

A Preliminary Study on Data Augmentation of Deep Learning for Image Classification

Benlin Hu, Cheng Lei, Dong Wang, Shu Zhang, Zhenyu Chen*
State Key Lab of Novel Software Technology, Nanjing University, China
*zychen@nju.edu.cn

Abstract—Deep learning models have a large number of free parameters that need to be calculated by effective training of the models on a great deal of training data to improve their generalization performance. However, data obtaining and labeling is expensive in practice. Data augmentation is one of the methods to alleviate this problem. In this paper, we conduct a preliminary study on how three variables (augmentation method, augmentation rate and size of basic dataset per label) can affect the accuracy of deep learning for image classification. The study provides some guidelines: (1) it is better to use transformations that alter the geometry of the images rather than those just lighting and color. (2) 2-3 times augmentation rate is good enough for training. (3) the smaller amount of data, the more obvious contributions could have.

Index Terms—Data Augmentation, Deep Learning, Quality Assurance

I. INTRODUCTION

Deep learning is powerful, but they usually need to be trained on massive amounts of data to perform well, which can be considered as a major limitation [1] [2]. Deep learning [3] models trained on small dataset show the low performance of versatility and generalization from the validation and test set. Hence, these models suffer from the problem caused by over-fitting.

A quantity of methods have been proposed to reduce the over-fitting problem [4]. Data augmentation, which increases both the amount and diversity of data by “augmenting” them, is an effective strategy that we can reduce over-fitting on models and improve the diversity of the dataset and generalization performance [5]. In the application of image classification, there have already been some general augmentations, like flipping the image horizontally or vertically, translating the image by a few pixels and so on.

In order to have an in-depth investigation and propose a guideline about how to use the augmentation methods properly, we perform a preliminary experiment to summarize guidelines and testify the explanation of the phenomenon. In the experiment, 10 state of the art augmentation methods are studied on two popular datasets, which representing gray images and color images.

The datasets used in the experiment are the MNIST [6] and CIFAR-10 [7]. MNIST consists of 60000 handwritten digit images in the training set and 10000 in the test set, which are in gray-scale with 10 classes with image dimensions of $28 \times 28 \times 1$. CIFAR-10 contains 50000 training and 10000 testing $32 \times 32 \times 3$ color images with 10 classes as well.

The main contributions of this paper are summarized as follows. A preliminary experiment has been conducted to find out how the accuracy of a deep learning model can be affected by the three variables: **augmentation method, augmentation rate and size of basic dataset per label**. According to the experimental results, some guidelines based on the three variables have been summarized to use augmentation methods properly. The further experimental results show that the simple augmentation methods, such as altering the geometry of the image, have better effectiveness than the complicated ones.

II. DATA AUGMENTATION

Based on the existing work, 10 state of the art augmentation methods, shown in Fig. 1, for image classification are adopted in the experiment. As for the selection of experimental models, ResNet-20 and LeNet-5 are chosen for CIFAR-10 and MNIST to conduct our experiment, respectively. Table I lists the methods and their own descriptions and ranges we used to augment data.

The reason why different augmentation methods can improve the model performance is because they mimic the image with different *features* when taking photos. Furthermore, these augmentation methods have different effects in improving accuracy, because there are differences existing in the quantity and quality of these *features* in the datasets. Therefore, it is valuable to study how these *features* can influence on the model performance.

The original dataset with c classes is denoted as M . The basic dataset used to train model is defined as $N \subseteq M$, and the size of dataset per label in the basic dataset N is $cnum$. f and $rate$ are the augmentation method and its rate we chose, respectively. The dataset, which is augmented from N by the augmentation method f , is named augmented dataset N' . Obviously, the size of basic dataset $|N|$ is $c * cnum$ and the size of augmented dataset $|N'|$ is $c * cnum * rate$. The model, which are trained by the basic dataset N and augmented dataset N' , are denoted as m_{basic} and $m_{augmented}$, respectively. The accuracy of m_{basic} is defined as acc_{basic} , and the counterpart of $m_{augmented}$ is $acc_{augmented}$. The augmented test dataset is denoted as T' , which is augmented from test dataset T by the augmentation method f . acc_{aug_test} is the accuracy of the m_{basic} evaluated by augmented test dataset T' . And we use *feature rate*, defined as $(acc_{basic} - acc_{aug_test})$, to represent the *features* that the model m_{basic} has not yet learned from the basis dataset N .

