

Age and Gender Prediction From Face Images Using Attentional Convolutional Network

Amirali Abdolrashidi¹, Mehdi Minaei², Elham Azimi³, Shervin Minaee⁴

¹University of California, Riverside

²Sama Technical College, Azad University

³New York University

⁴Snap Inc

Abstract—Automatic prediction of age and gender from face images has drawn a lot of attention recently, due it is wide applications in various facial analysis problems. However, due to the large intra-class variation of face images (such as variation in lighting, pose, scale, occlusion), the existing models are still behind the desired accuracy level, which is necessary for the use of these models in real-world applications. In this work, we propose a deep learning framework, based on the ensemble of attentional and residual convolutional networks, to predict gender and age group of facial images with high accuracy rate. Using attention mechanism enables our model to focus on the important and informative parts of the face, which can help it to make a more accurate prediction. We train our model in a multi-task learning fashion, and augment the feature embedding of the age classifier, with the predicted gender, and show that doing so can further increase the accuracy of age prediction. Our model is trained on a popular face age and gender dataset, and achieved promising results. Through visualization of the attention maps of the train model, we show that our model has learned to become sensitive to the right regions of the face.

I. INTRODUCTION

Age and gender information are very important for various real world applications, such as social understanding, biometrics, identity verification, video surveillance, human-computer interaction, electronic customer, crowd behavior analysis, on-line advertisement, item recommendation, and many more. Despite their huge applications, being able to automatically predicting age and gender from face images is a very hard problem, mainly due to the various sources of intra-class variations on the facial images of people, which makes the use of these models in real world applications limited.

There are numerous works proposed for age and gender prediction in the past several years. The earlier works were mainly based on hand-crafted features extracted facial images followed by a classifier. But with the great success of deep learning models in various computer vision problems in the past decade [1]–[5], the more recent works on age and gender predictions are mostly shifted toward deep neural networks based models.

In this work, we propose a deep learning framework to jointly predict the age and gender from face images. Given the intuition that some local regions of the face have more clear signals about the age and gender of an individual (such as beard and mustache for male, and wrinkles around eyes and mouth for age), we use an attentional convolutional network

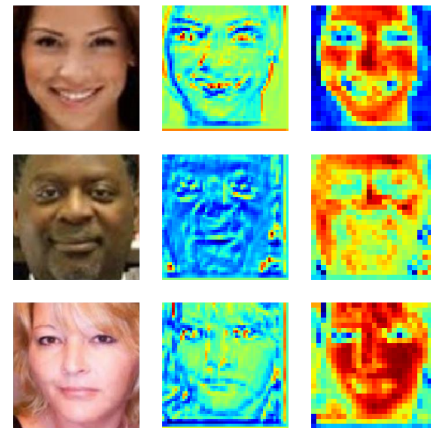


Fig. 1. Three sample face images, and their corresponding attention map outputs, from two different layers. Color scale is from blue (lowest) to red (highest).

as one of our backbone models, to better attend to the salient and informative part of the face. Figure 1 provide three sample images, and the corresponding attention map outputs of two different layers of our model for these images. As we can see, the model outputs are mostly sensitive to the edge patterns around facial parts, as well as wrinkles, which are important for age and gender prediction.

As predicting age and gender from faces are very related, we use a single model with multi-task learning approach to jointly predict both gender and age bucket. Also, given that knowing the gender of someone, we can better estimate her/his age, we augment the feature of the age-prediction branch with the predicted gender output. Through experimental results, we show that adding the predicted gender information to the age-prediction branch, improves the model performance. To further improve the prediction accuracy of our model, we combine the prediction of attentional network with the residual network, and use their ensemble model as the final predictor.

Here are the contributions of this work:

- We propose a multi-task learning framework to jointly predict the age and gender of individuals from their face images.
- We develop an ensemble of attentional and residual

networks, which outperforms both individual models. The attention layers of our model learn to focus on the most important and salient parts of the face.

- We further propose to feed the predicted gender label to the age prediction branch, and show that doing this will improve the accuracy of age prediction branch.
- With the help of the attention mechanism, we can explain the predictions of the classifiers after they are trained, by locating the salient facial regions they are focusing on each image.

