

DAGER: Deep Age, Gender and Emotion Recognition Using Convolutional Neural Networks

Afshin Dehghan Enrique G. Ortiz Guang Shu Syed Zain Masood
{afshindehghan, egortiz, guangshu, zainmasood}@sighthound.com

Computer Vision Lab, Sighthound Inc., Winter Park, FL

Abstract. This paper describes the details of Sighthound’s fully automated age, gender and emotion recognition system. The backbone of our system consists of several deep convolutional neural networks that are not only computationally inexpensive, but also provide state-of-the-art results on several competitive benchmarks. To power our novel deep networks, we collected large labeled datasets through a semi-supervised pipeline to reduce the annotation effort/time. We tested our system on several public benchmarks and report outstanding results. Our age, gender and emotion recognition models are available to developers through the Sighthound Cloud API at <https://www.sighthound.com/products/cloud>

1 Introduction

Facial attribute recognition, including age, gender and emotion, [1,2,3,4,5,6,7] has been a topic of interest among computer vision researchers for over a decade. One of the key reasons is the numerous applications of this challenging problem which range from security control, to person identification, to human-computer interaction. Due to the release of large labeled datasets, as well as the advances made in the design of convolutional neural networks, error rates have dropped significantly. In many cases, these systems are able to outperform humans [5]. However, this still remains a difficult problem and existing commercial systems fall short when dealing with real world scenarios. In this work, we present an end-to-end system capable of estimating facial attributes including age, gender and emotion with low error rates. In order to support our claims, we tested our system on several benchmarks and achieved results better than the previous state-of-the-art. The contributions of this work are summarized below.

- We present an end-to-end pipeline, along with novel deep networks, that not only are computationally inexpensive, but also outperform competitive methods on several benchmarks.
- We present large datasets for age, emotion and gender recognition that are used to train state-of-the-art deep neural networks.
- We conducted a number of experiments on existing benchmarks and obtained leading results on all of them.

2 System Overview

The pipeline of our system is shown in Figure 1. Our first deep model is trained on a large dataset of four million images for the task of face recognition. This model serves as the backbone to our facial attribute recognizers and is used to fine-tune networks for four tasks: real age estimation, apparent age estimation, gender recognition and emotion recognition. What follows explains each of the steps in more detail.

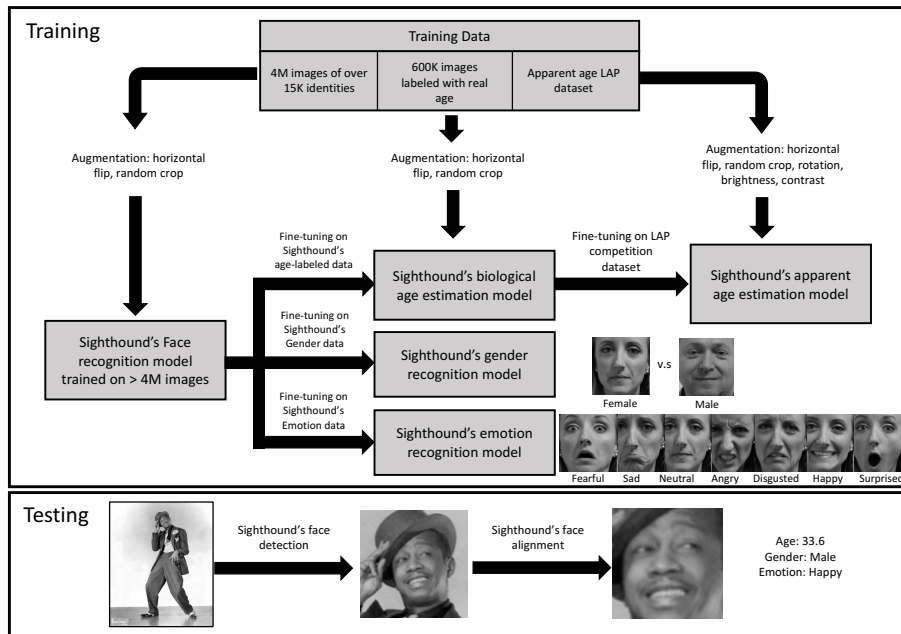


Fig. 1: This figure shows the pipeline of our system. Images are collected from different sources and labeled for different tasks. Over 4 million images of more than 40,000 people are collected for the task of facial recognition. All images are labeled with their corresponding gender label and part of the data is annotated with emotion. These images are later pruned using a semi-automated process with a team of human annotators in the loop. The images are pre-processed next to extract the faces and align them. The aligned images are then fed to our proprietary deep network for training.

