
DD2424 Group 118:

The mechanisms, powers and limitations of some Data Augmentation techniques

Anton Stråhle

Jan Alexandersson

Fredrika Lundahl

Abstract

To obtain good results in deep learning the quality and quantity of the data is crucial. A smart and cheap way of increasing and improving the available data is data augmentation. In this project we have tried different augmentation techniques such as Mixup, Fourier transforms and more ordinary manipulations such as rotations and brightness adjustments on top of a decently performing CNN as well as some experiments on ResNet50. We have tried to see under which circumstances the augmentations have effect and how big that effect is. The data set used consists of about 25 000 colored images of 200 bird species, mostly from the United States. The data has further been divided into subsets of different size and difficulty. Our results show that the Fourier transform does not seem to have a positive effect on the accuracy. We also see that basic data augmentations techniques such as flips and rotations seem to increase the performance under all circumstances, whilst Mixup and random erasing which manipulates the data seem to require a sophisticated network such as ResNet50 where they performed well, or they might impact the performance in a negative manner.

1 Introduction

One of the major drawbacks of supervised learning is the need of immense amounts of labeled data, producing such in the quantity needed for optimal results can be both costly and extremely time demanding, if even possible at all, which may be the case with medical data where only a limited number of known and active cases may exist. (Shorten & Khoshgoftaar, 2019).

Data augmentation provides a partly solution to this issue when it comes to image classification and after being briefly introduced to some data augmentation techniques during the lectures we wanted to explore this topic further and investigate the payoff (or the lack thereof) from some data augmentation techniques as well as getting an understanding of when it is worth the effort. We are interested in interpretability and have used this project to in an experimental way find the how, when and why.

Our main focus has been to try out Mixup, Random erasing and Fourier transformation in practice, but we have also chosen to apply some more standard augmentations such as rotations simple adjustments to have something to compare with.

To clarify, our project does not aim to obtain the highest possible testing accuracy but rather aims to show the impact of data augmentations on the accuracy.

For our experiments we created a basic CNN architecture as well as implemented a fully trained version of ResNet50 in order to also examine the effects on data augmentations techniques in the case of transfer learning.

To create different settings we divided the original data into different subsets with different sizes and characteristics, some random and some more targeted such as birds with bright colors.

2 Related Work

2.1 Basic Data Augmentation

A survey on Image Data Augmentation for Deep Learning by Shorten & Khoshgoftaar(2019) is a profound survey paper providing guidelines for different data augmentations and applications. The power of data augmentations is a red line throughout the text but the authors also emphasize the need of choosing the right augmentations for the specific dataset, and points out that in situations with very little data, augmentations may even lead to further overfitting, which is the direct opposite of its purpose. This is something we would like to investigate further, as being aware of problems may be as important as advantages.

2.2 Mixup

The concept of Mixup was introduced by Zhang et al(2018) and was published as a conference paper at ICLR 2018. Mixup is really fascinating because of its' creativeness and ability to improve performance while being very simple. As a motivation to the augmentation the authors mention that networks often pay too much attention to contradictory cases and tend to memorize instead of paying attention to the general features. Mixup aims to confuse the network enough for it to focus more on the general and essential parts and has proven to be successful on for example CIFAR 10 and CIFAR 100.

It creates virtual training samples by combining two images by

$$\begin{aligned}\tilde{x} &= \lambda x_i + (1 - \lambda)x_j, & \text{where } x_i, x_j \text{ are input vectors} \\ \tilde{y} &= \lambda y_i + (1 - \lambda)y_j, & \text{where } y_i, y_j \text{ are one-hot encoded labels}\end{aligned}$$

where (x_i, y_i) and (x_j, y_j) are two randomly drawn samples from the training data, and λ is a probability, that is $\lambda \in [0, 1]$. Some examples where Mixup where performed is shown in Figure 2. Usually, λ is randomly drawn from a Beta(α, α) distribution, for each pair of images which are to be combined. Examples and further applicational details are described under Methods.