The structure of the remaining parts of this paper is as follows. Section II provides an overview of some of the previous works on age and gender prediction. Section III provides the details of our proposed framework, and the architecture of our multi-task learning model. Section IV, provides a quick overview of the dataset used in our framework. Then, in Section V, we provide the experimental studies, the quantitative performance of our model, and also visual evaluation of model outputs. Finally the paper is concluded in Section VI.

II. RELATED WORKS

Face is one of the most popular biometrics (along with fingerprint, iris, and palmprint [6]–[11]), and face recognition and facial attributes/characteristics prediction have attracted a lot of attention in the past few decades [12]–[14]. Age and gender prediction from face images, as a specific face analysis problem have also drawn attention in the recent years, and there have been several previous works for age and gender prediction from face images. Here we provide an overview of some of the most promising works.

In [15], Levi and Hassner proposed a simple convolutional neural network architecture with 5 layers, to perform age and gender prediction. Despite the simplicity of this model, they achieved promising results on Adience benchmark for age and gender estimation. Figure 2 shows the architecture of the model proposed in [15].

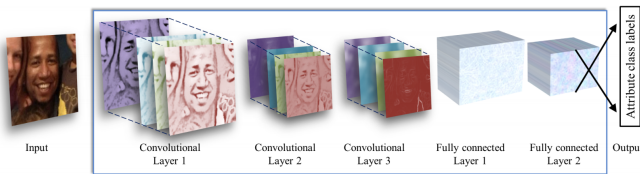


Fig. 2. The architecture of the work by Levi and Hassner, courtesy of [15].

In [16], Duan et al. proposed a hybrid structure, which combines Convolutional Neural Network (CNN) and Extreme Learning Machine (ELM), to perform age and gender classification. CNN is used to extract the features from the input images, while ELM classifies the intermediate results. They were able to obtain reasonable accuracy on the MORPHII and Adience Benchmarks.

In [17], Ozbulak et al. analyzed the transferability of existing deep convolutional neural network (CNN) models for age and gender classification. The generic AlexNet-like architecture and domain specific VGG-Face CNN model are fine-tuned with the Adience dataset prepared for age and gender

classification in uncontrolled environments. Not surprisingly, they were able to obtain promising results on the features learned from these popular architectures.

In [18], Lapuschkin et al. compared four popular neural network architectures, studied the effect of pre-training, evaluated the robustness of the considered alignment preprocessings via cross-method test set swapping, and intuitively visualized the model’s prediction strategies in given pre-processing conditions using the Layer-wise Relevance Propagation (LRP) algorithm. They were able to obtain very interesting relevance maps for some of the popular model architectures, as shown in Figure 3. [15].

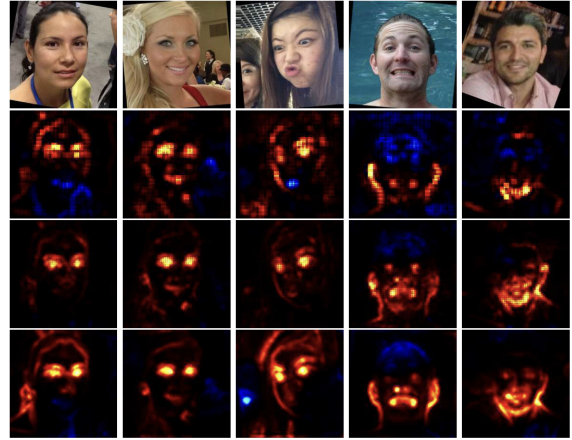


Fig. 3. From top to bottom: Input image, followed by relevance maps for the best performing CaffeNet, GoogleNet and the VGG16 model for gender prediction. Courtesy of [18].

In [19], Rodriguez et al. proposed a novel feedforward attention mechanism that is able to discover the most informative and reliable parts of a given face for improving age and gender classification. More specifically, given a downsampled facial image, the proposed model is trained based on an end-to-end learning framework to extract the most discriminative patches from the original high-resolution image. They were able to obtain promising results on several benchmarks, including Adience, Images of Groups, and MORPH II. The block-diagram of this work is shown in Figure 4.