2.1 Training

Below we describe different steps involved in training our models in more detail.

- **Data collection:** Data collection plays an important role in training any deep neural network (DNN). In this paper, we aim to label data for three

separate tasks: age, gender and emotion recognition. Collecting labeled data for some tasks, such as real age estimation, is much more challenging compared to popular classification [8,9] or detection [10] problems. This disparity is due to the fact that human error in estimating real age is large (sometimes greater than the computer vision estimations) and one cannot rely on human annotators to label faces with their corresponding real age. However, at Sighthound we have collected a large dataset of faces with their corresponding age, gender and emotion labels. To our knowledge, our datasets are the largest or among the largest in either the academic or commercial world. Below we provide some statistics on the data used for training our models.

Face recognition: The base model for our facial recognition is trained on over four million images of more than 40,000 individuals. The large variation in images of each identity make our deep model robust to common challenges in face recognition. Our face recognition model is available to developers through the Sighthound Cloud API ¹.

Age estimation: Recently there have been some efforts in collecting data with corresponding age labels [5,11,4]. Among those, the dataset proposed by Rothe et al. in [4] is the largest dataset that contains 523,051 images and is available for research purposes. However, the dataset is not carefully annotated and contains many mistakes. Additionally the distribution of the data across different ages is highly unbalanced. This led to the authors using only half of the data for training in the original paper [4]. To better address this problem we collected a large dataset of $\sim 600,000$ images with corresponding age labels. In contrast to previous works, our dataset has a more balanced distribution across different ages. For example we have over 120,000 people in our dataset with labeled ages over 70 or younger than 20 years of age. We used a team of human annotators to further clean our dataset through a semi-supervised procedure.

Gender and emotion recognition: Our four million faces labeled for the task of face recognition are also labeled with their corresponding gender. To better improve our model, we added tens of thousands of images of different ethnicities as well as age groups. Additionally, we also annotated part of our data with emotion labels for the task of emotion recognition.

- **Data pre-processing:** We pre-process each image before feeding them to our DNNs. These pre-processing steps include face detection, facial landmark detection and alignment. We used Sighthound's face detection which is available through our cloud API. If more than one face is detected in an image, we choose the most centered one. (This is especially the case in the ChaLearn v2 dataset which contains multiple faces). In the Chalearn dataset we were able to detect all faces using a combination of techniques, but all using the Sighthound Cloud APIs.² Given the face bounding boxes, we detect

¹ <https://www.sighthound.com/products/cloud>

² Not all of the face detection bounding boxes generated by the cloud API were perfectly accurate, mostly due to occlusion or low resolution of some images. However,

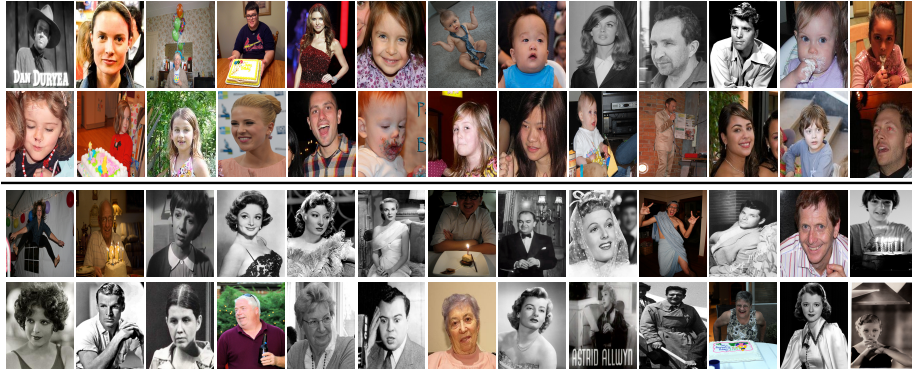


Fig. 2: Qualitative results of our method on ChaLearn LAP challenge. The top half of the figure shows some sample images where our absolute error is less than half a year and the bottom half shows images in which our error was more than 5 years. As shown, the performance of our method drops mostly for gray-scale/old-style images where our network tends to over estimate the age for those images.

68 facial landmarks and use those for alignment. Finally the aligned faces are all cropped and resized to a fixed size. In contrast to some previous works, which do not use any face alignment [4], we found this to be important in our final accuracy numbers.

