



Research paper

Automatic Facial Expression Recognition Method Using Deep Convolutional Neural Network

S. H. Erfani*

School of Engineering, Damghan University, Damghan, Iran.

Article Info

Article History:

Received 24 August 2019

Revised 28 August 2020

Accepted 29 November 2020

DOI: 10.22044/jadm.2020.8801.2018

Keywords:

Automatic Facial Expression Recognition, Convolutional Neural Network, Facial Emotion Analysis, Deep Learning, Effective Computing.

*Corresponding author:
sh.erfani@du.ac.ir (S. H. Erfani).

Abstract

Facial expressions are part of human language and are often used to convey emotions. Since humans are very different in their emotional representation through various media, the recognition of facial expression becomes a challenging problem in machine learning methods. Emotion and sentiment analysis also have become new trends in social media. Deep Convolutional Neural Network (DCNN) is one of the newest learning methods in recent years that model a human's brain. DCNN achieves better accuracy with big data such as images. In this paper an automatic facial expression recognition (AFER) method using the DCNN. In this work, a way is provided to overcome the overfitting problem in training the DCNN for AFER, and also an effective pre-processing phase is proposed that improved the accuracy of facial expression recognition (FER). Here the results for recognition of seven emotional states (neutral, happiness, sadness, surprise, anger, fear, disgust) have been presented by applying the proposed method on the two largely used public datasets JAFFE and CK+. The results show that in the proposed method, the accuracy of AFER is better than traditional FER methods and is about 98.59% and 96.89% for JAFFE and CK+ datasets, respectively.

1. Introduction

Facial expression is an effective way of human communication. Facial expression recognition is the key technology for realizing human-computer interaction to make an emotional computing system. The facial expression has a broad application prospect in many research fields, such as virtual reality, video conference, customer satisfaction survey [1], Internet of Things (IoT). Facial expression recognition is the most important way of human emotional expression in daily emotional communication, just like the tone of voice [2]. For the two past decades, it has been a very important research field in computer vision and image recognition [3]. Nevertheless, facial expression recognition is still a challenging task [4, 5, 6]. The recognition of facial expressions is not an easy problem for machine learning methods since people can vary significantly in the way they show their expressions. Even images of the same person in the same facial expression may

vary in brightness, background, and pose, and these variations are emphasized if considering different subjects [5].

Artificial neural network (ANN) technology has been shown to have advantages over traditional methods in pattern recognition, regression, and categorization [7, 8]. In recent years, Deep Convolutional Neural Network (DCNN) has attracted increasing attention in machine learning and artificial intelligence, and many types of DCNN related algorithms have been successfully applied to image recognition tasks [9, 10]. DCNN's computation intensive tasks can run on GPU which result in high performance at very low power consumption [11]. They have also yielded high performance for some challenges such as the CNN-based model proposed by Kim et al. [12]. CNN is extensively used for facial feature extraction for determining age [13], gender [14], etc.

In this paper, we propose a method for automatic FER using a deep convolutional neural network. The main contributions of the paper include the following items:

- Proposing an effective pre-processing phase improved the accuracy of FER.
- Providing a way to overcome the overfitting problem in training the deep convolutional neural network for FER.
- A deep convolutional neural network is created to extract effective features from training set images.

The rest of the paper is organized as follows: Section 2 gives a background on existing works done by researchers on FER. Section 3 presents the proposed method. The experimental results are presented in Section 4. Finally, Section 5 concludes the paper.

2. Related Works

Several methods have been reported that automatically recognize facial expressions. Some recent approaches for FER have focused on uncontrolled situations, such as the not frontal face, images partially overlapped, spontaneous expressions, etc [5]. It is still a challenging problem [15, 16]. Lucey et al. [17] manually labeled 68 facial points in keyframes and used a gradient descent Active Appearance Model to fit these points in the remaining frames. Liu et al. [16] proposed a novel approach called Boosted Deep Belief Network (BDBN). BDBN is composed of a set of classifiers, named by the authors as weak classifiers. Each weak classifier is responsible for classifying one expression.

