The former is tuning the classifiers to adapt to the distribution in target domains. The later is tuning the embedding functions (feature extractors) to find a common embedding space where the distributions of tow domains are aligned. As for the embedding distribution of source domain, areas near the cluster centers are the high density places. Aligning the distributions of two domains is similar to uniformly assigning the embeddings of target domain to the classes of source domain, so both approaches are relied on an assumption that source and target domains are shared classes. Oppositely, in face recognition tasks, the classes are different in domains. It is not reasonable to apply the approaches above to solve this unique problem. Despite the different preconditions, the techniques of embedding alignment is still widely adopted in most existing works, [2][3][4][5][6], to mitigate the domain shift of face recognition.

In this paper, to avoid using embedding alignment, we use self-supervised learning to maximize the self-similarity of each sample. In this way, the embedding distribution of target domain is not aligned to the clusters of source domain, but organizes itself. Tested in IJB-A [7], IJB-B [8], and IJB-C [9] datasets, the proposed approach shows its ability of domain transfer. It is interesting that we find the improvement is achieve by lowering inter-class similarity rather than enlarging intra-class similarity. Unlike prior works, [2][3][10], using larger models, like ResNet [11] or VGG [12], we use MobileFaceNet [13] as the backbone. In spite of smaller capacity, the testing results in experiments are competitive and even better. The contributions of this research can be summarized into two parts. First, we propose a novel method to perform domain transfer for face recognition. Second, we further analyze the mechanism beyond the learning algorithm.

This work was supported by the Ministry of Science Technology under Grant MOST 108-2638-E-009-001-MY2.

2. RELATED WORK

2.1. Domain adaptation for face recognition

Because of the unique issues of face recognition, few researches focus on domain adaptation for face recognition. Some prior works, [4][5][6], adopt embedding alignment methods to be the core algorithm directly, like Maximum Mean Discrepancy (MMD), [14], or adversarial learning, [15], which does not really reduce the domain shift existing in face recognition. To compensate embedding alignment, pseudo labels generated by clustering on target domain are utilized to train another cluster distribution are proposed in [2] and [3]. Arachchilage et al. [10] focus on the task of video frames. With the identity consistency in a video, clustering can be more accurate. After clustering, triplets are mined to fine-tune the model with a modified triplet loss. However, clustering relies on prior knowledges of class number or cluster margin, which may be a non-trivial work when setting these parameters.

2.2. Self-supervised learning

To leverage unlabeled data, some auxiliary tasks are designed in training losses based on the prior knowledges of data. For vision tasks, the simplest task is self-similarity. Chen et al. [16] augment an image, and adopt contrastive learning to minimizing their distances while maximizing the distances among different images in the embedding space. The method in [17] achieves better performance by adopting cosine similarity with stop-gradient operation. Its transfer ability is briefly illustrated, but it is limited. Therefore, we mainly refer this research to design an adapting loss for face recognition.

3. PROPOSED APPROACH

There are two datasets in our case. The source dataset (training dataset) is denoted as $X^s = \{x_i^s, y_i^s\}_{i=1}^N$, where x_i^s is a facial image in source domain, y_i^s is one-hot encoding label for x_i^s , and N is the number of source images. The target dataset is denoted as $X^t = \{x_j^t\}_{j=1}^M$, where x_j^t is a facial image in target domain, and M is the number of target images. For convenience, an arbitrary image, a source image and a target image are denoted as x , x^s and x^t respectively.

3.1. Learning embeddings

High recognition accuracy replies on a good embedding function. The purpose of an embedding function $f(x)$ is mapping images on a lower dimensional space where embeddings from same classes are closer while embeddings from different classes are farther. We can train an embedding function with an extra classifier, $\hat{y}(f(x))$, by cross entropy, which is simple

and robust. The cross entropy loss is defined as:

$$L_c = - \sum_{i=1}^N y_i^s \log \hat{y}(f(x_i^s)) \quad (1)$$

Since training dataset is usually large, to accelerate the convergence, focal loss [18] is adopted to encourage the learning of hard samples, so the loss can be modified as:

$$L_c = - \sum_{i=1}^N (1 - \hat{y}(f(x_i^s)))^\gamma y_i^s \log \hat{y}(f(x_i^s)) \quad (2)$$

where γ is a parameter for weighting down the loss caused by easy samples, and it is set to 2 according to [18]. Also, metric learning policies, like [19][20][21], can be applied to train a better embedding function.

