
PatchShuffle Regularization

Guoliang Kang, Xuanyi Dong, Liang Zheng, Yi Yang
Center of AI, University of Technology Sydney, Australia
{Guoliang.Kang@student., Xuanyi.Dong@student., Liang.Zheng@, Yi.Yang@}uts.edu.au

Abstract

This paper focuses on regularizing the training of the convolutional neural network (CNN). We propose a new regularization approach named “PatchShuffle” that can be adopted in any classification-oriented CNN models. It is easy to implement: in each mini-batch, images or feature maps are randomly chosen to undergo a transformation such that pixels within each local patch are shuffled. Through generating images and feature maps with interior orderless patches, PatchShuffle creates rich local variations, reduces the risk of network overfitting, and can be viewed as a beneficial supplement to various kinds of training regularization techniques, such as weight decay, model ensemble and dropout. Experiments on four representative classification datasets show that PatchShuffle improves the generalization ability of CNN especially when the data is scarce. Moreover, we empirically illustrate that CNN models trained with PatchShuffle are more robust to noise and local changes in an image.

1 Introduction

The trend of the architectures of deep convolutional neural networks (CNNs) is to become wider [28, 29] and deeper [6, 7, 23, 26]. However, millions of parameters make CNNs prone to overfitting when training data is not sufficient. In practice, plenty of regularization approaches have been adopted to improve the generalization ability of CNNs, such as weight decay [15], dropout [8, 25], batch normalization [10], *etc.* This paper provides an alternative regularization option during CNN training with application in image classification.

Overfitting is a long-standing issue in the machine learning community. The nature of overfitting is that the model adapts to the noise rather than capturing the underlying key factors of variations existing in the data [30]. For image classification, when lacking sufficient training data, the learned model may be misled by the irrelevant local information which can be regarded as noise. Moreover, in most classification tasks, it is the overall structure rather than the detailed local pixels that has a large influence on the performance of a CNN model. The model should be able to identify the input image correctly if pixels within local structures change in a way that does not destroy the overall view of an image. From another perspective, if we consider that human will probably not be confused about the the image content under moderate extent of local blur, it is expected that a data model such as CNN should behave similarly.

In this paper, we propose a new regularization approach: PatchShuffle, which is a beneficial supplement to existing regularization techniques [10, 14, 15, 25]. In the training stage, an image or feature map within a mini-batch is randomly chosen to undergo either of the two actions: 1) keep unchanged, or 2) be transformed in such a way that pixels within each patch are shuffled. On the one hand, when applied on the images, the shuffled images have nearly the same global structures with the original ones but possess rich local variations, which are expected to benefit the training of CNNs. On the other hand, when PatchShuffle is applied on the feature maps of the convolutional layers, it can be viewed as implementing model ensemble. In fact, locally shuffling the pixels within a patch is equivalent to shuffling the convolutional kernels given unshuffled patches. Thus at each iteration,

the model is trained from different kernel instantiations. PatchShuffle can also be considered to enable weight sharing within each patch. By shuffling, the pixel instantiation at a specific position of an image can be viewed as being sampled from its neighboring pixels within a patch with equal probability. Therefore, across different iterations, patch pixels with different original locations share the same weight.

One might argue that PatchShuffle is a type of data augmentation technique since new images are generated. However, being applied on *a very small percent* of images/feature maps in a mini-batch, we speculate that PatchShuffle is more of a regularization method than data augmentation¹. We also differentiate PatchShuffle from dropout. The latter samples activations from the hidden units, while PatchShuffle samples from all the possible permutations of pixels within patches and no hidden units are discarded.

In summary, the PatchShuffle regularization has the following merits.

- An efficient method that costs negligible extra time and memories. It can be easily adopted in a variety of CNN models without changing the learning strategy.
- A complementary technique to existing regularization approaches. On four representative classification datasets, PatchShuffle further improves the classification accuracy when combined with multiple regularization techniques.
- Improving the robustness of CNNs to data that is noisy or losses partial information. For example, when adding salt-and-pepper noise to the MNIST dataset, our approach outperforms the baseline by more than 20 percent.