To obtain a better representation of facial expressions several deep learning techniques that work on data representations such as DCNNs have been developed [18]. DCNN uses several layers leading to accurate feature learning. Nwosu et al. [18] used a system based on deep convolutional neural network by using facial parts. Their proposed method uses a two-channel convolutional neural network in which Facial Parts are used as input to the first convolutional layer, the extracted eyes are used as input to the first channel while the mouth is the input into the second channel. Burkert et al. [19] also proposed a method based on convolutional neural networks. The authors claim that their method is independent of any handcrafted feature extraction. Liu et al. [20] proposed an Action Unit (AU) that inspired deep networks to explore a psychological theory that expressions can be decomposed into multiple facial expression action units. Softmax Regression-based Deep Sparse Autoencoder

Network (SRDSAN) was proposed by Chen et al. [3] to recognize facial emotion in human-robot interaction. Jain et al. [21] proposed a model based on single DCNNs, which contain convolution layers and deep residual blocks. The convolutional neural networks used in [22], has a conceptual similarity to the one used to build the FACS model. The input is decomposed into features and each deeper layer has a more complex representation which builds upon the previous layer. The final feature representations are then used for classification. Dennis H. et al. [23] in their work, applied a face image to two channels of CNN, the information from the two networks was combined to generate a 94.4% recognition accuracy.

The proposed method is compared to some of the CNN-based methods [3, 18, 21, 22, 23] that their results are available and achieved better accuracy.

3. Proposed Method

In this paper, an automatic facial expression recognition method using DCNN is proposed. DCNNs are powerful models, which are capable to capture effective information, especially for images. The proposed method is consisting of three main steps: 1) Pre-processing, 2) Training, and 3) Prediction. These steps for recognizing facial expressions using the proposed method are shown in Figure 1 and are described in Algorithm 1.

Algorithm 1

1. Pre-processing:

- Face detection using the Viola-Jones standard method and rescaling the face image to 140×140 pixels.
- Finding the center of face image using Canny edge detection and extracting ROI image.
- Adding the average intensity of ROI image to it for normalization.

2. Training:

- The proposed DCNN is trained using 70% of normalized ROI dataset images that are rotated to ± 45 and ± 75 degrees.

3. Prediction:

- The facial emotion of the input image is extracted using pre-trained DCNN.

3.1. Pre-processing

One of the most challenging problems in machine learning is the overfitting problem that occurs when using small datasets [24]. Most facial expression datasets are small in size. To overcome

the overfitting problem, for each image in the dataset, four rotated images are created by turning each image to 45 and 75 degrees clockwise and counterclockwise (± 45 , ± 75). Therefore, the number of dataset images increases by 5 times (4 rotated images and original ones).