3.2. Self-supervised learning for domain transfer

According to the studies in [16], minimizing the distances between the embeddings from random cropped and color distorted images can achieve the best result. However, applying color distortion may cause the racial bias, and we do not find obvious difference between using mirroring and random cropping in experiments. For each image, x , we only use it and its mirror, x' , for training, so a self-supervised learning loss called SimSiam [17] can be expressed as:

$$L_s = \frac{1}{2} [D(h(z'), \varphi(z)) + D(h(z), \varphi(z'))] \quad (3)$$

$$D(p, z) = - \frac{p^T z}{\|p\| \|z\|} \quad (4)$$

where $z = f(x)$, $z' = f(x')$, φ is a stop-gradient operation, and $D(p, z)$ is a function for computing the cosine similarity between embeddings p and z . Specially, h is a head which remap the embeddings in another view, and its effectiveness has been discussed in [16] and [17].

To perform domain adaptation, we apply the loss on both source and target datasets and sum them up by a ratio which we call it adapting ratio, ρ . Thus, the loss is expressed as:

$$L_a = (1 - \rho)L_s(x^s) + \rho L_s(x^t) \quad (5)$$

where ρ controls the weights of SimSiam loss on two domains. For the purpose of domain transfer, it is reasonable to set the value of ρ larger than 0.5 a little bit. Since the loss only focus on self-similarity, it is necessary to use cross entropy loss to maintain the inter-class discrepancy. The total loss is defined as:

$$L = L_c + L_a \quad (6)$$

We call it SimSiam Adapting (SSA) Loss, and use it to do the experiments in the next section.

4. EXPERIMENTS

4.1. Datasets

In the experiments, CASIA-WebFace [22] is utilized to be the source dataset. It contains 10,575 identities and 494,414 images which are the pictures of celebrities on Internet. As for the testing datasets, we use IJB datasets, including IJB-A/B/C [7][8][9]. IJB datasets contain mixture of images and videos in the wild. The conditions are all challenging due to blurry frames and large pose variations. IJB-A [7] contains 500 identities with 5,396 images and 20,412 video frames. IJB-B [8] extended from IJB-A contains 1,845 identities with 11,755 images and 7,018 videos, and 10,044 non-face images. IJB-C [9] extended from IJB-B contains 3,531 identities with 21,294 images and 11,779 videos, and 10,040 non-face images.

We follow the protocols in [7] to evaluate the performance of verification, open-set and close-set identification, and there metrics are True Positive Rate vs. False Positive Rate (TPR@FAR), True Positive Identification Rate vs. set False Positive Identification Rate (TPIR@FAIR), and top-K accuracy (Rank-K) respectively.

4.2. Implementation detail

Cross entropy loss (2) is applied to train MobileFaceNet [13] from scratch by CASIA-WebFace dataset [22]. The training epoch and batch size is set to 50 and 128 respectively. The learning rate is set to 0.1 and is divided by 10 every 12 epoch. The trained model is the baseline for domain adaptation, and its performance is the baseline in the experiments.

We treat all images in IJB-A dataset [7] as the target dataset, so most images in IJB-B and IJB-C are not covered. As the unlabeled data from a target dataset are mixed in, the learning rate is decayed to be 0.0001 to protect the learned knowledge. The training epoch, batch size and learning rate schedule are preserved, but the epoch is counted based on the size of the target dataset, which induces less training iterations. For each batch, 128 images are sampled from each dataset.

The head, h , used in SimSiam loss is a two-layer neural network, and the number of hidden neurons is set to 32 which is 4 times less than the embedding dimension of MobileFaceNet [13]. According to [17], we add batch normalization and rectified linear unit in the hidden layer to guarantee its performance.

4.3. Ablation study

It is more efficient to use smaller dataset to evaluate the trained models in ablation study, so only IJB-A dataset [7] is adopted in this part of experiments.