2.3 Random Erasing

Random erasing was first introduced by Zhong et al (2017). It is another photometric augmentation which puts a random sized (within some range) rectangle over some randomly chosen part of the image with some probability, see Figure 3. This has lead to "consistent improvements", for example on CIFAR10 and CIFAR100. A mentioned advantage, aside from that it prevents overfitting, is that it makes the network more robust to occlusion, which often happens when the picture is taken in lively environments.

3 Data

In this project we have worked with a dataset consisting of images of different species of birds, the Bird Species Dataset published on Kaggle. Each image has the format $224 \times 224 \times 3$ and the images are cropped such that the bird covers at least 50% of the pixels. Most pictures include the whole body of the bird. There are a total of 200 species and the training data consist of 27503 images. The data is not balanced but does however contain at least 100 training images for each species. Both the validation set and the test set consists of 5 images of each species. It should also be said that around 80% of the images are of male birds and 20% of female birds which, by the nature of birds, may look entirely different, which has sometimes caused some trouble when the data set has been used.

We will not always work with the full dataset but instead use the following subsets:

- All birds (200 species)
- Bright colored birds (84 species)
- Dull colored birds (41 species)

The dull and bright colored species have been determined by inspection. For each subset there is a random selection of 15 birds and a small and medium sized subset, giving the number of training images per bird species, not including augmented images. We have used 5 for small and 50 for medium, giving 6 subsets in total, 7 with the complete data.

4 Methods

To get a deeper understanding of the impacts different kinds of data augmentations we have chosen four different categories: *Basic Data Augmentations* which consists of transformations such as rotations and flips, *Mixup*, which blends images and labels for different images, *Random Erasing*, which erases a random part of an image and lastly *Fourier Transforms*, which we consider as an interesting experiment.

4.1 Basic Data Augmentations

The idea behind Data Augmentation is to increase the relevant training data using the available data. This is obtained by manipulating the image such that it appears to be a new image with new valuable information to the network. For this project we have chosen some of the most common and simple Data Augmentation techniques to see what difference simple changes can make (or not), and to be able to compare with more sophisticated methods.

Figure 1 demonstrates some simple augmentations, there are 4 versions of the same picture where rotation, flipping, brightness and shearing (a kind of stretching) have been applied. Other common techniques are location shifts and zooming, but since the pictures are cropped such that at least 50% of the pictures are covered by the bird, applying any of those techniques often lead to the bird being beheaded or losing a big body part, which might not be a relevant image in this dataset as most images appear to be professional photos with the bird centered and having its whole body in the picture.



Figure 1: Rotation, shearing, flipping and brightness adjustments applied to a picture

4.2 Mixup

As mentioned under Related Work, Mixup creates virtual training example by combining two images by

$$\begin{aligned}\tilde{x} &= \lambda x_i + (1 - \lambda)x_j, & \text{where } x_i, x_j \text{ are input vectors} \\ \tilde{y} &= \lambda y_i + (1 - \lambda)y_j, & \text{where } y_i, y_j \text{ are one-hot encoded labels}\end{aligned}$$

where $\lambda \sim \text{Beta}(\alpha, \alpha)$. This distribution seems like a reasonable choice since it has the right support(0 to 1) and the Beta-distribution is among the most natural distributions to consider when working with probabilities. An example of Mixup performed on images can be seen in Figure 2. The $\text{Beta}(\alpha, \alpha)$ -distribution is also symmetric around 0.5, which may be a desirable property, however this should not matter since we combine our randomly drawn examples with weights λ and $1 - \lambda$ and it should not matter if, for example $\lambda = 0.2$ or $\lambda = 0.8$ since it would yield the same two weights, but in different order, but since our examples are randomly drawn the order of the weight should not have an impact. The above labeling is what we call fractional labeling. We will also consider performing Mixup on the input vectors of the training images but letting \tilde{y} keep the label of the example with the highest weight. That is,

$$\tilde{y} = I_{\{\lambda \leq 0.5\}} y_i + (1 - I_{\{\lambda \leq 0.5\}}) y_j,$$

where $I_{\{\lambda \leq 0.5\}} = 1$ if $\lambda \leq 0.5$ and 0 otherwise. This we will call majority vote labeling.



Figure 2: Three examples of Mixup performed on images of different bird species.