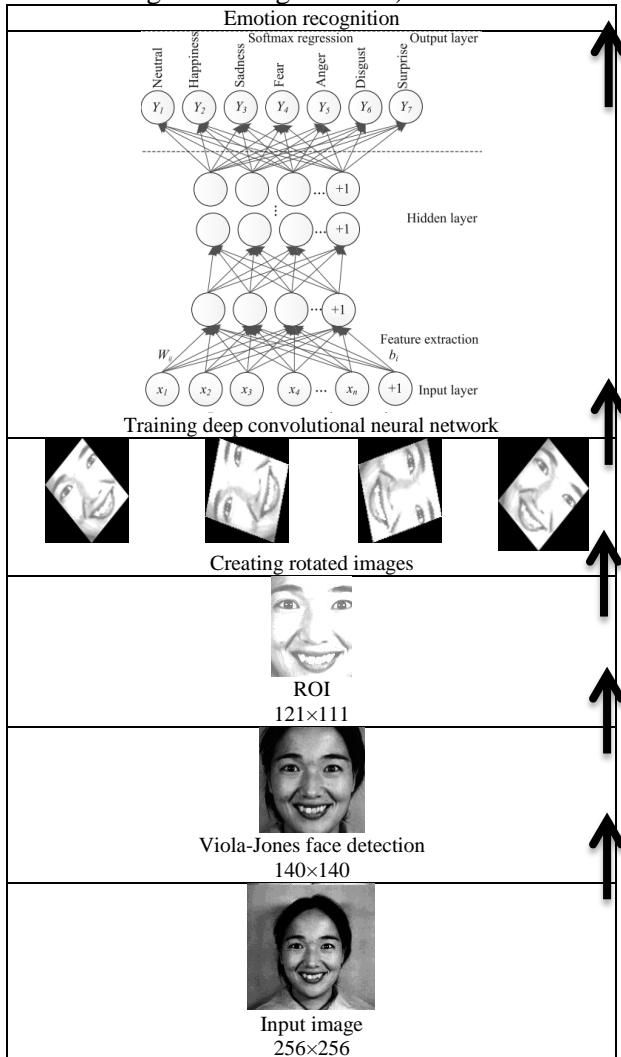


Figure 1. Structure of the proposed method

As presented in Algorithm 1, in the pre-processing step face detection is performed using the Viola-Jones method as a standard method for face and facial region detection/extraction to extract the face area [25]. Then this face area is rescaled to a 140×140 pixel image. After that, to achieve an effective face area, canny edge detection is performed on it. The average rows and columns of detected edge pixels are found and are called AvgRow and AvgCol respectively as the center of the face area image. Then the effective face area is found by selecting pixels from AvgRow-60 to AvgRow+60 in rows and from AvgCol-50 to AvgCol+50 in columns of the input image and is named Region of Instance (ROI). Therefore, in this step, we have a smaller image with 121×111 pixels. Having smaller input images will speed up

learning. After that, the average intensity of ROI is added to the image for normalization. Then four rotated versions of the normalized ROI image are created by turning it for ± 45 and ± 75 degrees.

3.2. Training

DCNN of the proposed method consists of 3 layers, 2 convolutional layers, and one fully connected layer. Each convolutional layer is followed by a ReLU layer, a max-pooling layer, and a normalization layer. The fully connected layer presents the 7-way class predictor which relates to 7 different facial expressions. Training DCNN of the proposed method is done using 70% of dataset images. The input layer of DCNN is the same size as ROI image 121×111 pixels. The optimization of the proposed network is carried out using the stochastic gradient descent method (sgdm). The training starts with a 0.1 learning rate and then it is decreased by a factor of 10 whenever there is no improvement in the validation set accuracy result. The architecture of the proposed deep convolutional neural network is described in Table 1.

Table 1. Architecture of proposed deep convolutional neural network

Type	Filter size/ Stride	Filter number
Input layer	256x256	
Conv1	7x7/2	96
crossChannelNormalization1	4x4	
Maxpool1	3x3/2	
Conv2	5x5/2	256
crossChannelNormalization2	4x4	
Maxpool2	3x3/2	
FC	7	

3.3. Prediction

ROI of the test image is extracted using the Viola-Jones object detection method and then the extracted ROI image is created using the preprocessing step and then it is fed to the pre-trained DCNN to predict its expression.

4. Experimental Results

The proposed method implements on a 2.60GHz Intel Core i7-4510U CPU, 16GB RAM, and developed on MATLAB 2017 software. The performance of the proposed method is evaluated using the classification accuracy experiments performed on the JAFFE and the CK+ datasets. Figure 2 presented the prediction accuracy of the proposed method for training validation on the CK+ and Figure 3 presented prediction accuracy on the JAFFE dataset. These charts clearly show the smooth performance of the proposed method.