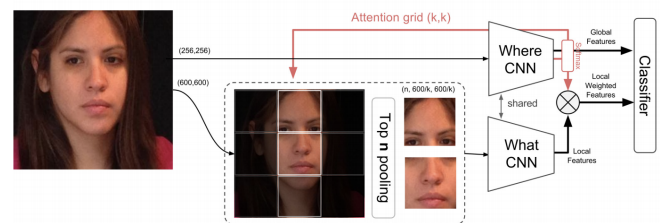


Fig. 4. The proposed attention model by Rodriguez et al. Courtesy of [19].

Some of the other promising works for age and gender predictions includes: adversarial spatial frequency domain critic learning [20], Region-SIFT and multi-layered SVM [21], and Landmark-Guided Local Deep Neural Networks [22].

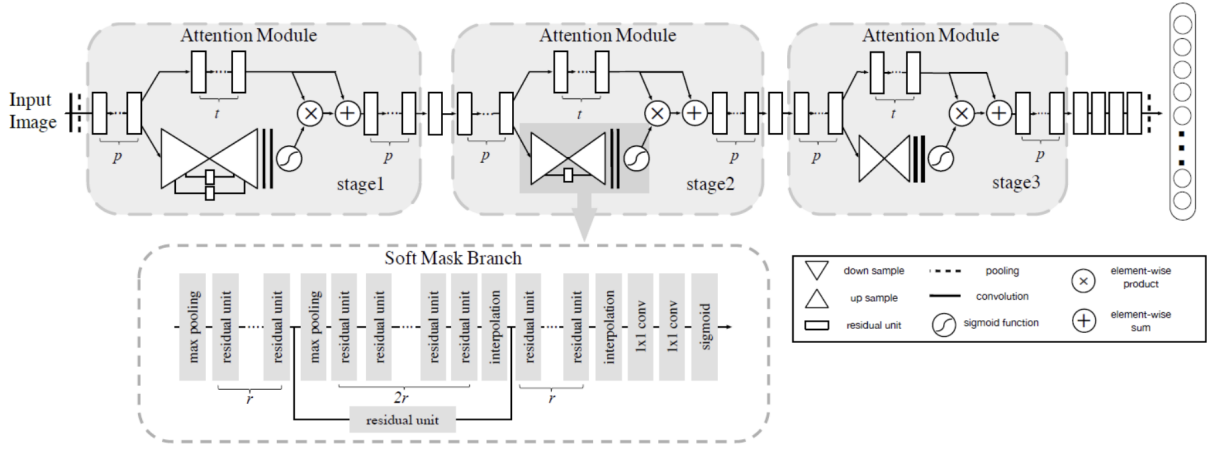


Fig. 5. The architecture of residual attention network, courtesy of [23].

III. THE PROPOSED FRAMEWORK

In this section we provide the details of the proposed age and gender prediction framework. We formulate this as a multi-task learning problem, in which a single model is used to predict both gender and age-buckets simultaneously. In another word, a single convolutional neural network with two heads (output branches) is used to jointly predict age and gender. Figure 6 shows the block-diagram of a simple multi-task learning model for joint age and gender prediction.

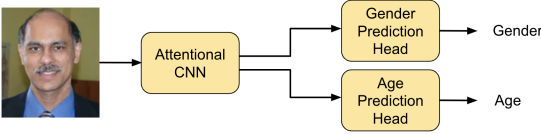


Fig. 6. The block-diagram of a multi-task learning network for joint age and gender prediction.

Given the intuition that knowing one's gender, can enable us to better predict her/his age, we augment the input feature of the age prediction part of the model, with the the predicted-gender from the other head. This can help us the age model to have access to a rough estimation of gender. Through experimental study, we show that doing so improves the performance of the age prediction. Figure 7, provides the block diagram of the proposed model architecture.

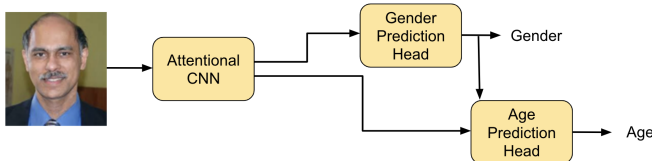


Fig. 7. The block-diagram of the framework for joint age and gender prediction, in which the predicted gender is used along with the neural embedding as the input for the age prediction branch.

