

Face Emotion Recognition Using Dataset Augmentation Based on Neural Network

Mengyu Rao
Fuzhou University

Ruyi Bao
University of Nottingham

Liangshun Dong*
Shanghai Jiao Tong University

ABSTRACT

Face expression plays a critical role during the daily life, and people cannot live without face emotion. With the development of technology, many methods of facial expression recognition have been proposed. However, from traditional methods to deep learning methods, few of them pay attention to the hybrid data augmentation, which can help improve the robustness of models. Therefore, a method of hybrid data augmentation is highlighted in this paper. The hybrid data augmentation is a method of combining several effective data augmentation. In the experiments, the technique is applied on four basic networks and the results are compared to the baseline models. After applying this technique, the results show that four benchmark models have higher performance than those previously. This approach is simple and robust in terms of data augmentation, which makes it applied in the real world in the future. Besides the results show versatility of the technique as all of our experiments get better results.

CCS CONCEPTS

• **Computing methodologies** → Artificial intelligence.

KEYWORDS

Deep learning, Computer vision, Facial expression recognition, Facial emotion

ACM Reference Format:

Mengyu Rao, Ruyi Bao, and Liangshun Dong*. 2022. Face Emotion Recognition Using Dataset Augmentation Based on Neural Network. In *2022 The 6th International Conference on Graphics and Signal Processing (ICGSP 2022)*, July 1–3, 2022, Chiba, Japan. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3561518.3561519>

1 INTRODUCTION

Facial expression is one of the most external indications of a person's feelings and emotions. In daily conversation, according to the psychologist, only 7% and 38% of information is communicated through words and sounds respective, while up to 55% is through facial expression [13]. It plays an important role in coordinating interpersonal relationships. Ekman and Friesen [6] recognized six

essential emotions in the nineteenth century depending on a cross-cultural study [5], which indicated that people feel each basic emotion in the same fashion despite culture. As a branch of the field of analyzing sentiment [7], facial expression recognition offers broad application prospects in a variety of domains, including the interaction between humans and computers [4], healthcare [10], and behavior monitoring [15]. Therefore, many researchers have devoted themselves to facial expression recognition. In this paper, an effective hybrid data augmentation method is used. This approach is operated on two public datasets, and four benchmark models see some remarkable results.

2 RELATED WORKS

2.1 VggNet

The VGG model [16] was posted by the Visual Geometry Group team at Oxford University. The primary goal of this architecture is to demonstrate how the its final performance can be impacted by increasing network depth. In VGG, 7×7 convolution kernels are replaced by three 3×3 convolution kernels, and 5×5 convolution kernels are replaced by two 3×3 convolution kernels. The main goal of the change is to make sure that the depth of the network and the impact of the neural network can be ameliorated with the condition of the same perceptual field.

2.2 ResNet

The ResNet [2] model won first place in the ImageNet competition [1] held in 2015. The problem that deepening the model can decrease the accuracy was solved by this work. Due to the proposed residual block, it is easy to learn the identity mapping, even though stacked. If there are numerous blocks, redundant blocks can also learn the identity mapping with the help of the residual block. Furthermore, it improves the effectiveness of SGD optimization, which can optimize the network in deeper. What is more, no additional parameters and computational complexity are introduced. Only a very simple addition operation is performed and the complexity is negligible compared to the convolution operation. The ResNet architecture is shown in Figure 1.



Figure 1: The structure of ResNet

2.3 Xception

The Xception [3] model is an upgraded version of the InceptionV3 [17] model. Chollet F offers a new structure of deep convolutional neural network named Xception that replaces the Inception module

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICGSP 2022, July 1–3, 2022, Chiba, Japan

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9637-0/22/06...\$15.00

<https://doi.org/10.1145/3561518.3561519>

with a depthwise separable convolution. The residual network and the depthwise separable convolution are the fundamental components of this network. Xception is typically composed of 36 convolutional layers grouped into 14 blocks, with 12 blocks in the middle containing all linear residual connections. Simultaneously, the model holds the properties of depthwise separable convolution[9] since the model executes spatial layer-by-layer convolution on every channel of the inputs individually, and then conducts point-by-point convolution on the output.