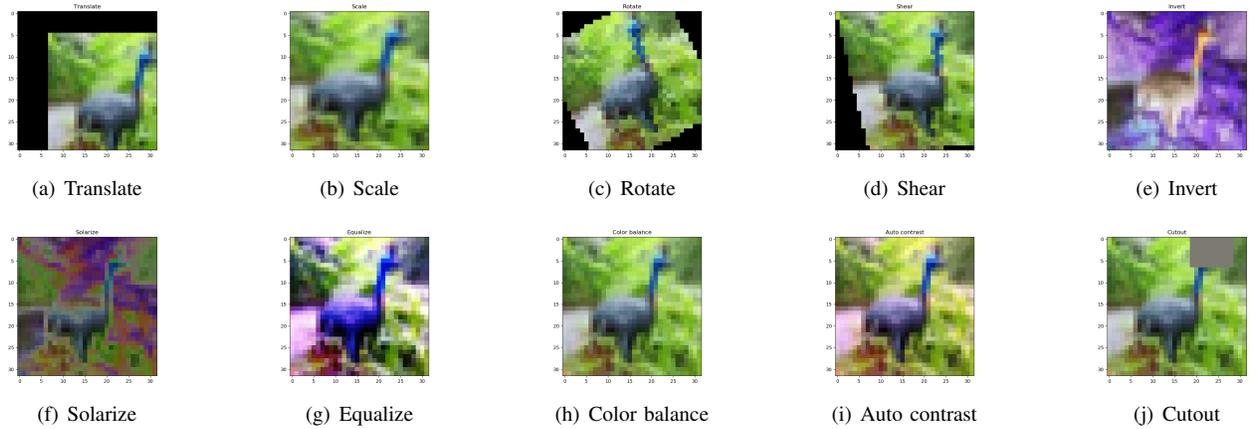


Fig. 1. Examples of augmentation methods

TABLE I
AUGMENTATION METHODS IN EXPERIMENT

| Name | Description | Range |
|----------------------|--------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------|
| <i>Translate</i> | Translate the image in the horizontal and vertical direction with rate <i>magnitude</i> | $[-0.3, 0, 3]$ |
| <i>Scale</i> | Zoom in or Zoom out the image with rate <i>magnitude</i> and select the center of the scaled image | $[-0.5, 0.5]$ |
| <i>Rotate</i> | Rotate the image <i>magnitude</i> degrees | $[-30^\circ, 30^\circ]$ |
| <i>Shear</i> | Shear the image along the horizontal or vertical axis with rate <i>magnitude</i> | $[-0.3, 0.3]$ |
| <i>Invert</i> | Invert the pixels of the image | |
| <i>Solarize</i> | Invert all pixels above a threshold value of <i>magnitude</i> | $[0, 256]$ |
| <i>Equalize</i> | Equalize the image histogram | |
| <i>Color balance</i> | Adjust the color balance of the image. A <i>magnitude</i> = 0 gives a black and white image, while <i>magnitude</i> = 1 gives the original image | $[0.1, 1.9]$ |
| <i>Auto contrast</i> | Maximize the contrast of the image | |
| <i>Cutout</i> | Set a random square patch of side-length <i>magnitude</i> pixels to gray | $[0, 20]$ |

III. EXPERIMENT

A. Experiment Design

Our experiment is designed as follows:

- 1) The training set N is randomly chosen from M in accordance with $cnum$. The basic dataset N is used to train the model m_{basic} .
- 2) On the basis of N , each data is augmented $rate$ times by the augmentation methods f .
- 3) A new model $m_{augmented}$ is trained on the augmented dataset N' .
- 4) The model m_{basic} is evaluated with accuracy acc_{aug_test} by using augmented test dataset T' .
- 5) The improvement of accuracy and the *feature rate* of the model are used as evaluation criteria.

In order to obtain convincing results without loss of generality, each experiment is repeated 10 times, and the average values are taken as the final experimental results.

B. Experimental Results

The experimental results are shown in Fig. 2 and 3. For MNIST, the three methods of *Translate*, *Scale* and *Rotate* are good and relatively stable, while the *Shear* method is good but unstable. For CIFAR-10, almost all augmentation methods can improve accuracy. Among them, the *Rotate*, *Scale*, *Shear* and *Translate* have outstanding performance, which is 6%-11% higher than the control group. The augmentation methods are roughly sorted in descending order of improvement as follows: (1) *Translate*, *Scale*, *Shear*. (2) *Rotate*. (3) *Cutout*, *Solarize*. (4) *Invert*, *Equalize*, *Auto Contrast*, *Color balance*.