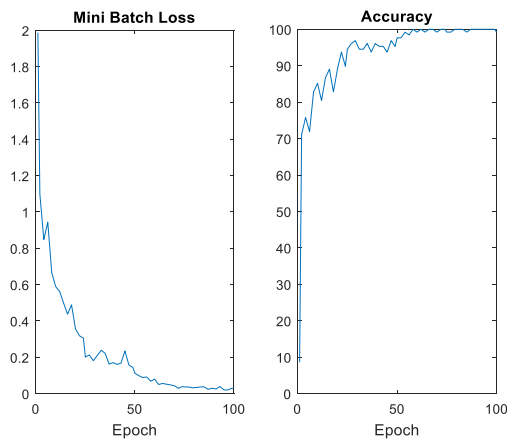


Figure 2. Loss and Prediction Accuracy of the proposed method for training on the CK+ dataset

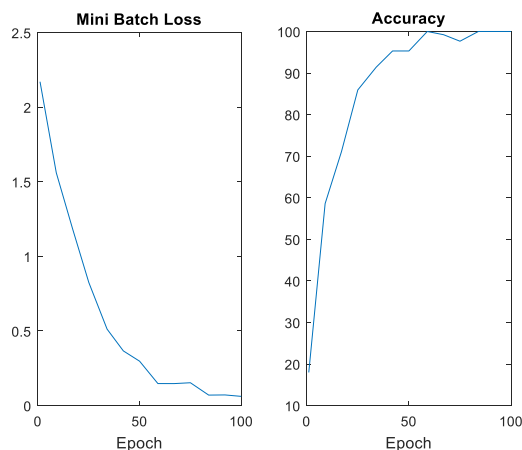


Figure 3. Loss and Prediction Accuracy of the proposed method for training on the JAFFE dataset

4.1. Datasets

To train a deep convolutional neural network, a huge amount of labeled data is required to handle the curse of dimensionality [26]. Several facial expression datasets are publicly available such as Japanese Female Facial Expression (JAFFE) [27] datasets and Cohn-Kanade (CK+) [17] and were used in this paper.

4.1.1. JAFFE Dataset

The Japanese Female Facial Expression (JAFFE) dataset contains 213 facial expression images of female facial expressions corresponding to 10 distinct subjects. Each image is stored at a resolution of 256×256 pixels and 8-bit gray level. Each subject in the dataset is represented with 7 categories of expression (neutral, happiness, sadness, fear, anger, disgust, and surprise), and each subject has 2–4 images per expression.

The network is trained for classification using the training subset of 70% while the testing subset of 30% is used to test the probability that a given facial image belongs to a particular facial expression class. The average recognition

accuracy is used to evaluate the performance of the network. Table 2 shows the recognition accuracy of the proposed method using the training weights obtained with the best accuracy. The recognition accuracy achieved using this method is 98.59%. Table 3 shows the confusion matrix of the recognition accuracy for seven facial expressions on the JAFFE dataset. The proposed method gives a high recognition accuracy of 100% for the surprise, disgust, neutral and fear while Anger, happy and sad had a lower accuracy of 96.77%.

4.1.2. CK+ Dataset

The CK+ expression dataset is used as the other experimental samples, which includes more than 100 performers from different regions, with different colors, ages, and genders, and contains expression image sequences starting from the neutral emotional state and finishing at the expression apex [3]. The neutral emotional state images and the 1-3 last images from each sequence are selected as samples. An example sequence of happy expression is shown in Figure 4. Here we use 902 images that increased to 4510 images (by turning them to ± 15 and ± 30 degrees) 70% of that is selected randomly for training and 30% remains of that for testing. The contempt expression in the dataset was not used in the experiment. Table 4 shows the confusion matrix of the recognition accuracy for seven facial expressions on the CK+ dataset. (e.g. University of Shahrood, shahrood, Iran).

Email address is compulsory for the corresponding author.



Figure 4. An example sequence images of happy expression from the CK+ dataset. The first image has a neutral expression and the 15th image has a happy expression at its peak intensity.

