

# AGE RANGE ESTIMATION USING MTCNN AND VGG-FACE MODEL

Dipesh Gyawali<sup>1</sup>, Prashanga Pokharel<sup>2</sup>, Ashutosh Chauhan<sup>3</sup>, Subodh Chandra Shakya<sup>4</sup>

InfoDevelopers Pvt. Ltd.

Lalitpur, Nepal

(<sup>1</sup>dipesh9393, <sup>2</sup>pokharel997, <sup>3</sup>Cashutosh6711, <sup>4</sup>subshakya591)@gmail.com

**Abstract**—The Convolutional Neural Network has amazed us with its usage on several applications. Age range estimation using CNN is emerging due to its application in myriad of areas which makes it a state-of-the-art area for research and improve the estimation accuracy. A deep CNN model is used for identification of people's age range in our proposed work. At first, we extracted only face images from image dataset using MTCNN to remove unnecessary features other than face from the image. Secondly, we used random crop technique for data augmentation to improve the model performance. We have used the concept of transfer learning in our research. A pretrained face recognition model i.e VGG-Face is used to build our model for identification of age range whose performance is evaluated on Adience Benchmark for confirming the efficacy of our work. The performance in test set outperformed existing state-of-the-art by substantial margins.

**Index Terms**—Age Range Estimation, CNN, MTCNN, Transfer Learning, VGG-Face

## I. INTRODUCTION

The face images are used for different applications recently in face recognition[1], surveillance system[15], emotion identification[16], and smart attendance for multifarious purposes. Identifying the age range using CNN has become a challenging task as the age range among the adults seems to be similar than the child and aged people[8]. The facial features among the child and old people vary drastically from the adult people which makes it easier to identify the age of child and old people[7]. But the identification among the adult age becomes cumbersome due to the presence of almost the same facial features.

By far, CNN has become most effective for processing image data in object detection, face recognition, and image classification. Image data consists of a lot of parameters which are difficult for a basic neural network to process. A lot of matrix calculation needs to be handled for 2D image data where CNN is a most effective neural network as different convolutional layers in CNN are used to extract the feature map of images and the convolutional process with the filter make it easy to process image data.

In our research, we have identified the age range among various ages of people using novel approach. Transfer learning has been considered more effective than building our own custom CNN model for a limited number of dataset. The model which was pre-trained for image classification and face recognition[3] is used for estimating the age range.

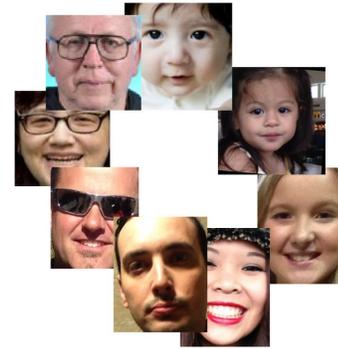


Fig. 1. Face images of different age groups on Adience Benchmark[11]

## II. RELATED WORKS

A dataset including one million faces of celebrities are prepared which eventually improved face recognition accuracy[1]. Face recognition using new dataset across various identities minimizing the label noise[2], using either a single image or set of faces[3], analysis and survey[4] is performed. Age Estimation is done using VGG-Face model[5], the ranking SVM[6] and finding sufficient embedding space by applying manifold learning methods which models data with multiple linear regression function[7] effectively. A craniofacial growth model is proposed which models growth related variations in shape of human faces[8]. Different faces are divided into several minute regions for the extraction of Local Binary Pattern(LBP) to use as a face descriptor by concatenating into feature vectors[9]. Two new approaches, Ranking-CNN[10] based on rank relationship and deep CNN architecture for low quality face images[24] are proposed for identifying age from different perspectives. Using unfiltered faces on Adience Benchmark[11] and in the presence of limited learning data using deep CNN[12], identification of gender and age is performed. Various research challenges in addition with recent survey are delineated for research in face recognition along with age estimation[13]. The effects of aging on performance of age invariant face recognition are surveyed[14] and analyzed using deep features analysis[21].

The above works showed that a handful of research tasks have already been carried out using different metrics and methods. Although a plethora of works have been done, a significant improvement in the identification of age range across adult age group due to similar facial features had

not been achieved. In this experiment, we have initially extracted only the facial features from joint face alignment using MTCNN network and later identified age range using VGG-Face model which significantly reduced overfitting and improved identification among adult age group as compared to previous works. Also, a comparative analysis has been done across multiple face models to analyze efficiency of our work using unconstrained face images.