2 Related Work

We briefly review several aspects that are closely related to this paper, *i.e.*, data augmentation, regularization, and transformation equivariant and invariant networks.

Data augmentation. The direct strategy against overfitting is to train CNNs on more data. Data augmentation addresses this problem by creating new data from existing data to augment the training set. Data augmentation is widely adopted in the training of deep neural networks [3, 5, 6, 14, 18]. An effective way to perform data augmentation is to do various transformations, such as flipping, translation, cropping, *etc.* From the perspective of generating more images for training, PatchShuffle shares some properties with data augmentation methods.

Regularization. Regularization is an effective way to reduce the impact of overfitting. Various types of regularization methods have been proposed [8, 10, 15, 24, 25, 27]. PatchShuffle relates to two kinds of regularizations. 1) Model ensemble. It adopts model averaging in which several separately trained models vote on the output given a test sample. The voting procedure is robust to prediction errors made by individual classifiers. Many methods implicitly implement model ensemble, such as dropout [8, 25], stochastic depth [9] and swapout [24]. Architectures are averaged by dropout through randomly discarding a group of hidden units, each of which has different widths of layers. Stochastic depth averages architectures with various depths through randomly skipping layers. Swapout samples from abundant set of architectures with dropout and stochastic depth as its special case. 2) weight sharing. It forces a set of weights to be equal [20] and has been used in the architecture of deep convolutional neural networks [17]. Networks regularized by weight sharing always have transformation invariant properties. For example, through a weight sharing framework, Ravanbakhsh *et al.* [21] propose the permutation equivariant layer that gains robustness to permutations of the input. PatchShuffle is more of a regularization method because the generated images/feature maps share the global structures with the original ones and PatchShuffle is applied on a very small amount of images/feature maps.

Transformation equivariant and invariant networks. PatchShuffle regularization is also related to the family of transformation equivariant and invariant networks. Deep symmetry networks [4] generalize vanilla CNN architecture to model arbitrary symmetry groups. In [2], a series of rotation equivariant operations are proposed, such as cyclic slicing, pooling and rolling. Ravanbakhsh *et al.*

¹In [16], data augmentation is considered as belonging to regularization. We differentiate the two concepts in this paper: data augmentation enlarges the training set to a large extent, while regularization makes more elaborate data changes without noticeable enlarging the data volume.

[21] propose a kind of permutation equivariant layer that is robust to the permutations of the inputs through designed weight sharing. All of the aforementioned neural networks mainly focus on the transformations of the whole images and aim to enable the neural networks to be robust to several specific parametric transformation types. They are problem-driven and not easily generalized to other datasets. Moreover, few investigate and exploit various kinds of transformations in the regularizing of deep neural networks in a general sense.

In a recent work, Shen *et al.* [22] propose using patch reordering to achieve the rotation and translation invariance. They divide the feature maps into non-overlapping local patches and reorder the patches according to the ℓ_1 or ℓ_2 norm of the activations of the patches. Their work is similar to PatchShuffle in that they also break the original arrangement of an image or feature maps during training, but critical differences should be clarified. 1) Shen *et al.* [22] reorder the *patches*, while we shuffle the *pixels* within each local patch, which does not destroy the global structure. 2) Shen *et al.* [22] perform *ranking* according to specific heuristic rule, while PatchShuffle does the *shuffle* operation randomly. Randomness is proved to be useful to regularize the training of CNNs by explicitly performing model averaging [8, 9, 14, 24, 25]. 3) Shen *et al.* [22] concentrate on the rotation and translation invariance of models, whereas we adopt PatchShuffle as a regularizer.