According to Table 2, the CK+ dataset achieved recognition accuracy of 96.89%. This result is lower than that achieved using JAFFE. These results in several angry, happy, neutral, and surprise expressions being confused with sad expressions but the fear, disgust, and sad expressions have 100% accuracy.

Table 2. Comparison of recognition accuracy between the proposed method and some other CNN-based method.

Method	Accuracy JAFFE	Accuracy CK+
Chen et al. [3] (2018)	89.12%	89.03%
Nwosu et al. [18] (2017)	97.71%	95.72%
Jain et al. [21] (2019)	95.23%	93.24%
Santiago et al. [22] (2016)	95.60%	--
Hamester et al. [23] (2015)	94.40%	--
Proposed method	98.59%	96.89%
Simple Trained DCNN	86.85%	89.30%

The proposed method by using the rotation of images in the training phase learning features using DCNN results in 98.59% and 96.89% recognition accuracy for JAFFE and CK+ datasets respectively. According to Table 2, the proposed method gains higher results in comparison with some other CNN-based methods [3, 18, 21, 22, 23] that were mentioned in Section 2 and there are reports of the accuracy of these methods on JAFFE and CK+ datasets.

In Addition, in this work a pre-processing step is proposed that has a positive effect on extracting the effective features of the input image for training DCNN. To present the effect of the pre-

processing step, a DCNN with the same architecture of proposed DCNN is trained on the JAFFE dataset and also CK+ dataset without pre-processing phase and we called it Simple Trained DCNN. The accuracy of Simple Trained DCNN results is 86.85% for JAFFE and 89.30% for the CK+ dataset.

Table 3(a) and 4(a) show the normalized confusion matrix of the simple trained DCNN method on the JAFFE and CK+ dataset respectively. Tables 3(b) and 4(b) are about the normalized confusion matrix of the proposed method on the JAFFE and CK+ dataset respectively. According to these confusion matrices, we can see that the proposed method by using the pre-processing phase can learn useful features.

As shown in Table 5, the training time and test time of the proposed method are higher than the Simple Trained DCNN method, which is due to an increase in the number of training images in the pre-processing phase of the proposed method.

Table 3. Normalized confusion matrix of (a) Simple Trained DCNN and (b) Proposed Method on JAFFE dataset. (ANgry, DIsgust, FEAr, HAPpy, NEUtral, SAd, SURprise).

	AN.	DI.	FE.	HA.	NE.	SA.	SU.
AN.	76.67	10	0	0	0	13.33	0
DI.	10.34	82.76	0	0	0	6.9	0
FE.	3.13	0	93.75	0	0	3.13	0
HA.	0	0	0	83.87	9.68	6.45	0
NE.	0	0	0	0	86.67	13.33	0
SA.	0	0	0	0	0	100	0
SU.	0	0	6.67	3.33	3.33	3.33	83.33

(a)

	AN.	DI.	FE.	HA.	NE.	SA.	SU.
AN.	96.77	0	0	0	0	0	0
DI.	3.33	100	0	0	0	0	0
FE.	0	0	100	0	0	0	0
HA.	0	0	0	96.77	0	3.22	0
NE.	0	0	0	0	100	0	0
SA.	0	3.22	0	0	0	96.77	0
SU.	0	0	0	0	0	0	100

(b)

Table 4. Normalized confusion matrix of (a) Simple Trained DCNN and (b) Proposed Method on CK+ dataset. (ANgry, DIsgust, FEAr, HAPpy, NEUtral, SAd, SURprise).

	AN.	DI.	FE.	HA.	NE.	SA.	SU.
AN.	50	7.14	0	0	35.71	7.14	0
DI.	0	72.22	0	0	27.78	0	0
FE.	0	0	14.29	0	71.43	0	0
HA.	4.76	0	0	85.71	9.52	0	0
NE.	0	0	0	0	100	0	0
SA.	12.5	0	0	0	37.5	50	0
SU.	0	0	0	0	16	0	84

(a)

	AN.	DI.	FE.	HA.	NE.	SA.	SU.
AN.	72.73	0	0	0	27.27	0	0
DI.	0	100	0	0	0	0	0
FE.	0	0	100	0	0	0	0
HA.	0	0	0	94.12	5.887	0	0
NE.	0	0	0	0	100	0	0
SA.	0	0	0	0	28.57	71.43	0
SU.	0	0	0	0	4.76	3.33	95.24

(b)

Table 5. Training time and test time of the proposed method in comparison with Simple Trained DCNN.