### III. DATA PREPARATION

MTCNN (Multi-task Cascaded Convolutional Networks) is a network[17] which is used in our proposed work to extract the facial features. MTCNN contains 3 CNN steps, each step called as Proposal Network (P-net), Refine Network (R-net) and Output Network (O-net). All the input images are fed to Proposal Network (P-Net), which is a CNN model. The candidate window inside the image with their bounding box regression vector is received as output from P-net. We create an image pyramid, in order to detect faces of all different sizes as shown in Fig. 2.

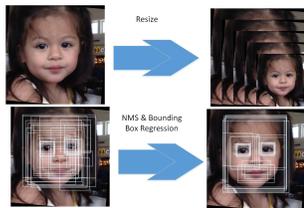


Fig. 2. Image pyramids and NMS and bounding box regression (P-Net)

After feeding those candidates to Refined Network (R-Net), we have identified regression vector of bounding box and used Non-Max Suppression (NMS) to find and integrate densely overlapped candidates as shown in Fig. 3. The outcomes from this network are passed to O-Net.



Fig. 3. NMS and bounding box regression (R-Net)

Finally, Output Network (O-Net) provides three output, coordinate of bounding box, the coordinate of five facial landmarks, and the confidence level of each box as shown in Fig. 4. This bounding box image is saved as a new image which is then passed into the proposed model.



Fig. 4. Output-network

Images in the dataset are rescaled to 256x256 pixels. These images are extracted of size 224x224 pixels using random crop[18]. Horizontal flip and rotation is applied to the extracted images for data augmentation purpose. All age groups were divided into 8 classes. For each class, images were splitted 80-20 for training and validation. Finally this dataset is fed to the network.

### IV. PROPOSED SYSTEM

A large number of images are required for building our own custom model to ameliorate its performance. In order to prevent overfitting, we have used transfer learning for building our CNN network. The model is well validated and tested after optimizing it using appropriate hyperparameters. The VGG-Face model performs well on Adience Benchmark[11] due to its less number of CNN layers which can be useful for small dataset.

#### A. Architecture

For the purpose of age estimation, the pre-trained model used for face recognition task[3] is considered extremely useful. The age estimation is generally performed by observing the facial features of different age groups[7]. Due of this, for the estimation of age robustly, model previously trained on face recognition can be used. By the use of pre-trained models, there will be less chance of overfitting. Also, we have used MTCNN for extracting the facial features while training in our dataset.

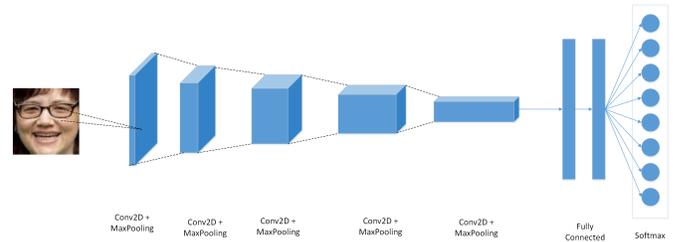


Fig. 5. CNN Architecture

To identify a better age range, we have used pre-trained model VGG-Face. The model performed pretty well using VGG-Face due to its less number of CNN layers and appropriate tuning of hyperparameters. Eight convolutional layers comprising of Conv2D & Pooling and three fully connected layers makes up a overall of eleven layers of VGG-Face model. All the layers except fully connected softmax layers were used for training purposes. The layers from VGG-Face are kept frozen first and then additional layers are added to the VGG-Face model as per our requirements. Two fully connected layers are added in corresponding with one dropout layer and one softmax output layer to identify the age range for 8 classes. The layers and the respective dimensions used in our network are represented as in the Table I.