4.3 Random Erasing

Random erasing is a bit similar to Mixup in the sense that it aim at confusing the network a little by manipulating the image, hopefully just enough for it to see the more general traits and overfit and memorize less.

We wish to compare Mixup with Random Erasing and see which has the strongest effect on the birds dataset and if they may be used together.

We have chosen to let the box be different nuances of gray, so that the network does not focus on the color. Random colorization of the pixels in the box has proven to often be the most successful, however we thought that in this dataset using more colors might confuse, and also gray(white-black) was easier to implement.

We have constrained the box to cover between 10 and 20 percent of the image and let the probability of a box being placed be 0.3.

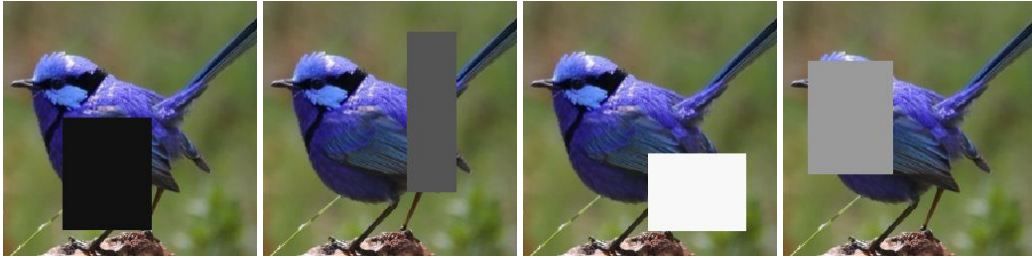


Figure 3: Four different random erasings for an image

4.4 Fourier Transformation

The 2D Fourier transform is used in for instance in image compression and gave us the idea of the possibility using it in deep learning and image classification. The 2D fast Fourier transform takes our image, with size $[224 \times 224 \times 3]$, as input and returns a complex valued matrix of the same size. From this output we can now get the amplitudes by taking the absolute value of each element in the matrix and the phase angles by computing $\arctan(y, x)$, where x is the real part and y is the imaginary part. An interesting property is that the phase angles are more important than the amplitudes and contain more information necessary to recreate the image again. As an example, in Figure 4 the Fourier transform was applied to two images which yielded amplitudes and phase angles respectively. The images were then recreated using the amplitudes from the first image and the phase angles from the second image and vise versa.

Using this as the inspiration we wanted to see if it is true that the phase angles contain more essential information about the image by applying image classification on the amplitudes and phase angles respectively to see if our hypothesis of higher accuracy on the phase angles is satisfied or not.

Initially we had no expectations that using the Fourier transform as augmentation would yield an increase the classification accuracy compared to using the raw data, however we thought that it would be an interesting experiment. In Figure 5 we can see an example of the Fourier transformation on one of the images in our dataset.

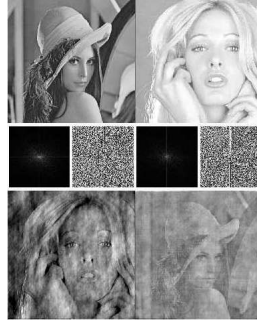


Figure 4: Example of recreation of images using amplitudes and phase angles of the Fourier transformation. Bottom left image use the amplitudes from the first image and the phase angles from the second image and bottom right image use the phase angles from the first image and the amplitudes from the second image

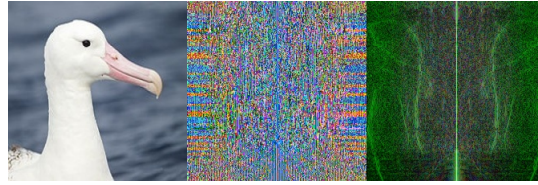


Figure 5: Left: Original image. Middle: Phase angles. Right: Amplitudes.

4.5 Underlying CNNs

In order to examine the effects of our data augmentation methods we wanted to observe them for some different networks. Due to restrictions in time and computational resources we were only able to focus on two. Below are the networks we examined.

- (i) Basic CNN (about 65 % accuracy on the whole data set)
- (ii) ResNet50 trained on ImageNet (about 95% accuracy on the whole data set)

For the Basic CNN the structure of the layers are shown in Figure 6.