4.3.1. Adapting ratio

To find out the proper adapting ratio, ρ , we vary its value, and compare with the baseline. The comprehensive results are listed in Table 1. It is obvious that the performances of verification at lower FPR and open-set identification are increased. Adopting SimSiam loss [17] directly on source dataset only, $\rho = 0.0$ can achieve some improvements, which shows its competency of generalization. Since larger weighting on target domain may limit the cluster learning relying on the supervision of source data and labels, overemphasis of learning self-similarity on target domain cannot preserve the inter-class discrepancy. As for verification at higher FPR and close-set identification, there is no improvement. Our hypothesis for this issue is discussed in the next part. According to the results in Table 1, we choose $\rho = 0.6$ as our best parameter setting.

4.3.2. Embedding analysis

How the embedding spaces are learned is discussed here. We compare the embedding distributions of the baseline, source-only case, $\rho = 0.0$, and the best case, $\rho = 0.6$. The averages of three similarities (mirror, intra-class, and inter-class) and embedding length are carried out in Table 2.

Although the mirror similarities are increased with the guidance of SSA, there is no obvious changes on intra-class similarities, but they are a little bit lower. Oppositely, the inter-class similarities are reduced, and the embedding lengths are also enlarged, which implies the discrepancies among identities are increased. Since higher TPR or TPIR at lower FPR and FPIR requires lower inter-class similarity, SSA can do better under these protocols. Due to giant negative pairs in the protocols, the significant improvements can be achieved by little decay on inter-class similarity. On the other hand, higher TPR at higher FPR or Rank-K requires much higher intra-class similarity, so this is the reason why SSA fails under these protocols.

Table 1. Evaluations on IJB-A [7] with different ρ .

Method	Verification TPR (%)			Identification TPIR (%)			
	FPR=0.001	FPR=0.01	FPR=0.1	FPIR=0.01	FPIR=0.1	Rank-1	Rank-10
Baseline	75.63	90.54	96.88	65.31	85.74	94.79	97.91
$\rho = 0.0$	79.37	90.93	96.86	71.59	87.85	94.57	97.85
$\rho = 0.5$	80.78	91.03	96.70	73.23	87.82	94.74	97.78
$\rho = 0.6$	82.13	91.45	96.39	75.77	87.92	93.99	97.63
$\rho = 0.7$	79.57	90.86	96.28	70.11	86.71	94.13	97.63
$\rho = 0.8$	78.29	90.08	95.99	68.12	85.76	93.93	97.19
$\rho = 0.9$	74.33	90.12	96.06	64.36	84.91	93.72	97.24

Table 2. Statistics of embedding metrics in IJB-A dataset [7].

Method	Similarity			Embedding Length
	Mirror	Intra-class	Inter-class	
Baseline	0.9478	0.7074	0.0728	113.89
SSA ($\rho = 0.0$)	0.9583	0.6915	0.0199	119.68
SSA ($\rho = 0.6$)	0.9547	0.6905	0.0166	119.55

Table 3. Verification performance on IJB-A [7], IJB-B [8], and IJB-C [9]. The bold texts stand for the highest TPIR in a column. The text with underline means it is better than baseline.

Method	IJB-A TPIR (%)				IJB-B TPIR (%)				IJB-C TPIR (%)			
	FPR=0.0001	FPR=0.001	FPR=0.01	FPR=0.1	FPR=0.0001	FPR=0.001	FPR=0.01	FPR=0.1	FPR=0.0001	FPR=0.001	FPR=0.01	FPR=0.1
Sohn et al. [5]	-	58.40	82.80	96.20	-	-	-	-	-	-	-	-
IMAN-A [3]	-	84.49	91.88	97.05	-	-	-	-	-	-	-	-
CDA(vgg-soft) [2]	-	74.76	89.76	98.19	-	-	-	-	-	-	-	-
CDA(res-arc) [2]	-	82.45	91.11	96.96	-	87.35	94.55	98.08	-	88.06	94.85	98.33
SoftMax ^a	52.23	75.63	90.54	96.88	68.91	83.61	93.43	98.22	74.04	86.44	94.59	98.54
ArcFace ^b [19]	72.60	84.82	92.18	96.11	77.29	87.18	94.25	98.26	81.33	89.75	95.35	98.48
SSA-SoftMax (ours)	<u>62.51</u>	<u>82.13</u>	<u>91.45</u>	96.39	<u>71.22</u>	<u>84.88</u>	<u>93.78</u>	98.21	<u>75.61</u>	<u>87.47</u>	<u>94.92</u>	98.47
SSA-ArcFace (ours)	<u>78.18</u>	<u>87.37</u>	<u>92.41</u>	95.84	<u>78.48</u>	<u>88.27</u>	<u>94.88</u>	<u>98.49</u>	<u>82.72</u>	<u>90.91</u>	<u>95.90</u>	<u>98.62</u>