The accuracy improvement is used to evaluate the performance of different augmentation methods. The average values and medians of the improvement are listed in Table II and Table III. When it comes to a small amount of data, the augmentation methods can still work well. The smaller the amount of data, the more obvious the contribution it has, which indicates that data augmentation has a large space and potential in the field where learnable data is scarce, such as rare disease diagnosis and large earthquake prediction.

TABLE II
THE FEATURE RATE AND THE ACCURACY IMPROVEMENT ON MNIST

| method | translate | scale | rotate | shear | invert |
|-----------------------|-----------|--------|--------|--------|--------|
| <i>fea. rate avg.</i> | 0.7695 | 0.1598 | 0.0791 | 0.4741 | 0.5045 |
| <i>impro. avg.</i> | 0.0067 | 0.0110 | 0.0055 | 0.0120 | 0.0046 |
| <i>impro. median</i> | 0.0061 | 0.0078 | 0.0086 | 0.0079 | 0.0034 |

TABLE III
THE FEATURE RATE AND THE ACCURACY IMPROVEMENT ON CIFAR-10

| method | translate | scale | rotate | shear | invert |
|-----------------------|-----------|----------|-----------|-----------|--------|
| <i>fea. rate avg.</i> | 0.2731 | 0.2640 | 0.2153 | 0.2175 | 0.3722 |
| <i>impro. avg.</i> | 0.1294 | 0.1339 | 0.1038 | 0.1348 | 0.0821 |
| <i>impro. median</i> | 0.1084 | 0.1149 | 0.0945 | 0.1098 | 0.0500 |
| method | solarize | equalize | col. bal. | auto con. | cutout |
| <i>fea. rate avg.</i> | 0.1599 | 0.1396 | 0 | 0.0133 | 0.063 |
| <i>impro. avg.</i> | 0.0617 | 0.0467 | 0.0224 | 0.0299 | 0.0469 |
| <i>impro. median</i> | 0.0438 | 0.0516 | 0.0194 | 0.0343 | 0.0417 |