TABLE I  
OUR PROPOSED CNN ARCHITECTURE USING VGG-FACE MODEL

Layer	0	1	2	3	4	5	6	7	8	9	10	11	12
type	input	conv	relu	conv	relu	mpool	conv	relu	conv	relu	mpool	conv	relu
name	-	conv1_1	relu1_1	conv1_2	relu1_2	pool1	conv2_1	relu2_1	conv2_2	relu2_2	pool2	conv3_1	relu3_1
support	-	3	1	3	1	2	3	1	3	1	2	3	1
filt dim	-	3	-	64	-	-	64	-	128	-	-	128	-
num filts	-	64	-	64	-	-	128	-	128	-	-	256	-
stride	-	1	1	1	1	2	1	1	1	1	2	1	1
pad	-	1	0	1	0	0	1	0	1	0	0	1	0
Layer	13	14	15	16	17	18	19	20	21	22	23	24	25
type	conv	input	conv	relu	mpool	conv	relu	conv	relu	conv	relu	mpool	conv
name	conv3_2	relu3_2	conv3_3	relu3_3	pool3	conv4_1	relu4_1	conv4_2	relu4_2	conv4_3	relu_3	pool4	conv5_1
support	3	1	3	1	2	3	1	3	1	3	1	2	3
filt dim	256	-	256	-	-	256	-	512	-	512	-	-	512
num filts	256	-	256	-	-	512	-	512	-	512	-	-	512
stride	1	1	1	1	2	1	1	1	1	1	1	2	1
pad	1	0	1	0	0	1	0	1	0	1	0	0	1
Layer	26	27	28	29	30	31	32	33	34	35	36	37	
type	relu	conv	relu	conv	relu	mpool	conv	relu	dropout	conv	relu	softmax	
name	relu5_1	conv5_2	relu5_2	conv5_3	relu5_3	pool5	fc6	relu6	dropout7	fc8	relu7	prob	
support	1	3	1	3	1	2	7	1	1	1	1	1	
filt dim	-	512	-	512	-	-	512	-	-	1000	-	-	
num filts	-	512	-	512	-	-	1000	-	-	100	-	-	
stride	1	1	1	1	1	2	1	1	1	1	1	1	
pad	0	1	0	1	0	0	0	0	0	0	0	0	

### B. Loss Function

We have multiple classes of age-range ranging from 0-2 to 60-100 with a total of eight classes. The actual and predicted probability distributions for all eight classes are identified whose average differences are summarized with a score calculated using cross entropy[19]. A total of eight nodes are present at the output layer for each age range in which 'softmax' activation is used to predict the probability across each classes. Due to a multiclass classification problem and effectiveness, we have used categorical cross entropy as a loss function[19].

$$CE = - \sum_{c=1}^M t_{o,c} \log(P_{o,c}) \quad (1)$$

where t is either 0 or 1 if label class c is the correct classification for observation o. M is a total number of classes. P is predicted probability of class c.

### C. Optimizer

Due to efficient computation, little memory requirements and presence of intuitive interpretations of hyperparameters, we have used Adam optimizer[20] in our research which is considered more effective comparatively with other optimizers.

### D. Training

We used transfer learning(VGG-Face) to train our model. The convolutional layer parameters of VGG-Face are not changed and kept frozen. We optimize the fully connected layers parameters by adding additional two fully connected layers, one dropout layer and a final output layer to the VGG-Face model. Number of filters in the first of two fully connected layers are 1000 and 100 consecutively. In these fully connected layers we use 'relu' activation function. We use 0.3

(30% chance of setting a neuron's output value to zero) values for one dropout layer. 'Softmax' activation is used in the final layer to classify the age range. 'Adam optimizer' is used as an optimizer with 'categorical cross entropy' as a loss function in our model.

### E. Prediction

After training the model, we have used the Adience Benchmark to evaluate the model performance. A test image is rescaled to 256x256 pixels. Then five images are extracted of size 224x224 pixels. The four images are obtained from four corners of the image and the final image is obtained from the center of the original test image as shown in Fig. 6. These five images are fed into our trained neural network to calculate the softmax probability output vector of each of the five images. These output scores vector of these five images were averaged. Average result helps to reduce the impact of poor quality and low-resolution images.



Fig. 6. Image Rescaling to 256x256 pixels initially. Five images are extracted from each four corners and centre which are resized to 224x224 for feeding into the network.