To further boost the accuracy of our model, we propose to use an ensemble of attentional network, and a residual convolutional network. Once these models are trianed, their

average predictions (output probabilities) are used as the final prediction. We will give more details on the architecture of each of these two models in the below parts.

A. Residual Attentional Network

As mentioned above, an important piece of our framework is the residual attentional network (RAN), which is a convolutional neural network with attention mechanism, which can be incorporated with state-of-art feed forward network architecture in an end-to-end training fashion [23]. This network is built by stacking attention modules, which generate attention-aware features that adaptively changes as layers go deeper into the network.

The composition of the Attention Module includes two branches: the trunk branch and the mask branch. Trunk Branch performs feature processing with Residual Units. Mask Branch uses bottom-up top-down structure softly weight output features with the goal of improving trunk branch features. Bottom-Up Step: collects global information of the whole image by downsampling (i.e. max pooling) the image. Top-Down Step: combines global information with original feature maps by upsampling (i.e. interpolation) to keep the output size the same as the input feature map. The full architecture of residual attention network is shown in Figure 5.

B. ResNet Model

Another model used in our framework is based on residual convolutional network (ResNet) [24]. ResNet is known to have a better gradient flow by providing the skip connection in each residual block. Here we use a ResNet architecture to perform gender classification on the input image. Then, the predicted output of the gender branch is concatenated with the last hidden layer of the age branch. The overall block diagram of ResNet18 model is illustrated in Figure 8. ResNet50 architecture is pretty similar to ResNet18, the main difference being having more layers.

C. Ensemble of Attentional and Residual Networks

To boost the prediction accuracy of our model, we use the ensemble of the above two architectures, and combine their

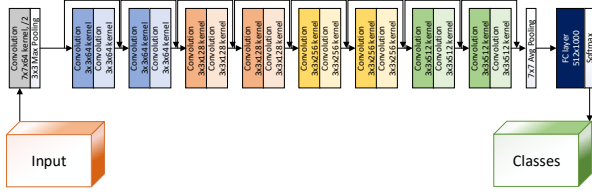


Fig. 8. The architecture of ResNet18 model. Courtesy of [24].

predictions. Through experimental results, we show that the ensemble model achieves higher performance than both of the individual models.

IV. UTK-FACE DATASET

In this section, we provide a quick overview of the dataset used in this work, UTK-Face dataset [25]. We used the cropped and aligned version of this dataset, which consists of 9780 images of size 200x200 (4372 male, 5408 female). The images are taken from the faces of people from various ethnic groups and ages (1 to 110 years). The age distribution of the used dataset is shown in Figure 9. It can be seen that the largest portion of the images belong to people under 2 years of age.

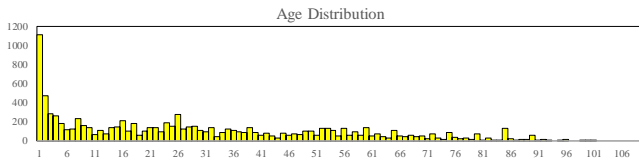


Fig. 9. Age distribution in the used UTK-Face dataset.

Also, Figure 10 denotes several images from UTK-Face dataset. As we can see from this figure, UTK-Face dataset has a good diversity in terms of age and gender. It is worth to note that, in order to make the age prediction simpler, we group people whose age is in some range into the same age-bucket (such as 0-10), and try to predict their age-bucket instead of the exact age. Doing so, on one hand makes age prediction simpler, and on the other hand turns this model from regression to classification (which does not need any clipping of the predicted values).

V. EXPERIMENTAL RESULTS

Before presenting the quantitative and qualitative performance of the proposed framework, let us first discuss about the hyper-parameter values used in our training procedure. These values are tuned based on the model performance on the validation set.

A. Model Hyper-parameter Values

The proposed model has been trained for over 100 epochs on an Nvidia Tesla GPU. The batch size is set to 16, and number of workers to 8. ADAM optimizer with a learning rate of 0.005 is used to optimize the loss function. PyTorch library is used for the implementation of experimental studies [26].