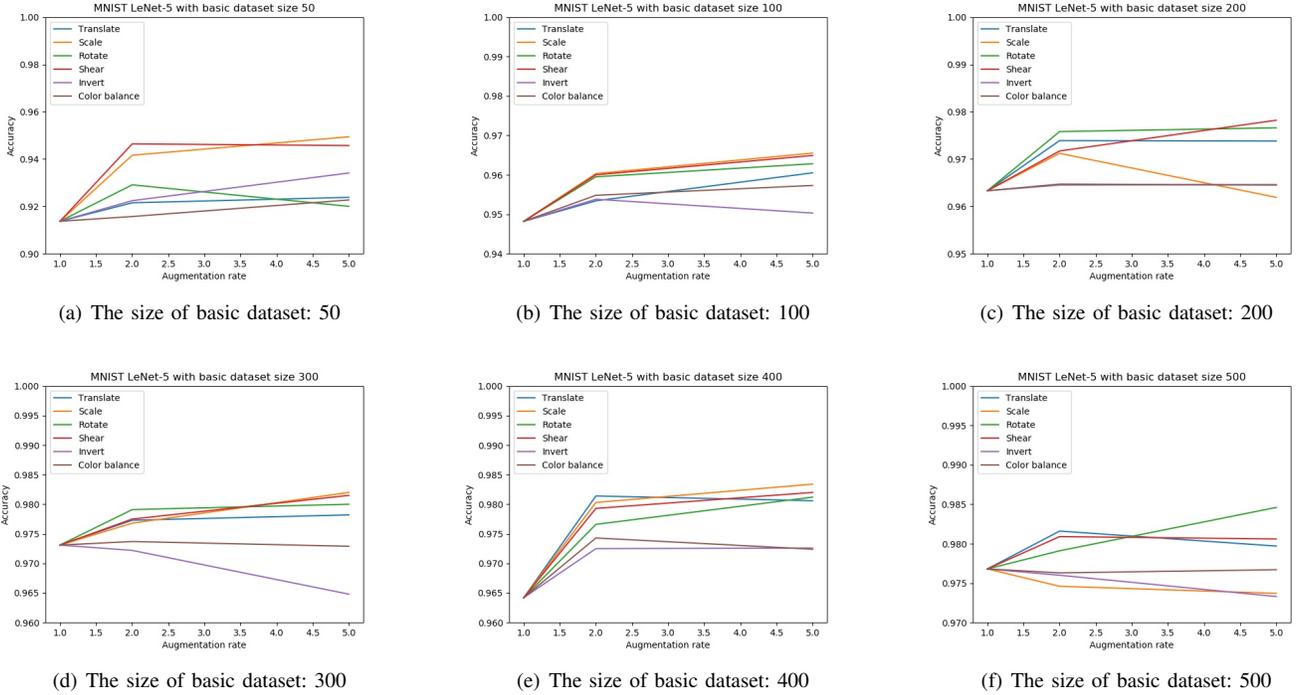


Fig. 2. The experimental results of LeNet-5 on MNIST

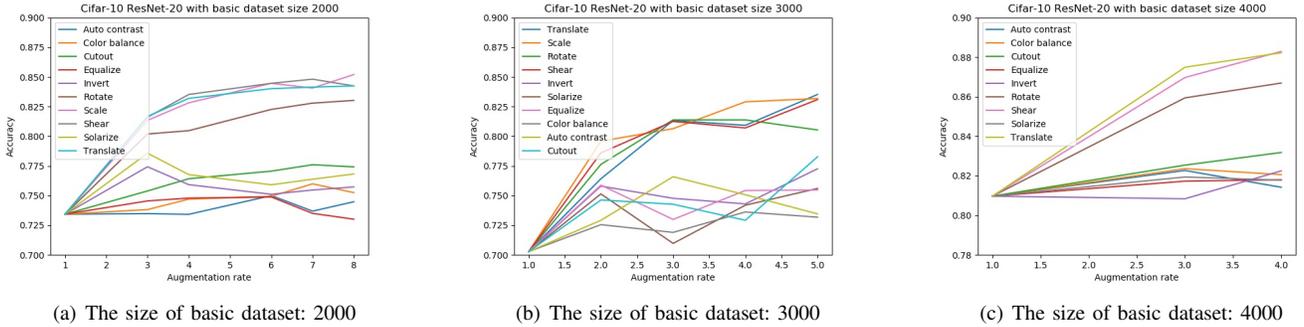


Fig. 3. The experimental results of ResNet-20 on CIFAR-10

Most augmentation methods can improve the accuracy by increasing the augmentation rate. However, the time, space and other costs increase rapidly, but the improvement of accuracy increases slowly. Small-scale experiments with higher augmentation rate also show that this is not cost-effective. In practice, unless there are extremely stringent requirements for accuracy, it is generally not necessary to take 5 times or more for a little improvement.

The results show that the simpler the augmentation method, the more obvious the improving effectiveness. This confirms the assumption when using a more complex model that it is better to use transformations that alter the geometry of the images rather than just lighting and coloring [8].

Figure 4 shows the results of positive correlation between the *feature rate* and the accuracy improvement on both

datasets, especially on CIFAR-10. High *feature rate* means that the model has not learned enough *features* of images when it was trained on the basic dataset N . And being trained on augmented dataset N' provides the model more *features*, thus the accuracy of the model improves. In the results of CIFAR-10, geometric methods generally better than photometric methods. This may inspire us that *features* could be one of the reasons why it is better to use transformations that alter the geometry of the images rather than just lighting and coloring.

IV. RELATED WORK

Image processing methods are implemented by PIL which accept an image as input and output a processed image [9]. The methods include ShearX/Y, TranslateX/Y, Rotate,

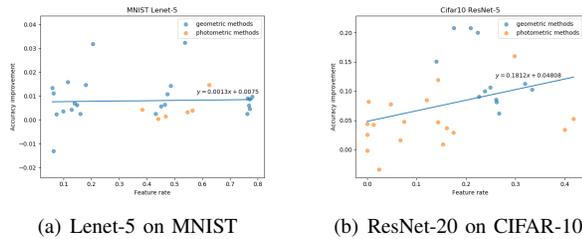


Fig. 4. The correlations between feature rate and the accuracy improvement

AutoContrast, Invert, Equalize, Solarize, Posterize, Contrast, Color, Brightness and Sharpness, as shown in Figure 1.

Various geometrical and photometric schemes are evaluated on a coarse-grained dataset using a relatively simple CNN [8]. The experimental results indicate that, under these circumstances, *cropping in geometric augmentation* significantly increases CNN task performance.

Affine transformation, a 2D geometric transform method, is based on reflecting the image, scaling and translating the image, and rotating the image by different angles. The affine augmentation method is very common and widely used for correcting geometric distortion introduced by perspective [10].

Cutout is originally considered as a targeted method for removing visual features with high activations in later layers of a convolutional neural network (CNN). However, the results in [11] [12] show that randomly selecting a rectangle region in an image and erasing its pixels with random values can be used to improve the overall performance of CNNs.