## V. EXPERIMENTS AND RESULTS

A wide known and flexible deep learning library i.e Keras consisting of Tensorflow as a backend engine is used in our proposed work for effectiveness. Training was done on NVIDIA GTX 2080 GPU with 4352 CUDA cores. Our

TABLE II  
THE ADIENCIE BENCHMARK

	0-2	4-6	8-13	15-20	25-32	38-43	48-53	60+	Total
Male	745	928	934	734	2308	1294	392	442	8192
Female	682	1234	1360	919	2589	1056	433	427	9411
Both	1427	2162	2294	1653	4897	2350	825	869	19487

TABLE III  
AGE ESTIMATION RESULT ON THE ADIENCIE BENCHMARK

Method	Exact Accuracy	1-off Accuracy
[11]	45.1	79.5
[12] using single crop	49.5	84.6
[12] using over-sample	50.7	84.7
[5]	59.9	90.57
Proposed work	<b>70.96</b>	<b>92.7</b>

training time took approximately four hours. We trained VGG-Face Model for 50 epochs with a batch size of 64. The model performance is evaluated using Adience benchmark. The Adience benchmark is a large-scale face dataset which consists of 19487 multiple face images with 8 classes of age range covering large variation in facial expression, pose, occlusion, resolution and illumination. The exact accuracy on Adience benchmark came out to be 70.96%. We have used this Adience benchmark for appraising the efficacy of our work. Different age range's label used in our research with the corresponding number of images per label on Adience Benchmark are shown in Table II.

Consequently, our proposed model outperformed the previous methods reported in [5], [11], and [12]. The accuracy was compared with the previous work done for age estimation in which our model performs better as compared to other performances. Table III shows exact accuracy and 1-off accuracy from the previous work and our proposed work using VGG-Face architecture. Our proposed work outperforms the previous work and confirm the efficiency of our proposed work. The use of MTCNN approach has significantly improved accuracy on Adience Benchmark. Fig. 7 shows the confusion matrix on Adience Benchmark.

The true label and the predicted label are denoted from indices of each rows and each columns respectively in the confusion matrix. The occurrences of prediction is given by number showing on each cell. In the leading diagonal of confusion matrix, the true label and predicted label are equal and other off-diagonal elements represent occurrence mislabeled by classifier. There are more correct predictions if there are higher values in the leading diagonal. The confusion matrix in Fig. 7. clearly shows probability of prediction. Our result from the confusion matrix depicts that the highest accuracy is 98.90% which of (0-2) age range as shown in Fig. 7. The highest accuracy is due to the distinctive features which enable the classifier to distinguish easily. The images of middle age groups consist of almost similar features which resulted in less accuracy as compared to (0-2) age range and (60+) age range. In our work, the age group of (48-53) and (25-32) is highly misclassified with 41.67% and 47.30%. The

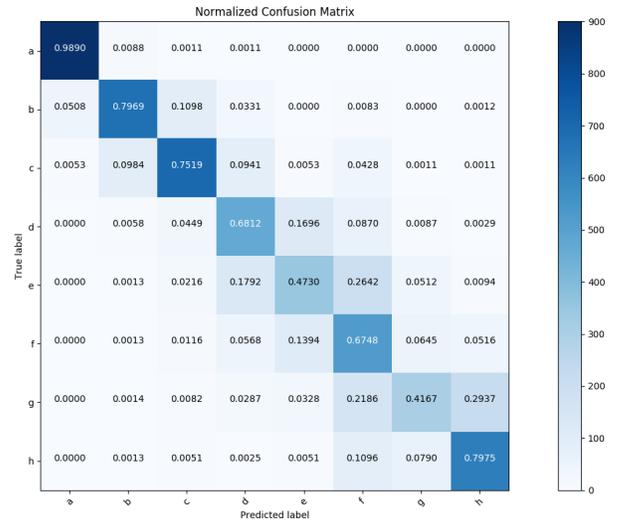


Fig. 7. Normalized Confusion Matrix for Adience Benchmark. The label is marked from a-h to make ease while building confusion matrix. The label is defined as a(0-2),b(4-6),c(8-13),d(15-20),e(25-32),f(38-43),g(48-53),h(60+).

TABLE IV  
CLASSIFICATION REPORT ON ADIENCIE BENCHMARK

Age Range From Adience Dataset	Precision	Recall	F1-Score	Support
0-2	0.95	0.99	0.97	1427
4-6	0.86	0.80	0.83	2162
8-13	0.81	0.75	0.78	2294
15-20	0.60	0.68	0.64	1653
25-32	0.58	0.47	0.52	4897
38-43	0.49	0.67	0.57	2350
48-53	0.66	0.42	0.51	825
60+	0.70	0.80	0.75	869
Accuracy			0.71	19487
macro avg	0.71	0.70	0.69	19487
weighted avg	0.72	0.71	0.71	19487

classification report on Adience Benchmark using VGG-Face model is shown in the Table IV.