Fig. 10. Sample images from UTK-Face database.

B. Model Quantitative Performance

In this part, we provide the model performance in terms of various metrics, such as classification accuracy for both branches (gender prediction, and age-bucket classification), as well as "average age-bucket absolute difference (AABD)". AABD calculates the absolute difference between the ground-truth and the predicted age-bucket of each person, and average them out among all test samples. As an example, if someone's age is 25 (its age-bucket being 20-30, therefore 3), and the predicted age by model is 4, then the AABD value will be 1.

As mentioned earlier, each model predicts a probability vector, which shows the likelihood of that image being in each of the possible classes (for example, male vs female, for gender model). In the ensemble model, the predicted probabilities from Attentional-CNN and ResNet are averaged and used to infer the predicted classes of each task. In Table I, we present the comparison between the proposed ensemble model, with each of the two individual models. As we can see the ensemble model outperforms both the individual models on both tasks (gender and age prediction). In Table II, we compare the performance of the proposed model with one of the previous works, on gender classification task.

TABLE I. COMPARISON OF ACCURACY OF DIFFERENT MODELS

Model	Age Range Acc	Gender Acc	Age Bucket Diff
Attention CNN	0.742	0.552	0.33
ResNet	0.900	0.965	0.14
The Ensemble Model	0.913	0.965	0.11

TABLE II. COMPARISON WITH A PREVIOUS WORK.

Model	Gender Acc
Automatic Gender Detection [27]	0.9485
Our model	0.965

C. Model Predicted Scores

Figure 11 shows the histogram of the predicted probability scores for gender by our model, for each class. Note that here we are showing the histogram of the likelihood of an image being female. It can be seen that in the majority of the cases, the model's prediction of gender is much higher in the extremes, i.e. the model makes its decision with remarkable

certainty. On the other hand, more uncertain decisions, which are shown in the middle of the chart, are much lower in frequency.

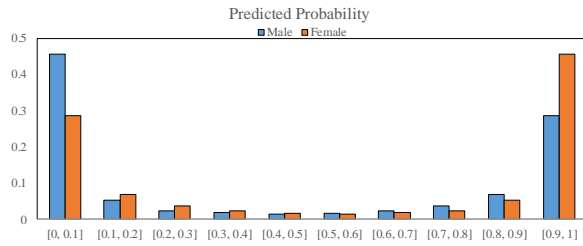


Fig. 11. Predicted probability distribution for gender classification

D. Model Predictions for Some Sample Images

We also present the model prediction for some of the sample images of our test set. In Figure 12, some of the model’s predictions are provided next to the true labels of the face image. It can be seen that, for most of these sample images, the true age label falls within the predicted age range by our model.

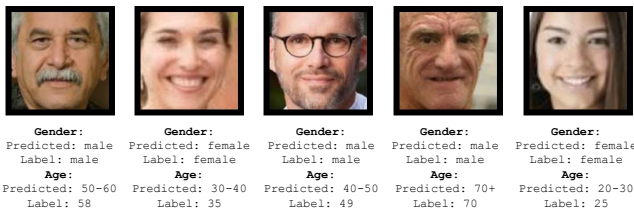


Fig. 12. Classification results for sample images

E. Model Confusion Matrix

To provide a more detailed analysis of the proposed model’s performance, we also present the confusion matrix of this model for each of the two tasks. Figure 13 and Figure 14 show the confusion matrix for age range and gender classification respectively. It can clearly be seen that a vast majority of the cases are predicted as the true label of their class, seen on the main diagonal of the matrix. It is worth noting the largest false positive percentage in the age range classification model, corresponds to the images belonging to the 30-40 age range which are mistaken for 20-30.

F. Model Attention Map Visualization

In Figure 15, we show the output of two of the attention maps (from the three attention modules in our model) of the trained model for a few sample images from the test set. As it can be seen from this figure, our model is learned to extract features from different salient parts of the face, such as the outline, eyes, and wrinkles, which sounds reasonable when trying to make prediction about someone’s age and gender.