We have chosen CNNs since they have shown to perform good, sometimes extremely good, on computer vision tasks. Batch Normalization has been a necessity in order to manage the gradients and to keep the network stable. We have further used dropout to avoid overfitting which we have anticipated to encounter, both considering the fact that we have 85 million parameters and also because we will try the network on subsets with very few training examples. Overall we have tried out some standard layer structures until reaching about 60% accuracy on the training data set.

We also chose to use ResNet50 in order to observe how the data augmentations techniques alter the performance in the case of transfer learning. In the case of ResNet50 we simply flattened the output added our own predictive layer in order to accomodate to our data.

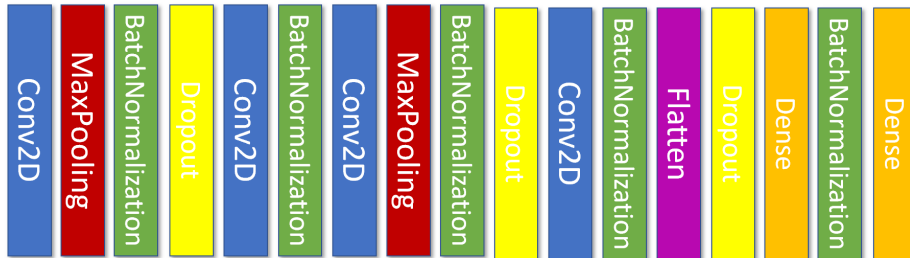


Figure 6: Layers of the Basic CNN

The networks were trained for 15 epochs in each experiment using a batch size of 100 and a learning rate of 0.01, in combination with momentum (0.9) and some decay ($1e-6$). We also implemented early stopping in order to cut down on the training times. The choices of hyperparameters were made by using generally recommended values and the slightly tuning them to fit our needs. As the aim of this project is not to maximize the performance of the network but rather to examine the effects of certain methods the search for good hyperparameters could be seen as very coarse.

5 Experiments

Recall that there are the same number of species in all datasets(15), what differs is the number of training images(excluding augmentations) for each species(5 or 50).

It should be said that each combination have only been run once and thus we do not know anything about the span of variation of the estimates. It would also have been desirable to try out different randomized selections of species and to test ResNet50 without pretraining, but that was far too time consuming for this project.

5.1 Basic Augmentation

In table 1 the results are shown for the runs with the CNNs on all subsets(except the full datasets is which is excluded since the medium ones have performed very well already) with and without data augmentations. For ResNet50, where the Small datasets performed extremely very well, we did not run anything for the Medium sized sets for that reason.

For all subsets for both our basic CNN and ResNet50 we can see a great gain in accuracy when applying data augmentations. The difference is the biggest for the basic CNN, but improvements for already good performing networks are hard earned and it is really good that such improvements can be made for an already such good performing network with such little effort.

The basic CNN had quite some problems with the RandomSmall dataset and seems to overfit, here the augmentations made the most difference. The ResNet50 did not have such problems, maybe because it is pretrained or was better at spotting the relevant parts. The ResNet50 also perform more evenly across different datasets.

Table 1: Accuracies for Basic Augmentations on Basic CNN and ResNet50

Data	No Aug Basic CNN	Basic Aug Basic CNN	No Aug ResNet50	Basic Aug ResNet50
RandomSmall	0.13	0.47	0.81	0.85
RandomMedium	0.65	0.79	0.85	not tested
BrightSmall	0.40	0.53	0.79	0.84
BrightMedium	0.84	0.90	0.99	not tested
DullSmall	0.27	0.38	0.83	0.89
DullMedium	0.57	0.81	0.99	not tested

As the Random datasets perform very similarly to the Dull ones we continue only with the Bright- and Dull subsets in our further experiments and only run one test on the ResNet50 which performs very evenly.

5.2 Mixup

In order to examine the effects of Mixup we created our own data generator which used two copies of our main data generator and combined the generated images in order to create the Mixup images as displayed in Figure 2.

In Table 2 below we showcase the results of a variety of parameter combinations for Mixup on the basic CNN and ResNet50. Note that the first row indicates the training without Mixup.