^aThis is baseline model trained by cross entropy on source dataset only.

^bThis is another baseline model trained by the guidance of margin penalty proposed in [19].

Table 4. Identification performance on IJB-A [7], IJB-B [8], and IJB-C [9]. The bold texts stand for the highest TPIR in a column. The text with underline means it is better than baseline.

Method	IJB-A TPIR (%)				IJB-B TPIR (%)				IJB-C TPIR (%)			
	FPIR=0.01	FPIR=0.1	Rank-1	Rank-10	FPIR=0.01	FPIR=0.1	Rank-1	Rank-10	FPIR=0.01	FPIR=0.1	Rank-1	Rank-10
Sohn et al. [5]	-	-	87.90	97.00	-	-	-	-	-	-	-	-
IMAN-A [3]	-	-	94.05	98.04	-	-	-	-	-	-	-	-
CDA(vgg-soft) [2]	66.85	85.32	94.89	99.23	-	-	-	-	-	-	-	-
CDA(res-arc) [2]	75.49	87.76	93.61	97.62	-	-	86.22	93.33	-	-	88.19	93.70
SoftMax ^a	65.31	85.74	94.79	97.91	59.77	77.10	88.01	95.22	58.86	76.49	88.93	95.07
ArcFace ^b [19]	78.99	88.75	94.58	97.50	66.22	81.23	89.49	95.37	70.49	81.87	90.61	95.60
SSA-SoftMax (ours)	<u>75.77</u>	<u>87.92</u>	93.99	97.63	58.74	<u>77.51</u>	86.98	94.41	57.24	<u>77.69</u>	87.92	94.50
SSA-ArcFace (ours)	<u>80.03</u>	<u>89.47</u>	94.26	97.40	64.47	<u>82.48</u>	89.31	95.23	69.80	<u>83.68</u>	90.50	95.47

^aThis is baseline model trained by cross entropy on source dataset only.

^bThis is another baseline model trained by the guidance of margin penalty proposed in [19].

4.4. Benchmark comparison

To guarantee the effectiveness of SSA, it is compared with the state-of-the-arts focusing on domain adaptation for face recognition. We also use a margin penalty method (ArcFace) proposed in [19] to train our backbone to be another baseline. The comprehensive evaluations on IJB-A [7], IJB-B [8], and IJB-C [9] are listed in Table 3 and Table 4.

From Table 3, we can observe that SSA can successfully improve the baselines on all benchmarks under almost all FPRs especially under lower FPRs. With the guidance of ArcFace [19], the performance can be better. However, in the identification protocols, Table 4, the improvements only exists on open-set protocols. Such issue has been discussed in 4.3. Since the models are adapted on IJB-A only [7], the performances of on the open-set protocols of IJB-B [8] and IJB-C [9] are not that good, but it can be further refined by adapting more data from these datasets. Expect for close-set identification protocols, compared with the state-of-the-arts, our approach shows its good performance by not only the improvements but also the much lighter backbone.

5. CONCLUSION

We focus on the unique problem of domain discrepancy in face recognition whose classes in domains are non-overlapping. Self-Supervised Adapting (SSA) loss is proposed in this pa-

per. By adding an adapting ratio between the self-similarity losses on source and target domain, SSA can successfully improve the baseline models both verification and open-set identification protocols. Interestingly, we find that this progress is achieved by reducing inter-class similarities rather than increasing intra-class similarities through the analysis on the embedding distributions. Compared with other adapting methods under comprehensive protocols, SSA shows its competitive performance. However, it seems that SSA cannot preserve or even improve the intra-class similarity on target domain, so some advanced researches should be done in the future to compensate this problem.