3 PatchShuffle Regularization

3.1 PatchShuffle Transformation

Formulations. Let us consider a matrix \mathbf{X} with the size of $N \times N$ elements. A random switch r controls whether \mathbf{X} needs to be transformed (PatchShuffled). Supposing the random variable r subjects to a Bernoulli distribution $r \sim \text{Bernoulli}(\epsilon)$, i.e., $r = 1$ with probability ϵ and $r = 0$ with probability $1 - \epsilon$, the resulted matrix $\tilde{\mathbf{X}}$ can be represented as

$$\tilde{\mathbf{X}} = (1 - r)\mathbf{X} + rT(\mathbf{X}), \quad (1)$$

where $T(\cdot)$ denotes the PatchShuffle transformation. When \mathbf{X} is partitioned into a block matrix with non-overlapping patches of $n \times n$ elements, i.e.,

$$\begin{pmatrix} \mathbf{x}_{11} & \mathbf{x}_{12} & \cdots & \mathbf{x}_{1,N/n} \\ \mathbf{x}_{11} & \mathbf{x}_{12} & \cdots & \mathbf{x}_{1,N/n} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_{N/n,1} & \mathbf{x}_{N/n,2} & \cdots & \mathbf{x}_{k,N/n} \end{pmatrix}, \quad (2)$$

the PatchShuffle transformation acts on each patch and can be formulated as follows,

$$\tilde{\mathbf{x}}_{ij} = \mathbf{p}_{ij} \times \mathbf{x}_{ij} \times \mathbf{p}'_{ij}, \quad (3)$$

where \mathbf{x}_{ij} denotes a patch located at the i -th row and j -th column of the block matrix \mathbf{X} . \mathbf{p}_{ij} and \mathbf{p}'_{ij} are permutation matrices. Pre-multiplying the patch \mathbf{x}_{ij} with \mathbf{p}_{ij} permutes the rows of \mathbf{x}_{ij} , whereas post-multiplying the patch \mathbf{x}_{ij} with \mathbf{p}'_{ij} results in the permutation of the columns of \mathbf{x}_{ij} .

In practice, we first split \mathbf{X} into non-overlapping patches with sizes of $n \times n$ elements. Within each $n \times n$ patch, the elements are randomly shuffled, as shown in Fig. 1. So each patch will undergo one of the $n^2!$ different permutations. For matrix \mathbf{X} , the number of possible permutations is $\prod_{i=1}^{N^2/n^2} n^2!$. For each patch, after finite times of shuffle, it will recover the original order. Note that although we describe the PatchShuffle transformation assuming \mathbf{X} and \mathbf{x}_{ij} are square, in practice, \mathbf{X} and its patches \mathbf{x}_{ij} don't need to be square. Our method can be trivially extended to the case of non-square matrix.

PatchShuffle on images. For CNNs, the input and output of a convolutional layer are feature maps (or images) which can be viewed as matrices. Thus the PatchShuffle transformation is readily applicable. Following Eq. (1), Eq. (2) and Eq. (3), PatchShuffle can be easily applied on images. Image samples are shown in Fig. 2. Intuitively, we observe that for an image, PatchShuffle transformation within an extent does not disable the recognition of corresponding object, which will benefit the training of a deep neural network.

PatchShuffle on feature maps. We also perform PatchShuffle transformation on the feature maps of the *all the convolutional layers*. Here, we treat each feature map as an image, and PatchShuffle

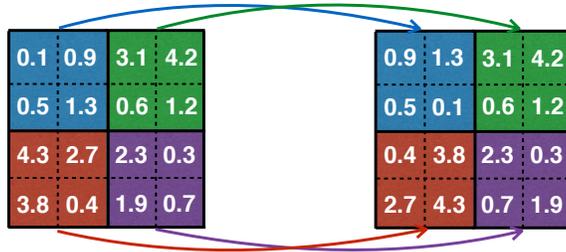


Figure 1: A cartoon illustration of PatchShuffle on a 4×4 matrix divided into four non-overlapping 2×2 patches (best viewed in color). Different patches are labeled with different colors. The shuffling of patches are independent from each other. Note that there is a possibility that pixels within a patch exhibit the original orders after shuffling, as illustrated in the upper-right green patch.