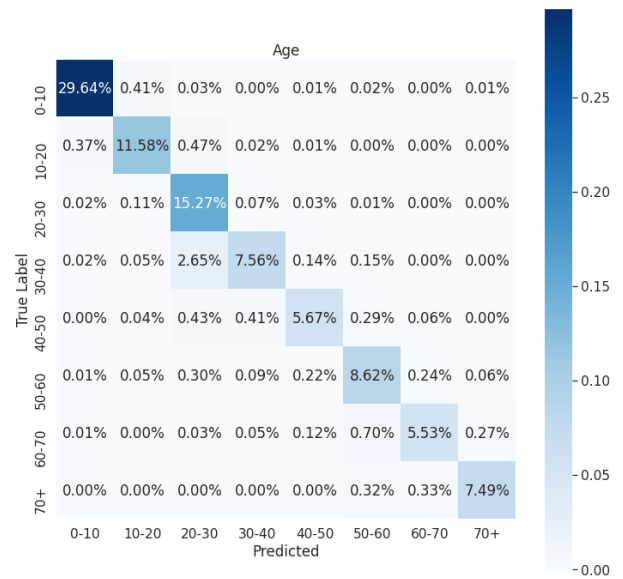


Fig. 13. Confusion matrix for age range

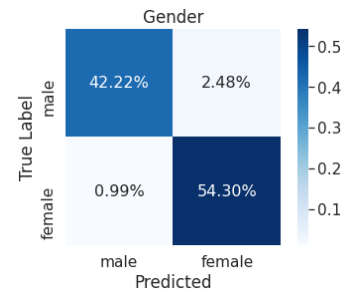


Fig. 14. Confusion matrix for gender outputs

VI. CONCLUSION

In this work we propose a multi-task learning framework, to simultaneously predict age and gender from face images. Our framework is based on an ensemble of ResNet-based model and an attention-based model. We trained and tested the proposed model on the UTKFace dataset consisting a large variety of faces from different ages, genders and ethnicities. Through experimental studies, we show that the prediction accuracy of the ensemble model (for both age and gender prediction tasks) surpasses in those of the separated models. We also showed that providing the prediction of the gender model as one of the input signal for the age-prediction branch, can improve the accuracy of predicted age values. Through visualization of the attention maps of the trained model, we show that the model learned to focus on the most salient part of the face, useful for predicting age and gender.

REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [2] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *Advances in neural information processing systems*, 2015, pp. 91–99.

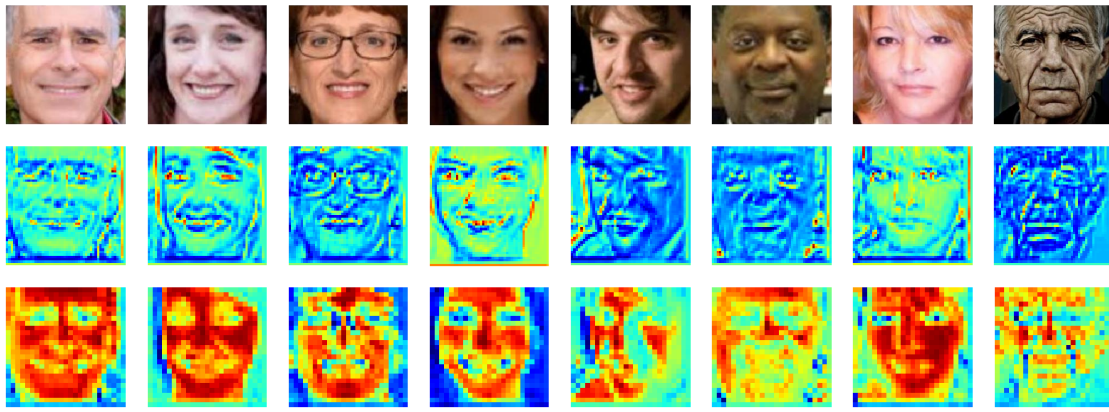


Fig. 15. The attention maps of our model, for eight sample images.