- Deep training:** As shown in Figure 1, we start by training a deep neural network for the task of face recognition using four million images of over 40,000 identities. Our face recognition model is not only computationally inexpensive (with feature extraction time of 70ms using just the CPU), but also achieves outstanding results on the LFW dataset [2]. This model serves as the backbone of our facial attribute recognition engine. We designed a highly optimized deep network architecture for accuracy and speed for each task. In some recent works [12], researchers try to design a network which performs all tasks at the same time, and they have shown marginal improvements. However, having separate networks for each task allowed us to design faster and more portable models for each task. Additionally running all models combined takes less time compared to the all-in-one model of [12] and we achieve better results. We should add that even though the network architecture is not the same for each task, all networks are trained first for the task of facial recognition using the four million image set.

when comparing with other methods on public datasets, we directly used the output of our face detection for age estimation without further adjustment.

3 Experiments

In this section, we report experimental results on several publicly available datasets as well as our internal datasets.

3.1 Real Age Estimation

For real age estimation, we report results on two publicly available datasets, the Group dataset [5] and the Adience dataset [11]. The images in these datasets are labeled with their corresponding age groups. In order to further evaluate our system on estimating the actual age and not only the age group, we collected an internal dataset of 3,800 images, on which we report the results of the proposed method in addition to other competitive algorithms.

Methods	MAE
Sighthound	5.76
Rothe et al. [4]	7.34
Microsoft. [13]	7.62
Kairos [14]	10.57
Face++ [15]	11.04

Table 1: This table shows the mean absolute error for our methods along with competitive approaches on our Sighthound dataset. As can be seen, our method outperforms the second best method by 1.58 years.

In Table 1, we present quantitative results on our Sighthound dataset, containing 3,800 test images. Each image is labeled with its corresponding age ranging from 10 to 90 years old. Unlike the Adience and Group datasets, our dataset includes the exact age labels for each image and not the age groups. We compare our results with [4], whose model is trained on the IMDB-Wiki dataset. Additionally, we compare our system with available commercial APIs: Microsoft [13], Face++ [15] and Kairos [14]. For quantitative comparison, we used the Mean Absolute Error (MAE) which is commonly used in the literature [4]. As can be seen, our method outperforms competitive approaches by a significant margin.³

Next, we provide results on the Group dataset, which contains 28,231 faces collected from Flickr and labeled with one of seven age categories roughly correspond to different life stages. Most faces are low-resolution making it more challenging for accurate age estimation. The median of faces are reported to have only 18.5 pixels between the eye centers. We followed the setup in [5] where

³ We compute the error rate for Microsoft and Face++ using the versions of their cloud API available in October 2016.

3,500 images are used for training and 1,050 images are used for testing. Both training and testing images are equally distributed across seven age groups. The age classification results in terms of accuracy along with a confusion table are reported in Table 2 and Figure 3. We can observe that Sighthound’s age estimation outperforms the latest research results.

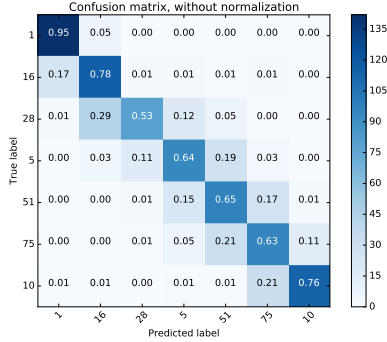


Fig. 3: Confusion table for the Group dataset.

Methods	top1	1-off
Sighthound	70.5%	96.2%
Hou et al. [16]	65.0%	96.1%
Rothe et al. [4]	62.3%	94.3%
Dong et al. [17]	56.0%	92.0%
Gallegher et al. [5]	42.9%	78.1%

Table 2: Age classification accuracy of the Group dataset. We report both the exact accuracy as well as the 1-off accuracy.

Finally we report results on the Adience benchmark. The entire Adience benchmark includes roughly 26,000 images of 2,284 subjects. However, some images are not annotated with corresponding age groups. Therefore the total number of images used for final testing is smaller than 26K. We used the standard 5-fold cross validation experiment defined for this set. When testing on each fold, the rest of the folds are used to fine-tune our model for the eight age groups defined in the dataset. The results of our method along with competitive approaches are shown in Table 3. Once again, our method improves on the best reported results on this dataset.