Figure 2: Samples of PatchShuffled real-world images. The original images are sampled from the STL-10 dataset [1]. The numbers above the shuffled images denote the patch size $n \times n$ (non-overlapping).

is performed on the feature maps independently. That is, each feature map is randomly chosen to undergo the PatchShuffle transformation, regardless of the original image or other feature maps. For the feature maps of lower or middle layers, the spatial structures of the image are preserved to a large extent, so we expect that applying PatchShuffle to these layers can regularize training. For the higher convolutional layers, recall that PatchShuffle enables weight sharing among neighboring pixels; this property is beneficial for the higher level feature maps where neighboring pixels have largely overlapping receptive fields projected onto the original image. We will verify this in the experiment part.

Discussions. PatchShuffle also faces the typical bias-variance dilemma. On the side of reducing the gap between training and test performance, PatchShuffle creates new images and feature maps, which increases the variety of the training data. However, on the side of bias, the data distributions of the new images and feature maps created are probably different from those of real-world data, which may induce more bias into the CNN model. Therefore, in the application of PatchShuffle, only a small percentage of the images/feature maps undergo PatchShuffle transformation (small ϵ) to achieve a bias-variance trade-off.

3.2 CNN Training and Inference

Objective function. We take PatchShuffling images for example. In this case, the training objective function can be formalized as,

$$\ell_s(\mathbf{X}, y, \boldsymbol{\theta}) = (1 - r)\ell(\mathbf{X}, y, \boldsymbol{\theta}) + r\ell(T(\mathbf{X}), y, \boldsymbol{\theta}), \quad (4)$$

where ℓ_s and ℓ denote the objective functions training with and without PatchShuffle, respectively. \mathbf{X} and $T(\mathbf{X})$ represent the original and PatchShuffled images, respectively. The label of the training sample is denoted as y , and $\boldsymbol{\theta}$ encode the weights of the neural network.

In Eq. (4), when the random switch r is set to 1, we have $\ell_s(\mathbf{X}, y, \boldsymbol{\theta}) = \ell(T(\mathbf{X}), y, \boldsymbol{\theta})$, which implies that network is chosen to be trained with PatchShuffle. When $r = 0$, we have $\ell_s(\mathbf{X}, y, \boldsymbol{\theta}) = \ell(\mathbf{X}, y, \boldsymbol{\theta})$, denoting that the network is trained without PatchShuffle. Taking the expectation over r

which follows a Bernoulli distribution, Eq. (4) becomes

$$\frac{1}{1-\epsilon} \mathbb{E}_r(\ell_s(\mathbf{X}, y, \boldsymbol{\theta})) = \ell(\mathbf{X}, y, \boldsymbol{\theta}) + \frac{\epsilon}{1-\epsilon} \ell(T(\mathbf{X}), y, \boldsymbol{\theta}), \quad (5)$$

where ϵ denotes the shuffle probability, and $\frac{\epsilon}{1-\epsilon} \ell(T(\mathbf{X}), y, \boldsymbol{\theta})$ works as a regularizer.

Training procedure. During the training process, PatchShuffle is applied to the training images and feature maps of all the convolutional layers, each with an independently sampled r and independent shuffling operations. Note that the sizes of non-overlap patches $H_p \times W_p$ and the shuffle probability ϵ are also not necessarily set as the same across layers. We use the same procedure when applying PatchShuffle to different layers. Without loss of generality, in Algorithm 1, we summarize the training procedure using PatchShuffle which is applied on the the feature maps of one convolutional layer.

Inference. At the test stage, the network performs a forward process *without any transformation* applied to the images or feature maps.

Algorithm 1 PatchShuffle on the feature maps of one convolutional layer in one iteration

Input:

- 1) feature maps $\mathbf{X}^{(c)}$ with spatial sizes $H \times W$, $c \in [1, C]$ where C denotes the number of channels of the feature maps in this convolutional layer.
- 2) shuffle probability ϵ , patch height H_p , and patch width W_p .