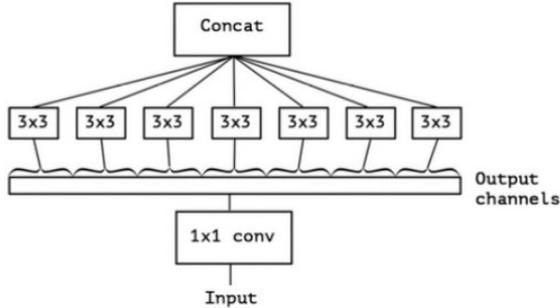


Figure 2: The network structure of Xception

3 APPROACH

3.1 HDA: hybrid data augmentation

3.1.1 Horizontal Flip. Geometric transformation is one of the most basic methods to augment data. Because of the particularity of the images, which means that facial expression images emoticons do not undergo a large degree of distortion and rotation in most cases, the horizontal flip (HF) is used to ensure that these images are consistent. Every image in the original dataset is horizontally flipped to create a mirror image. The formula of this method is shown in Eq. (1).

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} width - 1 \\ 0 \end{bmatrix} + \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix}$$

3.1.2 Gaussian Noise. Gaussian noise (GN) represents a kind of statistical noise that has a probability density function equivalent to that of the normal distribution. In the proposed approach, the training images are added with Gaussian noise to simulate noise that may happen in the real world so that the model can become robust against the original images. The formula of Gaussian noise can be written as Eq. (2).

$$GN(x, y) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)}\left(\frac{(x-\mu_1)^2}{\sigma_1^2} - \frac{2\rho(x-\mu_1)(y-\mu_2)}{\sigma_1\sigma_2} + \frac{(y-\mu_2)^2}{\sigma_2^2}\right)}$$

4 EXPERIMENTS AND RESULTS

4.1 Datasets

Following related work on facial emotion recognition, these experiments are conducted on the two benchmark public face emotion datasets: Ck+ dataset[12], and Fer2013 dataset [8]. Ck+: The most widely utilized laboratory-controlled dataset for facial expression

recognition is the Extended Cohn–Kanade (Ck+) Dataset[12] (some samples are shown in Figure 3). Sequences that change from neutral to peak expression are included in the Ck+ dataset. Extraction of the final 1 to 3 frames which have peak formation and the first frame of every sequence is the most common data selection approach for evaluation. Then, people are divided into n groups for person-independent n-fold cross-validation experiments, where n is typically between 5, 8, and 10. Fer2013: Fer2013 [8] is a large-scale

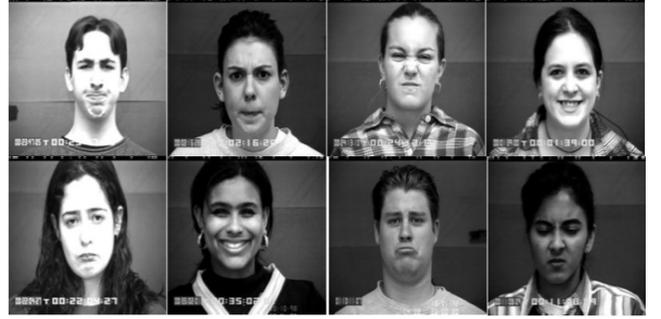


Figure 3: Some samples of the Ck+ dataset

dataset acquired automatically by the Google image search API (some samples are shown in Figure 4). 35887 images are contained in the Fer2013 dataset, and each image is labeled as one of the seven basic emotions. All of the images in this dataset are grayscale images. Furthermore, this dataset contains 547 disgusted images, 5121 fear images, 4953 angry images, 8989 happy images, 4002 surprised images, 6077 sad images, and 6198 neutral images.



Figure 4: Some samples of the Fer2013 dataset

Some information about the datasets are shown in Table 1.

Dataest	Number	Number of Emotion	Gender
JAFFE	213 images+	7	Female
CK+	593 videos	7	Female & Male
Fer2013	35886 images	7	Female & Male

Table 1: Datasets information

4.2 Experimental settings

The experiments are implemented via PyTorch [14], and the NVIDIA GTX 2080Ti with 4 CPU cores and 13 Gigabytes of RAM is used for experiments. The Adam optimizer [11] is used to train the networks with a learning rate of $2e-4$ and betas of 0.9 and 0.999. The best model is selected through the principle of selection of the best accuracy of several experiments.