Histogram equalization is introduced as a data augmentation method [13]. Histogram equalization, solarization and adjusting image color balance are common methods used in digital image processing. These methods can simulate the problems encountered when taking photos improperly, like harsh lighting combined with auto-white-balance will produce images that over or under exposed.

There are many existing studies discussed data augmentation methods. However, many of them focus on new data augmentation methods. The empirical study in this paper provides some guidelines to use data augmentation methods properly according to the application scenario.

V. CONCLUSION AND FUTURE WORK

Data augmentation methods generate new data by performing image processing methods such as rotation and translation on the training set. Therefore, these methods expand the training set in the amount and generalization degree. Our experimental results show that these augmentation methods work well even on small dataset. Applying higher augmentation rate is not cost-effective, because the marginal benefit is gradually reduced while the time, space and other costs increase linearly. Our study recommends that the best augmentation rate is 2-3 times. The simple augmentation methods such as translation, rotation, scaling and shearing can achieve good results, and in fact, much better than more complicated methods.

The results show that the performance of professional image processing methods such as *Equalize* and *Auto Contrast* is poor. A possible reason is that *features* pointed to by the simple method are more common in most images. However, models trained with training set which augmented by simple methods can get more learning data and stronger generalization ability, because an augmentation method means adding more *features* to the data, and then the model can eliminate the interference caused by the *features* and perform better. Visualization [14] may be used to understand the operating principle of augmentation methods.

In this paper, due to the consideration of time budget, the study is preliminary. Larger datasets and more complicated models will be studied in the future. There are some opportunities to improve effectiveness via combining different single augmentation methods. The mutation methods may be another strategy to [15] improve data augmentation. The data augmentation methods can also be considered in test set to reduce the labeling cost [16].

VI. ACKNOWLEDGEMENT

The work is partly supported by the National Natural Science Foundation of China (61832009, 61802171).

REFERENCES

- [1] Z.-H. Zhou and X.-Y. Liu, "Training cost-sensitive neural networks with methods addressing the class imbalance problem," *IEEE Transactions on Knowledge & Data Engineering*, no. 1, pp. 63–77, 2006.
- [2] H. Salehinejad, J. Baarbe, S. Sankar, J. Barfett, E. Colak, and S. Valaei, "Recent advances in recurrent neural networks," *CoRR*, vol. abs/1801.01078, 2018.
- [3] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *nature*, vol. 521, no. 7553, p. 436, 2015.
- [4] B. Wang and D. Klabjan, "Regularization for unsupervised deep neural nets," in *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [5] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems 25*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2012, pp. 1097–1105.
- [6] Y. LeCun, "The mnist database of handwritten digits," <http://yann.lecun.com/exdb/mnist/>, 1998.
- [7] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," Citeseer, Tech. Rep., 2009.
- [8] L. Taylor and G. Nitschke, "Improving deep learning using generic data augmentation," *CoRR*, vol. abs/1708.06020, 2017.
- [9] E. D. Cubuk, B. Zoph, D. Mané, V. Vasudevan, and Q. V. Le, "Autoaugment: Learning augmentation policies from data," *CoRR*, vol. abs/1805.09501, 2018.
- [10] C. C. Stearns and K. Kannappan, "Method for 2-d affine transformation of images," Dec. 12 1995, uS Patent 5,475,803.
- [11] T. Devries and G. W. Taylor, "Improved regularization of convolutional neural networks with cutout," *CoRR*, vol. abs/1708.04552, 2017.
- [12] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang, "Random erasing data augmentation," *CoRR*, vol. abs/1708.04896, 2017.
- [13] R. Dellana and K. Roy, "Data augmentation in cnn-based periorcular authentication," in *2016 6th International Conference on Information Communication and Management (ICICM)*, Oct 2016, pp. 141–145.
- [14] X. Zhang, Z. Yin, Y. Feng, Q. Shi, J. Liu, and Z. Chen, "Neuralvis: Visualizing and interpreting deep learning models," *arXiv preprint arXiv:1906.00690*, 2019.
- [15] W. Shen, J. Wan, and Z. Chen, "Munn: Mutation analysis of neural networks," in *2018 IEEE International Conference on Software Quality, Reliability and Security Companion (QRS-C)*. IEEE, 2018, pp. 108–115.
- [16] Q. Shi, J. Wan, Y. Feng, C. Fang, and Z. Chen, "Deepgini: Prioritizing massive tests to reduce labeling cost," *arXiv preprint arXiv:1903.00661*, 2019.