Output:

feature maps $\tilde{\mathbf{X}}^{(c)}$ with the same spatial sizes as $\mathbf{X}^{(c)}$, where $c \in [1, C]$.

Let $\mathbf{Z}^{(c)}$ represents the corresponding mapping between $\tilde{\mathbf{X}}^{(c)}$ and $\mathbf{X}^{(c)}$.

Forward process:

$h = \lceil H/H_p \rceil$, $w = \lceil W/W_p \rceil$

for $c = 1$ to C :

Generate the random switch r ;

if $r = 1$: shuffle each of the $h \times w$ patches of $\mathbf{X}^{(c)}$ and storing the mapping $\mathbf{Z}^{(c)}$

else: $\tilde{\mathbf{X}}^{(c)} = \mathbf{X}^{(c)}$

end

Backward process:

for $c = 1$ to C :

if $r = 1$: mapping from the gradients of $\tilde{\mathbf{X}}^{(c)}$ to those of $\mathbf{X}^{(c)}$ according to $\mathbf{Z}^{(c)}$.

else: copying from the gradients of $\tilde{\mathbf{X}}^{(c)}$ to those of $\mathbf{X}^{(c)}$

end

4 Experiments

In this section, we report results on four image classification datasets including CIFAR-10, SVHN, STL-10 and MNIST. CIFAR-10 [13] contains 50,000+10,000 (training+test) 32×32 color images of 10 object classes. SVHN [19] consists 73,257+26,032 (training+test) 32×32 color images for street view house numbers. STL-10 [1] contains 5,000+8,000 (training+test) 96×96 color images for 10 categories. MNIST [17] consists of 60,000+10,000 (training+test) 28×28 greyscale images of hand-written digits. We first show that our algorithm is robust to the change of hyper-parameters within a wide range. Then we demonstrates the improved generalization ability achieved by PatchShuffle on the benchmarks. Finally, we illustrate that CNNs trained by PatchShuffle are more robust to the noises such as salt-and-pepper, occlusions, *etc.* All experiments are implemented using *Caffe* [11].

4.1 Experiment Settings

In all of our experiments, we compare the CNN models trained with and without PatchShuffle. The training method without PatchShuffle is denoted as standard back-propagation (BP). For the same deep architecture, all the models are trained from the same weight initializations. Note that some popular regularization techniques (*i.e.*, weight decay, batch normalization and dropout) and various data augmentations (*i.e.*, flipping, padding and cropping) are employed in the experiments. The hyper-parameters for training with PatchShuffle and standard BP are all the same except that the

patch size $H_p \times W_p$ and shuffle probability ϵ are chosen through validation for PatchShuffle. Our experiments build on various CNN architectures, which are summarized as follows:

CNNs for CIFAR-10. Three CNN models are adopted in the experiments of CIFAR-10: Network in Network [18] (NIN), pre-activation ResNet-110 [7] (ResNet-110-PreAct), and the modification of original ResNet-110 [6] (ResNet-110-Modified). The architecture of ResNet-110-Modified is the same as the original ResNet-110 designed for CIFAR [6] except that it discards the ReLU unit after each summation of the shortcut and residual function. This small modification improves the performance of original ResNet to be comparable to that of pre-activation ResNet [7]. The training procedure of ResNet-110-PreAct and ResNet-110-Modified is the same with [6, 7], and that of NIN is the same with [18].

CNNs for SVHN. The CNNs adopted for SVHN are Plain-SVHN and ResNet-110-Modified. The architecture of Plain-SVHN is the same as provided in the *Caffe* examples² which is originally designed by [12] for the training of CIFAR-10. Because the architecture is general and the image resolutions are the same between SVHN and CIFAR-10, it can be employed for SVHN.

CNNs for STL-10. The CNN architecture adopted for STL-10 is denoted by ResNet-STL-10. It is similar to ResNet-110-Modified, but due to a higher resolution of the images, several modifications are made. 1) The kernel size for the first convolutional layer becomes 7 with the stride of 2. 2) Four residual stages are employed with each of them containing two residual units. The spatial sizes of feature maps are successively halved after each residual stage. The channels of the feature maps for four residual stages are 32, 64, 128 and 256, respectively. 3) Dropout is applied on the final fully-connected layer with dropout ratio of 0.5.