We have also used VGG-Face2 Model to analyze the efficacy of our task. Our proposed work using VGG-Face model proved to be far better than VGG-Face2 model which is shown in the Table V. There are more numbers of CNN layers on VGG-Face2 as compared to VGG-Face. As we had not enough face image data, the VGG-Face2 model was subjected to overfitting and could not perform well so that the performance across VGG-Face2 model seems very poor as compared to VGG-Face model.

Fig. 8 and Fig. 9 shows some of correctly classified and misclassified images on Adience Benchmark with their corresponding confidence level. The misclassified classes are mainly of adults due to same facial features. The age range of child and old age were almost classified correctly. The number of misclassified class could be improved further if we have high number of images for training. Different models can be used if there is copious amount of data in addition with proper pre-training process which might improve our proposed work.

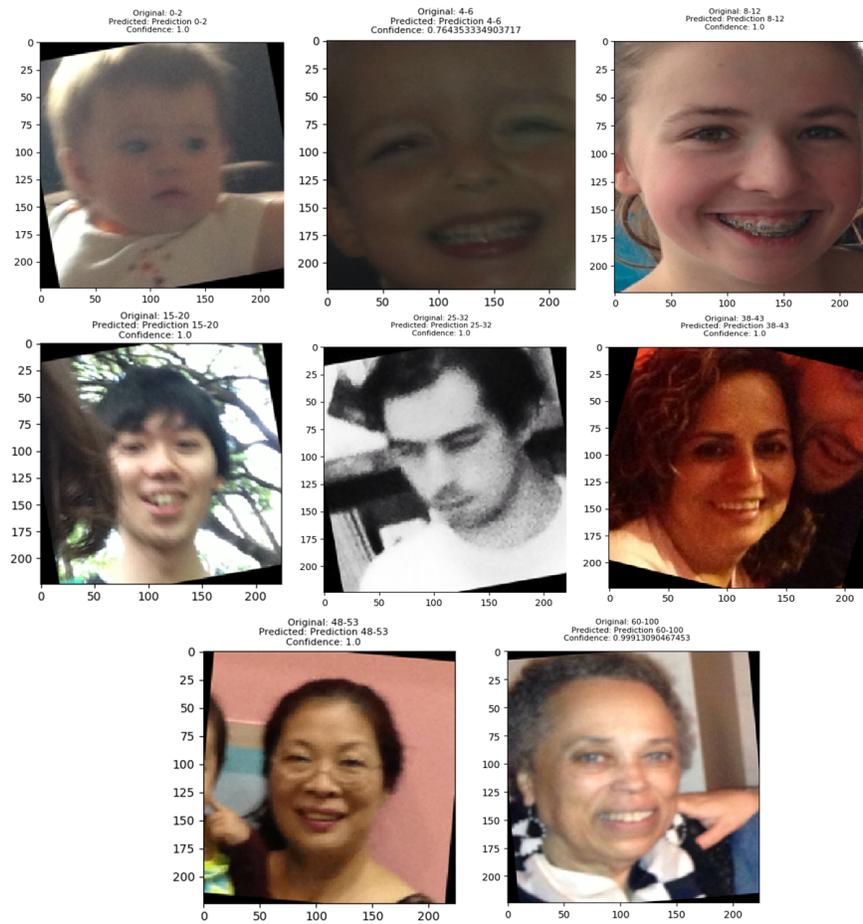


Fig. 8. Some of the correctly classified images on Adience Benchmark

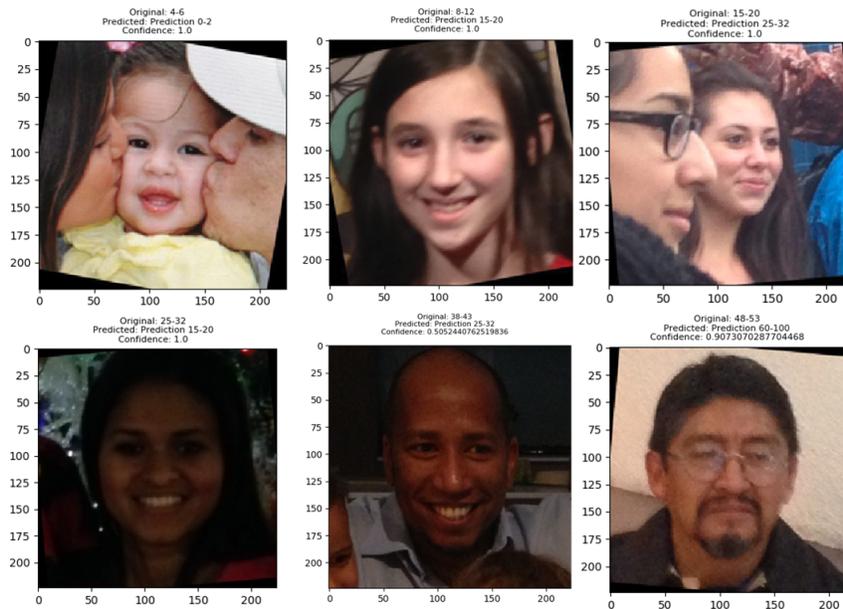


Fig. 9. Some of the misclassified images on Adience Benchmark

TABLE V  
EXACT ACCURACY ON DIFFERENT FACE MODELS

Label	Proposed model using VGG-Face	Model using VGG-Face2
0-2	98.90	80.40
4-6	79.69	36.52
8-13	75.19	29.64
15-20	68.12	52.92
25-32	47.30	10.37
38-43	67.48	13.08
48-53	41.67	49.89
60+	79.75	77.86
<b>Overall Exact Accuracy</b>	<b>70.96</b>	<b>42.02</b>

## VI. CONCLUSION

We have used the VGG-Face model and MTCNN to extract the facial features for estimating the age range for our research work. MTCNN helped to extract only the facial features from the image data which helps to determine the most germane features from the face. Our research work outperforms the previous work by almost 12% on the Adience Benchmark. Due to the small number of dataset for training and large number of layers for feature extraction, the VGG-Face2 model is subjected to overfitting which does not perform well while evaluating its performance. Usage of MTCNN and fine tuning the VGG-Face model significantly improved the network's performance for age range estimation.