Methods	Accuracy(Top-1)
Sighthound	61.3 ± 3.7%
Hou et. al. [16]	61.1 ± NR%
Eidinger et. al. [11]	45.1 ± 2.6%
Levi and Hassner [18]	50.7 ± 5.1%

Table 3: Real age estimation accuracy of the Adience benchmark. Sighthound outperforms other methods. (NR=Not Reported)

3.2 Apparent Age Estimation

The apparent age of a person could be very different from the real age of a person. Recently, thanks to the availability of the Chalearn LAP Apparent Age Estimation dataset and challenge [19], several researchers have focused on designing models that are focused on predicting the apparent age, rather than the

actual age. In the most recent version of the competition, the size of the dataset was extended to 7,591 images where 4,113 of them are used for training and 1,500 and 1,978 are used for validation and testing respectively. Each image in the dataset is annotated using at least 10 human votes and the mean (μ) and standard deviation (σ) of the votes are recorded and released with the dataset. Given the prediction for each image (x), the error for each image is computed using $\epsilon = 1 - \exp -\frac{(\hat{x}-\mu)^2}{2\sigma^2}$. This means the apparent age on an image with small standard deviation gets penalized more compared to one with larger standard deviation. The winner of the competition [6] used a multi CNN framework and achieved a test error of 0.2411. However, this method [6] has several limitations. Their network architecture is an order of magnitude slower in speed and an order of magnitude larger in size compared to ours. Additionally, the multiple CNNs (minimum of 88 forward passes for each image) in their pipeline makes it impossible to use their system in a real-time or even close-to real-time application. We should also mention that all the runner up approaches suffer from the same limitations. Even though our goal is to keep the computational complexity low, we still achieve the outstanding error rate of 0.319, which places our approach second.

Methods	Test Error	Score-level fusion
Sighthound	0.319	No
OrangeLabs [6]	0.2411	Yes
Palm-seu [20]	0.3214	Yes
CMP+ETH [21]	0.3361	Yes
WYU-CVL	0.3405	No
ITU-SiMiT	0.3668	Yes
Bogazici	0.3740	Yes
MIPAL-SNU	0.4565	Yes

Table 4: Results for ChaLearn [19] apparent age estimation 2016 challenge. Our fine-tuned system achieves a test error of 0.319 and obtains the second best place. One should note that our model uses only a single CNN, which is not the case for most top performing teams. Additionally, our base model is almost an order of magnitude faster than the base CNN model of top performing teams (VGG).

3.3 Emotion Recognition

There are several public datasets for emotion recognition. FER-2013 [22] and EmotiW are among the popular ones. The FER dataset contains low-quality gray scale images of size 48×48 which is not very representative of real world scenarios. Access to the EmotiW dataset was not granted to us. Thus we collected our own dataset of 2,156 images. Each image is labeled with one of the 7 labels of "happy", "sad", "neutral", "disgusted", "surprised", "fearful" and "angry". The data has a relatively equal distribution across the 7 emotions. We compared our method with the Microsoft Face API [13]. The results as well as the confusion

tables are shown in Table 5 and Figure 4 respectively. As shown, Sighthound’s emotion recognition system outperforms Microsoft by a 15% margin.⁴

Methods	Accuracy
Sighthound	76.1%
Microsoft [13]	61.3%

Table 5: Emotion recognition accuracy on Sighthound dataset.

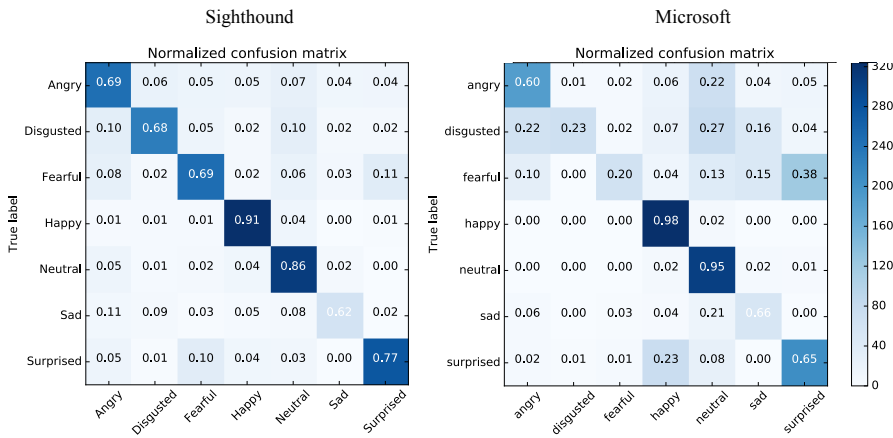


Fig. 4: Quantitative results in the form of a confusion matrix on Sighthound’s emotion recognition dataset. On the left we show the results of our system and on the right we show the results for the Microsoft Face API. Our method performs well on almost all emotions, while Microsoft’s performance drops significantly for emotions other than "happy" and "neutral".