CNNs for MNIST. The CNN architecture adopted on MNIST is provided in *Caffe* examples³.

4.2 The Impact of Hyper-parameters

When applying PatchShuffle to CNN training, we have two hyper-parameters to evaluate, *i.e.*, the patch size $H_p \times W_p$ and the shuffle probability ϵ . To demonstrate the impact of these two hyper-parameters on the performance of the model, we conduct experiments on CIFAR-10 based on ResNet-110-Modified under different hyper-parameter settings with PatchShuffle applied on the images. The results are compared with standard BP and shown in percentage in Table 1. Note that all the models are trained with the simple data augmentation as in [6, 7]: 4 pixels are padded on each side, and a 32×32 crop is randomly sampled from the padded image or its horizontal flip. Results are presented in Table 1 and Fig. 3. We arrive at two findings.

Probability	Patch Size	std-BP	Patch Size					Min.
			1×2	2×2	2×4	3×3	4×4	
0.01		6.33	6.10	6.17	6.26	6.34	6.40	6.10
0.05			6.05	5.66	5.84	5.93	6.06	5.66
0.10			6.01	5.95	6.08	5.86	5.93	5.86
0.15			6.25	6.27	5.99	6.29	6.10	5.99
0.20			6.16	6.09	6.32	6.62	6.22	6.09
0.30			6.08	6.56	6.83	7.17	6.86	6.08
Min.		6.33	6.01	5.66	5.84	5.86	5.93	5.66

Table 1: Test errors (%) on CIFAR-10 based on ResNet-110-Modified under different hyper-parameters.

First, PatchShuffle consistently outperforms standard BP under a wide range of hyper-parameters. On CIFAR-10, our best result reduces the classification error by 0.67% compared with standard BP.

Second, PatchShuffle is robust to parameter changes to some extent. When the shuffle probability and patch size increase, recognition error first decreases, touches the bottom, and then increases. In fact, within an extent, the increase of both parameters improve the variety of training sample without introducing too much bias. But under larger values, the benefit brought by diversity is gradually

²https://github.com/BVLC/caffe/blob/master/examples/cifar10/cifar10_full_train_test.prototxt

³https://github.com/BVLC/caffe/blob/master/examples/mnist/lenet_train_test.prototxt

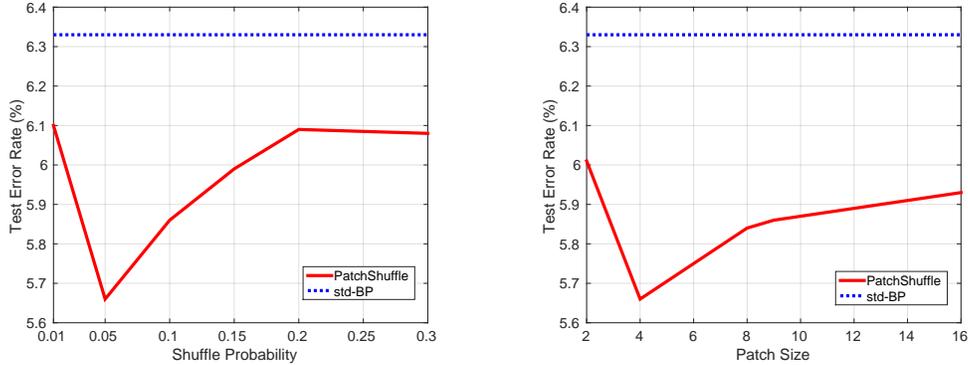


Figure 3: Test error (%) vs. different hyper-parameters. **Left**: shuffle probability. **Right**: patch size. We adopt the minimum value of test errors among different settings of one kind of hyper-parameters to represent the performance of another.

overtaken by the classifier bias, so error rate increases. In the following experiments, we use $\epsilon = 0.05$, patch size = 2×2 when not specified.