Table 2: Mixup with a Beta distribution on Basic CNN and ResNet50

Parameter	Label	Accuracy Basic CNN(%)	Accuracy ResNet50(%)
-	-	64.59	94.78
0.1	Majority	X	X
0.1	Fractional	X	X
0.2	Majority	58.29	92.14
0.2	Fractional	54.17	82.21
0.5	Majority	54.63	X
0.5	Fractional	53.54	X

Table 3: Mixup

Data	Acc Basic CNN	Acc ResNet50
BrightSmall	0.31	0.80
BrightMedium	0.73	0.97
DullSmall	0.31	0.82
DullMedium	0.57	0.99

5.3 Random Erasing

Table 4: Random Erasing

Data	Acc Basic CNN	Acc ResNet50
BrightSmall	0.34	0.89
BrightMedium	0.83	0.97
DullSmall	0.20	0.89
DullMedium	0.51	0.99

5.4 The Fourier Transform

In order to examine which components from the Fourier transform, that is the angles, the amplitudes or the combination of them, which generates the best results we examined their respective performances using the basic CNN architecture and the whole dataset. It should be noted that the training time increases drastically when taking both the angles and the amplitudes into account as the dimensions of the input increases.

Table 5: Basic CNN and Fourier Transform

Data	Accuracy (%)
Raw Images	64.59
Fourier Angles	28.84
Fourier Amplitudes	48.42
Fourier Angles & Amplitudes	47.42

For the basic CNN it seems as if using only the amplitudes generated the highest accuracy which seems to contradict what we initially thought about the phase angles containing more information about the content of the image. In all, it seems as if the Fourier transform of the input data is not a valid method in order to improve the performance of a CNN, or at least not in the case of our bird data. An interesting extension could be to apply this method to different datasets with more internal variety, such as for example CIFAR 100.

Due to computational constraints as well as the quite clear results from the performance on the Basic CNN we did not evaluate the effects of the Fourier transform on other datasets as well as for

ResNet50. A major reason for this is that the pre-trained weights from ImageNet are in no way compatible with our angles and amplitudes which would imply a complete re-training using the base architectures.

6 Conclusions

In our experiments we have learned that some basic augmentations such as rotations and flipping can make a large difference in the accuracy both when the dataset is performing poorly and good already.

SKRIV RESULTAT FÖR RANDOM ERASING

Our results for the Mixup did however disappoint since the accuracy got worse when performing Mixup compared to not. The reason for this is something we would like to investigate further since Mixup has been proven to be an effective method for others. Our guess is that it is because of the nature of our somewhat homogeneous data, the most birds have similar shape and look alike, and thus Mixup will not be as effective as if, for example, performed on CIFAR10. Another hypothesis is that Mixup may need a quite sophisticated network and a lot of data, ResNet50 might be such a network but here it performed good from start and did not leave much room for improvements.

What we have implemented is *input Mixup* where we do all augmentation before training and input the images in the network after performing Mixup, however Verma et al (2019) suggest a new algorithm called *Manifold Mixup* where Mixup is performed at an intermediate layer or final layer in the network. In an intermediate layer the feature spaces are more aligned than that of the input and it is suggested that Mixup will produce better augmented data if the Mixup is performed on this layer. Therefore, this would be a natural possible future extension to this project.

For our experiment with the Fourier transformation we learned that it was very ineffective as a possible data augmentation method and that our hypothesis that the results would be better when using the phase angles than the amplitudes were incorrect in this application.

Finally, all our experiments could use some further optimization in the parameter choices, batch size and number of epochs and also some different tests and on a different dataset, however due to the limitations in time and computational power, this was not possible.

References

Shorten & Khoshgoftaar (2019) A survey on Image Data Augmentation for Deep Learning, *Journal of Big data* 6

Verma et al (2019) Manifold Mixup: Better Representations by Interpolating Hidden States arxiv.org/pdf/1806.05236.pdf

Zhang et al (2018) Mixup: Beyond Empirical risk minimization *ICLR conference paper 2018*

Zhong et al (2017) Random Erasing Data Augmentation <https://arxiv.org/pdf/1708.04896.pdf>