3.4 Gender Recognition

We compared our gender recognition model on the Adience benchmark with other leading methods. The Adience benchmark contains 17,492 faces labeled with their corresponding gender. The faces are divided into 5 folds. However, we used the same model across all folds without further fine-tuning. Along with published state-of-the-art results, we compare our method with a couple of commercial APIs such as [15] and [14]. The results reported in Table 6 clearly

⁴ Microsoft’s API failed detecting a face in 193 images. To be fair to Microsoft we removed these images while evaluating their method.

show Sighthound’s enhanced gender recognition capability compared to recent research publications and commercial products.

Methods	Accuracy
Sighthound	91.00%
Microsoft. [13]	90.86%
Rothe et al. [4]	88.75%
Levi and Hassner [18]	86.80%
Kairos [14]	84.66%
Face++ [15]	83.04%

Table 6: Results for gender recognition on the Adience benchmark [11]. We compare our method against state-of-the-art research and commercial entities.

4 Conclusions

In this paper, we present an end to end system for age, gender and emotion recognition. We show that our novel deep architecture, along with our large, in-house collected data, can outperform competitive commercial and academic algorithms on several benchmarks.

References

1. Busso, Carlos, e.a.: Analysis of emotion recognition using facial expressions, speech and multimodal information. In: Proceedings of the 6th international conference on Multimodal interfaces. (2004)
2. Levi, G., Hassner., T.: Emotion recognition in the wild via convolutional neural networks and mapped binary patterns. In: Proceedings of the 2015 ACM on International Conference on Multimodal Interaction. ACM. (2015)
3. Pang, L., Ngo., C.W.: Mutlimodal learning with deep boltzmann machine for emotion prediction in user generated videos. In: Proceedings of the 2015 ACM on International Conference on Multimodal Retrieval. ACM. (2015)
4. Rothe, R., Radu Timofte, L.V.G.: Dex: Deep expectation of apparent age from a single image. In: International Conference on Computer Vision (ICCV),. (2015)
5. Gallagher, A.C., Chen., T.: Understanding images of groups of people. In: CVPR. (2009)
6. Antipov, Grigory, e.a.: Apparent age estimation from face images combining general and children-specialized deep learning models. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops. (2016)
7. Wang, X., Guo, R., Kambhamettu, C.: Deeply-learned feature for age estimation. In: WACV. (2015)
8. Dehghan, A., Masood, S.Z., Shu, G., Ortiz., E.G.: View independent vehicle make, model and color recognition using convolutional neural network. In: arXiv:1702.01721. (2017)

9. Deng, Jia, e.a.: Imagenet: A large-scale hierarchical image database. In: Computer Vision and Pattern Recognition. (2009)
10. Everingham, M., E.S.M.A.V.G.L.W.C.K.I.W.J., Zisserman, A.: The pascal visual object classes challenge: A retrospective . In: International Journal of Computer Vision. (2015)
11. Eiding, E., Enbar, R., Hassner., T.: Age and gender estimation of unfiltered faces. In: IEEE TRANSACTIONS ON INFORMATION FORENSICS AND SECURITY. (2013)
12. Ranjan, R., Sankaranarayanan, S., Castillo, C.D., Chellappa, R.: An all-in-one convolutional neural network for face analysis. In: arXiv:1611.00851. (2016)
13. Microsoft-Face-API.: (<https://www.microsoft.com/cognitive-services/en-us/face-api>.)
14. Kairos.: (<https://www.kairos.com/kairos-2.0/demos>)
15. Face++.: (<http://old.faceplusplus.com/demo-detect/>)
16. Hou, L., Yu, C.P., Samaras., D.: Squared earth mover's distance-based loss for training deep neural networks. In: arXiv. (2016)
17. Yuan, D., Liu, Y., Lian., S.: Automatic age estimation based on deep learning algorithm. In: Neurocomputing. (2016)
18. Levi, G., Hassner., T.: Age and gender classification using convolutional neural networks. In: CVPRW. (2015)
19. Escalera, Sergio, e.a.: Chalearn looking at people and faces of the world: Face analysis workshop and challenge 2016. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops. (2016)
20. Huo, Z., Yang, X., Xing, C., Zhou, Y., Hou, P., Lv, J., Geng, X.: Deep age distribution learning for apparent age estimation.. In: IEEE Conference on Computer Vision and Pattern Recognition Workshops. (2016)
21. Uříčář, M., Timofte, R., Rothe, R., Matas, J., Gool, L.V.: Structured output SVM prediction of apparent age, gender and smile from deep features. In: Proceedings of IEEE conference on Computer Vision and Pattern Recognition Workshops, Las Vegas, USA (2016)
22. Goodfellow, I.J.: Challenges in representation learning: A report on three machine learning contests. In: International Conference on Neural Information Processing. (2013)