4.3 Experimental evaluation metric

Face emotion recognition can be viewed as a multi-classification problem. Accuracy (Acc) is used as the evaluation metric in this paper, and the calculation formula is as follows: $Acc = (TP + TN) / (TP + TN + FP + FN)$. (3) In this formula, TP stands for a positive sample predicted by the model as a positive sample, TN stands for a negative sample predicted by the model as a negative sample, FP stands for a negative sample predicted by the model as a positive sample, and FN stands for a positive sample predicted by the model as a negative sample.

4.4 Performance on Ck+ dataset

This dataset consists of 8 emotions with a total of 981 trainable images, all of which are 640 pixels \times 490 pixels in size. 593 trainable images from 7 expressions are selected for the experiment (shown in Figure 5). However, the background of the volunteers in the pictures is larger than the face images. Therefore, the image size of the dataset is processed to 48 \times 48 pixels, which is convenient for the model input size to be uniform. The results of the Ck+ dataset



Figure 5: Images of different expressions in the Ck+ dataset

and the HDACK+ dataset are shown in Table 2, where both the Ck+ dataset and the HDACK+ dataset are divided in the ratio of 8:1:1 for training, validating, and testing respectively. The batch size used is 32, and 20 epochs are used for training.

Model	Ck+	HDACK+
Vgg19	74.19%	97.80%
Resnet18	90.32%	100%
Resnet50	95.70%	99.73%
Xception	83.87%	99.73%

Table 2: Accuracy Comparisons of Different Models on Ck+ dataset

Tables 2 and Figure 6 indicate the accuracies of the Ck+ dataset and those of the HDACK+ dataset. The dataset with data augmentation always has a higher performance in most models, and the ResNet18 model achieved 100% testing accuracy in the HDACK+ dataset. After analyzing the results, significant improvement in accuracy can be seen on all models, especially for Vgg19. The accuracy

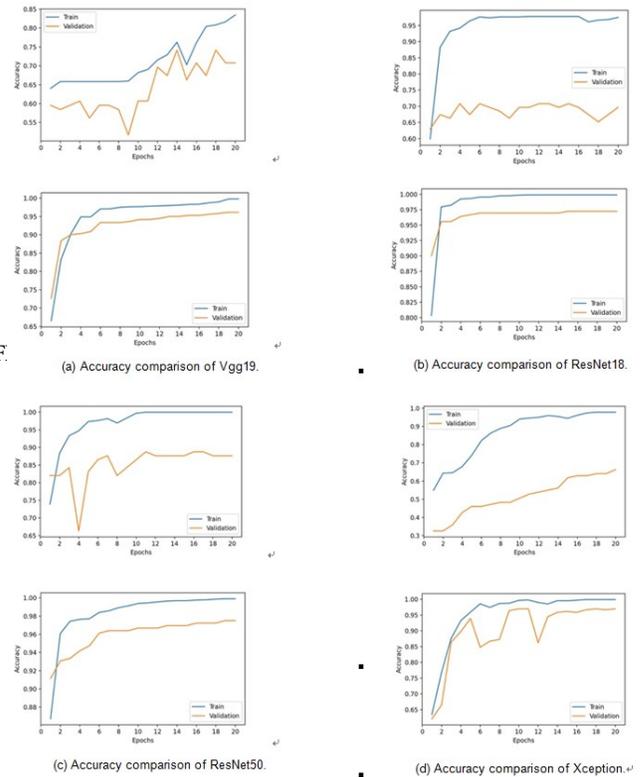


Figure 6: Accuracy comparison of different models on Ck+ dataset and HDACK+ dataset. (Best viewed in color)

comparison of Vgg19 shows that the accuracy improves by 23.61%. The ResNet50 network is the least improved one, but the accuracy is also 4.03% higher than the original data. Furthermore, the figures show that the data augmentation improves the feature learning for every model. For the left one of each column, the Ck+ dataset is used. The training accuracies have an upward tendency, but the validation accuracies always have some shocks, even increasing slowly. Especially for ResNet18, the curve always fluctuates around 70%. However, when comes to the right one of each column, the HDACK+ dataset is used. Both the curve of training accuracy and the curve of validation have an increasing toward, and the validation accuracies of all figures in the right position come to a peak of around 98%.