Dataset and method	CK+		JAFFE	
	Proposed method	Simple trained DCNN	Proposed method	Simple trained DCNN
Training time(sec)	1502793.99	375790.21	336761.11	132395.31
Test time(sec)	125.322370	37.975661	74.995641	23.842621

5. Conclusion

In this paper, an effective method for automatic recognition of facial expression using a deep convolutional neural network was proposed to address the problems of learning efficiency, where the Viola-Jones object detection method is used to extract the ROI images. Then in the training phase, to overcome the overfitting problem the number of images in datasets increased by five times by rotating them for ± 15 and ± 30 degrees. The facial features are extracted using the deep convolutional neural network. According to the results of the experiment on JAFFE and CK+ datasets, it can be concluded that in the proposed method the features can be learned more efficiently than the Simple Trained DCNN and the methods shown in Table 2.

References

- [1] J. Lia, D. Zhanga, J. Zhanga, T. Lia, Y. Xiaa, Q. Yana, and L. Xuna, "Facial Expression Recognition with Faster R-CNN", *Procedia Computer Science*, vol. 107, pp.135-140, 2017.
- [2] J. Przybyło, "Automatic recognition of facial expressions in the image and analysis of their suitability for control", *doctoral dissertation, AGH University of Science and Technology, Kraków*, 2008.
- [3] M. Rezaei, and V. Derhami, "Improving LNMFP Performance of Facial Expression Recognition via Significant Parts Extraction using Shapley Value", *Journal of AI and Data Mining*, vol. 7, no. 1, pp. 17-25, 2019.
- [4] L. Chen, M. Zhou, W. Su, M. Wu, J. She, and K. Hirota, "Softmax regression based deep sparse autoencoder network for facial emotion recognition in human-robot interaction", *Information Sciences*, vol. 428, pp. 49-61, 2018.
- [5] Q. Mao, Q. Rao, and Y. Yu, "Hierarchical bayesian theme models for multi-pose facial expression recognition", *IEEE Trans. Multimedia*, vol. 16, no. 4, pp. 861-873, 2017.
- [6] L. A. Teixeira, E. De Aguiar, A. De Souza, and T. Oliveira-Santos, "Facial Expression Recognition with Convolutional Neural Networks: Coping with Few Data and the Training Sample Order", *Pattern Recognition*, vol. 61, 2016.
- [7] L. Zhang, and D. Tjondronegoro, "Facial expression recognition using facial movement features", *IEEE Trans. Affect. Comput.*, vol. 2, no. 4, pp. 219-229, 2011.
- [8] C. J. Lin, W. L. Chu, C. C. Wang, G. K. Chen, and I. T. Chen, "Diagnosis of ball-bearing faults using support vector machine based on the artificial fish-swarm algorithm", *Journal of Low Frequency Noise, Vibration and Active Control*, pp. 1-4. 2019 doi: 10.1177/1461348419861822
- [9] B. L. Jian, C. C. Wang, C. T. Hsieh, Y. P. Kuo, M. C. Houg and H. T. Yau, "Predicting spindle displacement caused by heat using the general regression neural network", *The International Journal of Advanced Manufacturing Technology*, 2019 doi: 10.1007/s00170-019-04261-5
- [10] C. Shi, and C-M. Pun, "3d multi-resolution wavelet convolutional neural networks for hyper-spectral image classification", *Inf. Sci.*, vol. 420, pp. 49-65, 2017.
- [11] Y. Wang, X. Wang, and W. Liu, "Unsupervised local deep feature for image recognition", *Inf. Sci.*, vol. 351, pp.67-75, 2016.
- [12] V. Mayya, R. M. Pai, and M. M. Pai, "Automatic Facial Expression Recognition Using DCNN". *Procedia Computer Science*, vol. 93, pp. 453-461, 2016.
- [13] B. K. Kim, J. Roh, S. Y. Dong and S. Y. Lee, "Hierarchical committee of deep convolutional neural networks for robust facial expression recognition", *Journal on Multimodal User Interfaces*, vol. 10, no. 2, 2016.
- [14] X. Wang, R. Guo, and C. Kambhampettu, "Deeply-learned feature for age estimation", *In: Applications of Computer Vision (WACV)*, IEEE Winter Conference on., pp. 534-541, 2015.
- [15] G. Levi, and T. Hassner, "Age and gender classification using convolutional neural networks", *In: Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2015 IEEE Conference on., pp. 34-42, 2015.
- [16] M. K. A. E. Meguid, and M. D. Levine, "Fully automated recognition of spontaneous facial expressions in videos using random forest classifiers", *IEEE Transactions on Affective Computing*, vol. 5, no. 2, pp. 141-154, 2014.
- [17] C. Turan, and K. M. Lam, "Region-based feature fusion for facial expression recognition", *International Conference on Image Processing (ICIP)*, in: 2014 IEEE, pp. 5966-5970, 2014.
- [18] P. Lucey, J. F. Chon, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews, "The Extended Cohn-Kanade Dataset: A complete dataset for action unit and emotion-specified expression", *In: Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2010 IEEE Computer Society Conference on. IEEE, pp. 94-101, 2010.
- [19] L. Nwosu, H. Wang, L. Jiang, I. Unwala, X. Yang, and T. Zhang, "Deep Convolutional Neural Network for Facial Expression Recognition using Facial Parts", *IEEE 15th Intl Conf on Dependable, Autonomic and Secure Computing*, pp. 1318-1321, 2017.
- [20] P. Burkert, F. Trier, M. Z. Afzal, A. Dengel, and M. Liwicki, "Dexpression: Deep convolutional neural network for expression recognition", *CoRR* abs/1509.05371, 2015.