- [3] S. Minaee, A. Abdolrashidi, H. Su, M. Bennamoun, and D. Zhang, "Biometric recognition using deep learning: A survey," *arXiv preprint arXiv:1912.00271*, 2019.
- [4] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2017.
- [5] S. Minaee, Y. Wang, A. Aygar, S. Chung, X. Wang, Y. W. Lui, E. Fieremans, S. Flanagan, and J. Rath, "Mtbi identification from diffusion mr images using bag of adversarial visual features," *IEEE transactions on medical imaging*, vol. 38, no. 11, pp. 2545–2555, 2019.
- [6] Q. Zhao, L. Zhang, D. Zhang, and N. Luo, "Direct pore matching for fingerprint recognition," in *International Conference on Biometrics*. Springer, 2009, pp. 597–606.
- [7] S. Minaee and Y. Wang, "Fingerprint recognition using translation invariant scattering network," in *2015 IEEE Signal Processing in Medicine and Biology Symposium (SPMB)*. IEEE, 2015, pp. 1–6.
- [8] M. De Marsico, A. Petrosino, and S. Ricciardi, "Iris recognition through machine learning techniques: A survey," *Pattern Recognition Letters*, vol. 82, pp. 106–115, 2016.
- [9] S. Minaee, A. Abdolrashidi, and Y. Wang, "An experimental study of deep convolutional features for iris recognition," in *2016 IEEE signal processing in medicine and biology symposium (SPMB)*. IEEE, 2016, pp. 1–6.
- [10] S. Minaee and A. Abdolrashidi, "Highly accurate palmprint recognition using statistical and wavelet features," in *2015 IEEE Signal Processing and Signal Processing Education Workshop (SP/SPE)*. IEEE, 2015, pp. 31–36.
- [11] S. A. Mistani, S. Minaee, and E. Fatemizadeh, "Multispectral palmprint recognition using a hybrid feature," *arXiv preprint arXiv:1112.5997*, 2011.
- [12] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, and Y. Ma, "Robust face recognition via sparse representation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 31, no. 2, pp. 210–227, 2008.
- [13] O. M. Parkhi, A. Vedaldi, and A. Zisserman, "Deep face recognition," 2015.
- [14] S. Minaee, A. Abdolrashidi, and Y. Wang, "Face recognition using scattering convolutional network," in *2017 IEEE Signal Processing in Medicine and Biology Symposium (SPMB)*. IEEE, 2017, pp. 1–6.
- [15] G. Levi and T. Hassner, "Age and gender classification using convolutional neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2015, pp. 34–42.
- [16] M. Duan, K. Li, C. Yang, and K. Li, "A hybrid deep learning cnn-elm for age and gender classification," *Neurocomputing*, vol. 275, pp. 448–461, 2018.
- [17] G. Ozbulak, Y. Aytar, and H. K. Ekenel, "How transferable are cnn-based features for age and gender classification?" in *2016 International Conference of the Biometrics Special Interest Group (BIOSIG)*. IEEE, 2016, pp. 1–6.
- [18] S. Lapuschkin, A. Binder, K.-R. Muller, and W. Samek, "Understanding and comparing deep neural networks for age and gender classification," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2017, pp. 1629–1638.
- [19] P. Rodríguez, G. Cucurull, J. M. Gonfaus, F. X. Roca, and J. González, "Age and gender recognition in the wild with deep attention," *Pattern Recognition*, vol. 72, pp. 563–571, 2017.
- [20] S. S. Lee, H. G. Kim, K. Kim, and Y. M. Ro, "Adversarial spatial frequency domain critic learning for age and gender classification," in *2018 25th IEEE International Conference on Image Processing (ICIP)*. IEEE, 2018, pp. 2032–2036.
- [21] H. Kim, S.-H. Lee, M.-K. Sohn, and B. Hwang, "Age and gender estimation using region-sift and multi-layered svm," in *Tenth International Conference on Machine Vision (ICMV 2017)*, vol. 10696. International Society for Optics and Photonics, 2018, p. 106962J.
- [22] Y. Zhang and T. Xu, "Landmark-guided local deep neural networks for age and gender classification," *Journal of Sensors*, vol. 2018, 2018.
- [23] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang, "Residual attention network for image classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3156–3164.
- [24] 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.
- [25] 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*, 2017, pp. 5810–5818.
- [26] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," 2017.
- [27] R. R. Nair, R. Madhavankutty, and S. Nema, "Automated detection of gender from face images," 2019.