Furthermore, there are still some misclassified images with high confidence because of the same facial features, lightness, occlusions and multiple person in an image. If the number of images in various conditions are increased and pre-training tasks of CNN network are performed significantly, the model performance can increase further. Also, VGG-Face2 model performance can be improved by using high number of face images and fine-tuning the network. The exact accuracy is still low due to the small number of dataset available for age estimation which could be improved in future works.

## ACKNOWLEDGMENT

Our research is based upon work done at InfoDevelopers Pvt. Ltd. We would like to express gratitude towards our supervisor Mr. Sandesh Pandey for supervising and evaluating our research work. We are also thankful to Info Developers Private Limited for providing us the research platform and materials regarding our research procedures.

## REFERENCES

- [1] Guo Y., Zhang L., Hu Y., He X., Gao J. (2016) MS-Celeb-1M: A Dataset and Benchmark for Large-Scale Face Recognition. In: Leibe B., Matas J., Sebe N., Welling M. (eds) Computer Vision – ECCV 2016.
- [2] Cao Q., Shen L., Xie W., Parkhi O.M., Zisserman A., "VGGFace2: A dataset for recognising faces across pose and age", International Conference on Automatic Face and Gesture Recognition, 2018.
- [3] O.M Parkhi, A. Vedaldi, A. Zisserman, "Deep Face Recognition", British Machine Vision Conference, 2015.
- [4] S. Kumar, S. Singh and J. Kumar, "A study on face recognition techniques with age and gender classification," 2017 International Conference on Computing, Communication and Automation (ICCCA), Greater Noida, 2017, pp. 1001-1006, doi: 10.1109/CCAA.2017.8229960.