4.3 Classification Performance

CIFAR-10. For the training of NIN and ResNet-110-Modified, we apply PatchShuffle on the images only, while for the training of ResNet-110-PreAct, PatchShuffle is applied on the images and feature maps between two successive convolutional layers in each residual unit. The test errors are shown in percentage in Table 2. It can be seen that the models trained by PatchShuffle are consistently superior to those trained by standard BP with using the three CNN architectures.

We further evaluate the impact of the size of training set on recognition accuracy in CIFAR-10. We use the ResNet-110-Modified model. Results are shown in Table 2. In all of the experiments, we use the same hyper-parameter setting (*i.e.*, with patch size 2×2 , and shuffle probability 0.05). Although the hyper-parameter setting may not be optimal under small training sets, Table 2 still indicates that PatchShuffle improves the recognition accuracy, especial when the training set is small (see the results when training set size is 9,000).

Model	Size	
	std-BP	PatchShuffle
NIN[18]	10.43	10.09
ResNet-110-PreAct[7]	6.37	5.82
ResNet-110-Modified	6.33	5.66

Size	std-BP	
	std-BP	PatchShuffle
9,000	22.12	18.20
15,000	14.26	13.76
24,000	9.83	9.64
40,000	7.17	6.66
50,000	6.33	5.66

Table 2: Test errors (%) on CIFAR-10. **Left**: different CNN architectures. **Right**: different sizes of training data.

SVHN. PatchShuffle is applied on the images. The results are shown in Table 3, indicating that PatchShuffle consistently outperforms standard BP.

STL-10. Here we illustrate the impact of applying PatchShuffle on the feature maps of CNNs.

Model	std-BP	PatchShuffle
Plain-SVHN	6.24	5.85
ResNet-110-Modified	4.73	4.11

Table 3: Test errors (%) on SVHN. No data augmentation is adopted. No other data preprocessing is performed other than per-pixel mean computed over the training set is subtracted from each image.

The classification results on STL-10 are summarized in Table 4. The five-bit binary code denotes on what stages PatchShuffle is applied. The first bit denotes the input layer, and the other four bits

correspond to four residual stages. Applying PatchShuffle on a stage of ResNet means applying it on the feature maps between two adjacent convolutional layers of each residual unit in this stage. We set ϵ to 0.30 in this experiment.

Table 4 reveals that the generalization ability of the CNN models trained with PatchShuffle are significantly higher than using standard BP. More significant improvement over the baseline can be observed when using PatchShuffle on more convolutional layers. In addition, increasing the sizes of the patches also brings notable improvement in terms of the generalization performance. Note that all the models are trained *with dropout* applied on the output layer, which suggests that PatchShuffle can reduce overfitting beyond dropout.

Patch Size \ Layers	Layers					
	std-BP	10000	11000	11100	11110	11111
2×2	50.42	47.49	43.57	40.66	36.08	35.07
4×4		45.85	41.26	37.26	33.75	33.16
6×6		43.04	40.80	36.99	33.45	33.19

Table 4: Test errors (%) on STL-10 with different patch sizes and different choices of layers on which PatchShuffle is applied.

τ_1	std-BP*	PS*	std-BP	PS	τ_2	std-BP	PS
0.1	11.78	2.67	1.44	1.17	0.05	4.55	4.40
0.3	48.72	24.15	3.32	2.67	0.10	11.40	10.77
0.5	78.41	54.67	10.02	8.39	0.15	22.83	20.97
0.7	88.88	75.61	37.34	28.85	0.20	36.26	33.12
0.9	90.24	87.37	83.87	76.26	0.25	51.04	47.58

Table 5: Test errors (%) on MNIST under different levels of noise and occlusions. ‘‘PS’’ indicates PatchShuffle. All the test sets are polluted with noise or occlusions. ‘‘*’’ means that the training set is not polluted. None ‘‘*’’ denotes that the training set is also polluted with the same level of noise/occlusions with the test set. **Left:** salt-and-pepper noise. **Right:** occlusions.