4.5 Performance on Fer2013 dataset

After rejecting wrongly labeled frames, each image was registered and scaled to 48 \times 48 pixels. There are 28,709, 3,589 and 3,589 images respectively for training, validation, and testing with seven expression labels in the Fer2013 dataset (shown in Figure 7). The results of the Fer2013 dataset and the HDAFer2013 dataset are presented in Table 3, where both the Fer2013 dataset and the HDAFer2013 dataset are divided in the ratio of 8:1:1 for training, validating, and testing respectively. The batch size used is 32, and 20 epochs are used for training.



Anger Disgust Fear Happiness Neutral Sadness Surprise

Figure 7: Images of different expressions in the Fer2013 dataset

Model	Fer2013	Fer2013 _{DA}
Vgg19	62.13%	84.87%
Resnet18	65.67%	88.32%
Resnet50	62.55%	88.17%
Xception	60.635%	82.68%

Table 3: Accuracy Comparisons of Different Models on Fer2013 dataset

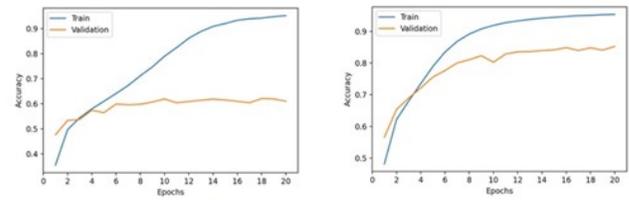
Tables 3 and Figure 8 indicate the accuracies of the Fer2013 dataset and those of the HDAFer2013 dataset. The dataset used data augmentation always has a higher performance in most models. The results show that improvement in accuracy can be seen in all models. The accuracy comparison of ResNet50 shows that the accuracy improves by 25.62%. The Xception network is the least improved, but the accuracy on it is also 22.33% higher than that on the original data. Furthermore, the figures show that the data augmentation improves the feature learning for every model. For the left one of each column, the Fer2013 dataset is used. The training accuracies have an upward tendency, but the validation accuracies always fluctuate around 63%. However, when comes to the right one of each column, the HDAFer2013 dataset is used. Both the curve of training accuracy and the curve of validation have an increasing toward, and the validation accuracies of all figures in the right position come to a peak above 82%.

5 CONCLUSION

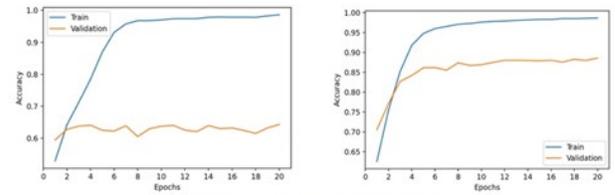
In this study, a hybrid data augmentation method of the dataset, which has a high performance in some models, is presented. Because convolutional neural networks require more samples for training to get accurate and robust result, the hybrid data augmentation method is used to enlarge the number of samples. After applying the technique, the numbers of images in both the Ck+ dataset and the Fer2013 dataset have increased, and four benchmark models have higher performance than those previously. This approach is simple and robust in terms of data augmentation, which makes it applicable in the real world in the future.

REFERENCES

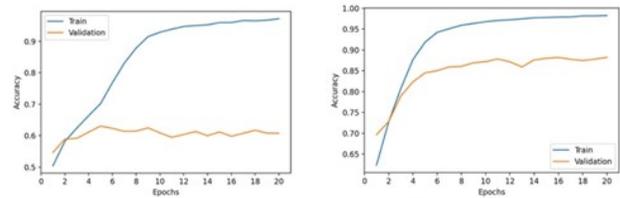
[1] [n.d.]. *ImageNet Large Scale Visual Recognition Challenge (ILSVRC)*. <https://www.image-net.org/challenges/LSVRC/> (2022, Jul 10).
 [2] J. Carreira, H. Madeira, and J.G. Silva. 1998. Xception: a technique for the experimental evaluation of dependability in modern computers. *IEEE Transactions on Software Engineering* 24, 2 (1998), 125–136. <https://doi.org/10.1109/32.666826>
 [3] François Chollet. 2017. Xception: Deep Learning with Depthwise Separable Convolutions. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 1800–1807. <https://doi.org/10.1109/CVPR.2017.195>



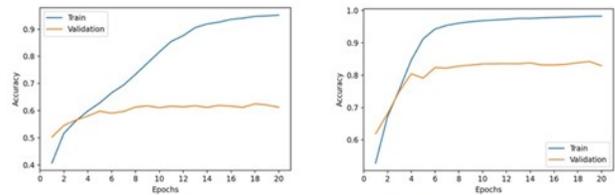
(a) Accuracy comparison of Vgg19.