- [5] Z. Qawaqneh, A. A. Mallouh, B. D. Barkana, "Deep Convolutional Neural Network for Age Estimation based on VGG-Face Model", arXiv preprint 1709.01664 (2017).
- [6] Cao D., Lei Z., Zhang Z., Feng J., Li S.Z. (2012) Human Age Estimation using Ranking SVM. In: Zheng W.S., Sun Z., Wang Y., Chen X., Yuen P.C., Lai J. (eds) Biometric Recognition. CCBR 2012. Springer, Berlin, Heidelberg.
- [7] Y. Fu, Y. Xu and T. S. Huang, "Estimating Human Age by Manifold Analysis of Face Pictures and Regression on Aging Features," 2007 IEEE International Conference on Multimedia and Expo, Beijing, 2007, pp. 1383-1386, doi: 10.1109/ICME.2007.4284917.
- [8] N. Ramanathan and R. Chellappa, "Modeling Age Progression in Young Faces," 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06), New York, NY, USA, 2006, pp. 387-394, doi: 10.1109/CVPR.2006.187.
- [9] A. Gunay and V. V. Nabyev, "Automatic age classification with LBP," in Computer and Information Sciences, 2008. ISCIS'08. 23rd International Symposium on, 2008, pp. 1-4.
- [10] S. Chen, C. Zhang and M. Dong, "Deep Age Estimation: From Classification to Ranking," in IEEE Transactions on Multimedia, vol. 20, no. 8, pp. 2209-2222, Aug. 2018.
- [11] E.Eidinger, R. Enbar, and T. Hassner, "Age and gender estimation of unfiltered faces," IEEE Transactions on Information Forensics and Security, vol. 9, pp. 2170-2179, 2014.
- [12] G. Levi and T. Hassner, "Age and gender classification using convolutional neural networks," 2015 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), Boston, MA, 2015, pp. 34-42.
- [13] R. R. Atallah, A. Kamsin, M. A. Ismail, S. A. Abdelrahman and S. Zerdoumi, "Face Recognition and Age Estimation Implications of Changes in Facial Features: A Critical Review Study," in IEEE Access, vol. 6, pp. 28290-28304, 2018, doi: 10.1109/ACCESS.2018.2836924.
- [14] Sawant, M.M., Bhurchandi, K.M. Age invariant face recognition: a survey on facial aging databases, techniques and effect of aging. Artif Intell Rev 52, 981–1008 (2019).
- [15] P. K. Mishra and G. P. Saroha, "A study on video surveillance system for object detection and tracking," 2016 3rd International Conference on Computing for Sustainable Global Development (INDIACom), New Delhi, 2016, pp. 221-226.
- [16] Gaind B., Syal V., Padgalwar S., "Emotion Detection and Analysis on Social Media", arXiv preprint 1901.08458 (2019).
- [17] Zhang K., Zhang Z., Li Z., Qiao Y., "Joint Face Detection and Alignment using Multi-task Cascaded Convolutional Networks", arXiv preprint 1604.02878 (2016).
- [18] R. Takahashi, T. Matsubara and K. Uehara, "Data Augmentation using Random Image Cropping and Patching for Deep CNNs," in IEEE Transactions on Circuits and Systems for Video Technology, doi: 10.1109/TCSVT.2019.2935128.
- [19] Zhang Z., Sabuncu M.R., "Generalized Cross Entropy Loss For Training Deep Neural Networks with Noisy Labels", arXiv preprint 1805.07836 (2018).
- [20] Kingma D. P., Ba J., "Adam: A Method for Stochastic Optimization", arXiv preprint 1412.6980 (2014).
- [21] Moustafa, A.A., Elnakib, A. & Areed, N.F.F. Age-invariant face recognition based on deep features analysis. SIViP (2020).
- [22] S. H. Nam, Y. H. Kim, N. Q. Truong, J. Choi and K. R. Park, "Age Estimation by Super-Resolution Reconstruction Based on Adversarial Networks," in IEEE Access, vol. 8, pp. 17103-17120, 2020, doi: 10.1109/ACCESS.2020.2967800.
- [23] Zhu, H., Zhang, Y., Li, G. et al. Ordinal distribution regression for gait-based age estimation. Sci. China Inf. Sci. 63, 120102 (2020).
- [24] K. Liu, H. Liu, S. Pei, T. Liu and C. Chang, "Age Estimation on Low Quality Face Images," 2019 IEEE International Conference on Artificial Intelligence Circuits and Systems (AICAS), Hsinchu, Taiwan, 2019, pp. 295-296, doi: 10.1109/AICAS.2019.8771612.