4.4 Robustness to the Noise

Finally, we show the robustness of PatchShuffle against noise and occlusions. In experiment, we add different levels of salt-and-pepper noise and occlusions on the MNIST dataset. Salt-and-pepper noise is added to the image by changing the pixel to white or black with probability τ_1 . For the occlusion, each pixel is randomly chosen to be imposed by a black block of certain size centered on it with probability τ_2 . The size of the block adopted in our experiment is 3×3 . Results are presented in Table 5.

A clear observation is that under increasing level of pollution, the performances of both standard BP and PatchShuffle drop quickly. Nevertheless, under each pollution level, our method yields consistently lower error rate than standard BP. For salt-and-pepper noise, the performance gap is largest (23.74%) under a noise level of 50%. For occlusion, our method exceeds standard BP by 3.46% under occlusion extent of 25%. These results indicate that Patchshuffle improves the robustness of CNNs against common image pollutions like noise and occlusion.

5 Conclusion

This paper introduces PatchShuffle, a new regularization method for generalizable CNN training. This method is efficient to compute, complementary to existing regularizers, and improves CNN’s robustness to noise and occlusions. We will explore Patchshuffle on more complex tasks in future, *e.g.*, object detection and language modeling.

References

- [1] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 215–223, 2011.
- [2] Sander Dieleman, Jeffrey De Fauw, and Koray Kavukcuoglu. Exploiting cyclic symmetry in convolutional neural networks. 2016.
- [3] Zhe Gan, Ricardo Henao, David Carlson, and Lawrence Carin. Learning deep sigmoid belief networks with data augmentation. In *Artificial Intelligence and Statistics*, pages 268–276, 2015.
- [4] Robert Gens and Pedro M Domingos. Deep symmetry networks. In *NIPS*, 2014.
- [5] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *ECCV*, 2016.
- [8] Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. 2012.
- [9] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q. Weinberger. Deep networks with stochastic depth. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part IV*, pages 646–661, 2016.
- [10] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- [11] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. In *ACM MM*, 2014.
- [12] A Krizhevsky. cuda-convnet, 2012. <https://code.google.com/p/cuda-convnet/>.
- [13] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Citeseer, 2009.
- [14] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [15] Anders Krogh and John A Hertz. A simple weight decay can improve generalization. In *NIPS*, 1991.
- [16] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. volume 521, pages 436–444. Nature Research, 2015.
- [17] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. volume 86, pages 2278–2324. IEEE, 1998.
- [18] Min Lin, Qiang Chen, and Shuicheng Yan. Network in network. In *ICLR*, 2014.
- [19] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, volume 2011, page 5, 2011.
- [20] Steven J Nowlan and Geoffrey E Hinton. Simplifying neural networks by soft weight-sharing. volume 4, pages 473–493. MIT Press, 1992.
- [21] Siamak Ravanbakhsh, Jeff Schneider, and Barnabas Poczos. Deep learning with sets and point clouds. 2016.

- [22] Xu Shen, Xinmei Tian, Shaoyan Sun, and Dacheng Tao. Patch reordering: A novel way to achieve rotation and translation invariance in convolutional neural networks. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA.*, pages 2534–2540, 2017.
- [23] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [24] Saurabh Singh, Derek Hoiem, and David A. Forsyth. Swapout: Learning an ensemble of deep architectures. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 28–36, 2016.
- [25] Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. volume 15, pages 1929–1958, 2014.
- [26] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–9, 2015.
- [27] Lingxi Xie, Jingdong Wang, Zhen Wei, Meng Wang, and Qi Tian. Disturblabel: Regularizing cnn on the loss layer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4753–4762, 2016.
- [28] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. 2016.
- [29] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *Proceedings of the British Machine Vision Conference 2016, BMVC 2016, York, UK, September 19-22, 2016*, 2016.
- [30] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *ICLR*, 2017.