(b) Accuracy comparison of ResNet18.



(c) Accuracy comparison of ResNet50.



(d) Accuracy comparison of Xception.

Figure 8: Accuracy comparison of different models on Fer2013 dataset and HDAFer2013 dataset. (Best viewed in color)

[4] Jia Deng, Gaoyang Pang, Zhiyu Zhang, Zhibo Pang, Huayong Yang, and Geng Yang. 2019. cGAN Based Facial Expression Recognition for Human-Robot Interaction. *IEEE Access* 7 (2019), 9848–9859. <https://doi.org/10.1109/ACCESS.2019.2891668>
 [5] Pual Ekman. 1994. Strong evidence for universals in facial expressions: A reply to Russell’s mistaken critique. *Psychological Bulletin* 115 2 (1994), 268–287.
 [6] Paul Ekman and W V Friesen. 1971. Constants across cultures in the face and emotion. *Journal of personality and social psychology* 17 2 (1971), 124–9.
 [7] Chenquan Gan, Lu Wang, Zufan Zhang, and Zhangyi Wang. 2020. Sparse attention based separable dilated convolutional neural network for targeted sentiment analysis. *Knowledge-Based Systems* 188 (2020), 104827. <https://doi.org/10.1016/j.kmosys.2019.06.035>
 [8] Ian J. Goodfellow, Dumitru Erhan, Pierre Luc Carrier, Aaron Courville, Mehdi Mirza, Ben Hamner, Will Cukierski, Yichuan Tang, David Thaler, Dong-Hyun Lee, Yingbo Zhou, Chetan Ramaiah, Fangxiang Feng, Ruifan Li, Xiaojie Wang, Dimitris Athanasakis, John Shawe-Taylor, Maxim Milakov, John Park, Radu Ionescu, Marius Popescu, Cristian Grozea, James Bergstra, Jingjing Xie, Lukasz Romaszko, Bing Xu, Zhang Chuang, and Yoshua Bengio. 2015. Challenges in representation learning: A report on three machine learning contests. *Neural Networks* 64 (2015), 59–63. <https://doi.org/10.1016/j.neunet.2014.09.005> Special Issue on “Deep Learning of Representations”.

- [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 770–778. <https://doi.org/10.1109/CVPR.2016.90>
- [10] Chaudhary Muhammad Aqduus Ilyas, Mohammad Ahsanul Haque, Matthias Rehm, Kamal Nasrollahi, and Thomas Baltzer Moeslund. 2018. Facial Expression Recognition for Traumatic Brain Injured Patients. In *VISIGRAPP*.
- [11] Diederik P. Kingma and Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. *CoRR* abs/1412.6980 (2015).
- [12] Patrick Lucey, Jeffrey F. Cohn, Takeo Kanade, Jason M. Saragih, Zara Ambadar, and I. Matthews. 2010. The Extended Cohn-Kanade Dataset (CK+): A complete dataset for action unit and emotion-specified expression. *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition - Workshops* (2010), 94–101.
- [13] Albert Mehrabian and Suzanne Ferris. 1967. Inference of attitudes from nonverbal communication in two channels. *Journal of consulting psychology* 31 3 (1967), 248–52.
- [14] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf>
- [15] Yassine Rabhi, Makrem Mrabet, Farhat Fnaiech, and Mounir Sayadi. 2018. A real-time emotion recognition system for disabled persons. *2018 4th International Conference on Advanced Technologies for Signal and Image Processing (ATSIP)* (2018), 1–6.
- [16] Karen Simonyan and Andrew Zisserman. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *International Conference on Learning Representations*.
- [17] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the Inception Architecture for Computer Vision. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2818–2826. <https://doi.org/10.1109/CVPR.2016.308>