- [21] P. Liu, S. Li, S. Shan, and X. Chen, "Facial expression recognition via a boosted deep belief network", in: *2014 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1805–1812, 2014.
- [22] D. K. Jain, P. Shamsolmoalib, and P. Sehdev, "Extended Deep Neural Network for Facial Emotion Recognition", *Extended Deep Neural Network for Facial Emotion Recognition*, 2019.
- [23] H. C. Santiago, T. Ren, and G. D. C. Cavalcanti, "Facial expression Recognition based on Motion Estimation", *International Joint Conference Neural Networks (IJCNN)*, 2016.
- [24] D. Hamester, P. Barros and S. Wermter, "Face Expression Recognition with a 2-Channel Convolutional Neural Network", *International Joint Conference on Neural Networks (IJCNN)*, 2015.
- [25] Z. Qawaqneh, A. Abu Mallouh, and B. D. Barkana, "Deep Convolutional Neural Network for Age Estimation based on VGG-Face Model", CoRR abs/1709.01664, 2017.
- [26] P. Viola, and M. Jones, "Rapid object detection using a boosted cascade of simple features". In *Proc. of CVPR*. 2001.
- [27] A. Krizhevsky, I. Sutskever and G. E. Hinton, "Imagenet classification with deep convolutional neural networks", *In Advances in neural information processing systems*, pp. 1097–1105, 2012.
- [28] M. Lyons, S. Akamatsu, M. Kamachi, and J. Gyoba, "Coding facial expressions with gabor wavelets", In: *Proceedings of the 3rd. International Conference on Face and Gesture Recognition; FG '98*. Washington, DC, USA: IEEE Computer Society, 1998.