6. REFERENCES

- [1] Mei Wang and Weihong Deng, “Deep visual domain adaptation: A survey,” *Neurocomputing*, vol. 312, pp. 135–153, 2018.
- [2] Mei Wang and Weihong Deng, “Deep face recognition with clustering based domain adaptation,” *Neurocomputing*, vol. 393, pp. 1–14, 2020.
- [3] Mei Wang, Weihong Deng, Jiani Hu, Xunqiang Tao, and Yaohai Huang, “Racial faces in the wild: Reducing racial bias by information maximization adaptation network,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 692–702.

- [4] Zimeng Luo, Jiani Hu, Weihong Deng, and Haifeng Shen, "Deep unsupervised domain adaptation for face recognition," in *13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*. IEEE, 2018, pp. 453–457.
- [5] Kihyuk Sohn, Sifei Liu, Guangyu Zhong, Xiang Yu, Ming-Hsuan Yang, and Manmohan Chandraker, "Unsupervised domain adaptation for face recognition in unlabeled videos," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3210–3218.
- [6] Sungeun Hong, Woobin Im, Jongbin Ryu, and Hyun S Yang, "SSPP-DAN: Deep domain adaptation network for face recognition with single sample per person," in *2017 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2017, pp. 825–829.
- [7] Brendan F Klare, Ben Klein, Emma Taborsky, Austin Blanton, Jordan Cheney, Kristen Allen, Patrick Grother, Alan Mah, and Anil K Jain, "Pushing the frontiers of unconstrained face detection and recognition: IARPA Janus Benchmark A," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1931–1939.
- [8] Cameron Whitelam, Emma Taborsky, Austin Blanton, Brianna Maze, Jocelyn Adams, Tim Miller, Nathan Kalka, Anil K Jain, James A Duncan, Kristen Allen, et al., "IARPA Janus Benchmark-B face dataset," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2017, pp. 90–98.
- [9] Brianna Maze, Jocelyn Adams, James A Duncan, Nathan Kalka, Tim Miller, Charles Otto, Anil K Jain, W Tyler Niggel, Janet Anderson, Jordan Cheney, et al., "IARPA Janus Benchmark-C: Face dataset and protocol," in *International Conference on Biometrics (ICB)*. IEEE, 2018, pp. 158–165.
- [10] Samadhi Wickrama Arachchilage and Ebroul Izquierdo, "SSDL: Self-supervised domain learning for improved face recognition," *arXiv preprint arXiv:2011.13361*, 2020.
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [12] Omkar M. Parkhi, Andrea Vedaldi, and Andrew Zisserman, "Deep face recognition," in *Proceedings of the British Machine Vision Conference (BMVC)*. September 2015, pp. 41.1–41.12, BMVA Press.
- [13] Sheng Chen, Yang Liu, Xiang Gao, and Zhen Han, "MobileFaceNets: Efficient CNNs for accurate real-time face verification on mobile devices," in *Chinese Conference on Biometric Recognition*. Springer, 2018, pp. 428–438.
- [14] Mingsheng Long, Yue Cao, Zhangjie Cao, Jianmin Wang, and Michael I Jordan, "Transferable representation learning with deep adaptation networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 12, pp. 3071–3085, 2018.
- [15] Yaroslav Ganin and Victor Lempitsky, "Unsupervised domain adaptation by backpropagation," in *International Conference on Machine Learning*. PMLR, 2015, pp. 1180–1189.
- [16] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
- [17] Xinlei Chen and Kaiming He, "Exploring simple siamese representation learning," *arXiv preprint arXiv:2011.10566*, 2020.
- [18] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2980–2988.
- [19] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou, "ArcFace: Additive angular margin loss for deep face recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4690–4699.
- [20] Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei 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.
- [21] Weiyang Liu, Yandong Wen, Zhiding Yu, Ming Li, Bhiksha Raj, and Le Song, "SphereFace: Deep hypersphere embedding for face recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 212–220.
- [22] Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z Li, "Learning face representation from scratch," *arXiv preprint arXiv:1411.7923*, 2014.