

MaxDropout: Deep Neural Network Regularization Based on Maximum Output Values

Claudio Filipi Goncalves do Santos
Federal University of Sao Carlos -
UFSCar, Brazil
cfsantos@ufscar.br

Danilo Colombo
Petróleo Brasileiro - Petrobras
colombo.danilo@petrobras.com.br

Mateus Roder
and João Paulo Papa
São Paulo State University -
UNESP, Brazil
{mateus.roder,joao.papa}@unesp.br

Abstract—Different techniques have emerged in the deep learning scenario, such as Convolutional Neural Networks, Deep Belief Networks, and Long Short-Term Memory Networks, to cite a few. In lockstep, regularization methods, which aim to prevent overfitting by penalizing the weight connections, or turning off some units, have been widely studied either. In this paper, we present a novel approach called MaxDropout, a regularizer for deep neural network models that works in a supervised fashion by removing (shutting off) the prominent neurons (i.e., most active) in each hidden layer. The model forces fewer activated units to learn more representative information, thus providing sparsity. Regarding the experiments, we show that it is possible to improve existing neural networks and provide better results in neural networks when Dropout is replaced by MaxDropout. The proposed method was evaluated in image classification, achieving comparable results to existing regularizers, such as Cutout and RandomErasing, also improving the accuracy of neural networks that uses Dropout by replacing the existing layer by MaxDropout.

I. INTRODUCTION

Following the advent of deeply connected systems and the new era of information, tons of data are generated every moment by different devices, such as smartphones or notebooks. A significant portion of the data can be collected from images or videos, which are usually encoded in a high-dimensional domain. Deep Learning (DL) techniques have been broadly employed in different knowledge fields, mainly due to their ability to create authentic representations of the real world, even for multimodal information. Recently, DL has emerged as a prominent area in Machine Learning, since its techniques have achieved outstanding results and established several hallmarks in a wide range of applications, such as motion tracking [1], action recognition [2], and human pose estimation [3], [4], to cite a few.

Deep learning architectures such as Convolutional Neural Networks (CNNs), Deep Autoencoders, and Long Short-Term Memory Networks are powerful tools that deal with different image variations such as rotation or noise. However, their performance is highly data-dependent, which can cause some problems during training and further generalization for unseen examples. One common problem is overfitting, where the technique memorizes the data either due to the lack of information or because of too complex neural network architectures.

Such a problem is commonly handled with regularization methods, which represent a wide area of study in the scientific

community. The employment of one or more of such techniques provides useful improvements in different applications. Among them, two well-known methods can be referred: (i) so-called “Batch Normalization” and (ii) “Dropout”. The former was introduced by Ioffe et al. [5] and performs data normalization in the output of each layer. The latter was introduced by Srivastava et al. [6], and randomly deactivates some neurons present in each layer, thus forcing the model to be sparse.

However, dropping neurons out at random may slow down convergence during learning. To cope with this issue, we introduced an improved approach for regularizing deeper neural networks, hereinafter called “MaxDropout”¹, which shuts off neurons based on their maximum activation values, i.e., the method drops the most active neurons to encourage the network to learn better and more informative features. Such an approach achieved remarkable results for the image classification task, concerning two important well-established datasets.

The remainder of this paper is presented as follows: Section II introduces the correlated works, while Section III presents the proposed approach. Further, Section IV describes the methodology and datasets employed in this work. Finally, Sections V and VII provide the experimental results and conclusions, respectively.

II. RELATED WORKS

Regularization methods are widely used by several deep neural networks (DNNs) and with different architectures. The main idea is to help the system to prevent the overfitting problem, which causes the data memorization instead of generalization, also allowing DNNs to achieve better results. A well-known regularization method is Batch Normalization (BN), which works by normalizing the output of a given layer at each iteration. The original work [5] showed that such a process speeds up convergence for image classification tasks. Since then, several other works [7], [8], [9], including the current state-of-the-art on image classification [10], also highlighted its importance.

As previously mentioned, Dropout is one of the most employed regularization methods for DNNs. Such an approach was developed between 2012 and 2014 [6], showing significant

¹<https://github.com/cfsantos/MaxDropout-torch>

improvements in neural network’s performance for various tasks, ranging from image classification, speech recognition, and sentimental analysis. The standard Dropout works by creating, during training time, a mask that direct multiples all values of a given tensor. The values of such a mask follow the Bernoulli distribution, being 0 with a probability p and 1 with a probability $1 - p$ (according to the original work [6], best value for p in hidden layers is 0.5). During training, some values will be kept while others will be changed to 0. Visually, it means that some neurons will be deactivated while others will work normally.

After the initial development of the standard Dropout, Wang and Manning [11] explored different strategies for sampling since at each mini-batch a subset of input features is turned off. Such a fact highlights an interesting Dropout feature since it represents an approximation by a Markov chain executed several times during training. Since the Bernoulli distribution tends to a Normal distribution when the dimensional space is high enough, such an approximation allows Dropout to its best without sampling.

In 2015, Kingma et al. [12] proposed the Variational Dropout, a generalization of Gaussian Dropout in which the dropout rates are learned instead of randomly attributed. They investigated a local reparameterization approach to reduce the variance of stochastic gradients in variational Bayesian inference of a posterior over the model parameters, thus retaining parallelizability. On the other hand, in 2017, Gal et al. [13] proposed a new Dropout variant to reinforcement learning models. Such a method aims to improve the performance and better calibration of uncertainties once it is an intrinsic property of the Dropout. In such a field, the proposed approach allows the agent to adapt its uncertainty dynamically as more data is seen.

Later on, Molchanov et al. [14] explored the Variational Dropout proposed by Kingma et al. [12]. The authors extended it to situations when dropout rates are unbounded, leading to very sparse solutions in fully-connected and convolutional layers. Moreover, they achieved a reduction in the number of parameters up to 280 times on LeNet architectures, and up to 68 times on VGG-like networks with a small decrease in accuracy. Such a fact points out the importance of sparsity for parameter reduction and performance overall improvement.

Paralleling, other regularization methods have been emerged, like the ones that change the input of the neural network. For instance, Cutout [15] works by literally cutting off a region of the image (by setting the values of a random region to 0). This simple approach shows relevant results on several datasets. Another similar regularizer is the RandomErasing [16], that works in the same manner, but instead of setting the values of the region to 0, it changes these pixels to random values.

By bringing the concepts mentioned above and works close to the proposed approach, one can point out that the MaxDropout is similar to the standard Dropout, however, instead of randomly dropping out neurons, our approach follows a policy for shutting off the most active cells, representing a

selection of neurons that may overfit the data, or discourage the fewer actives from extracting useful information.

III. PROPOSED APPROACH

The proposed approach aims at shutting out the most activated neurons, which is responsible for inducing sparsity in the model, at the step that encourage the hidden neurons to learn more informative features and extract useful information that positively impacts the network’s generalization ability.

For the sake of visualization, Figures 1a-c show the differences between the proposed approach and the standard Dropout, in which Figure 1a stands for the original grayscale image and Figures 1b and 1c denote their corresponding outcomes after Dropout and MaxDropout. It is important to highlight that Dropout removes any pixel of the image randomly, while MaxDropout tends to inactivate the lighter pixels.

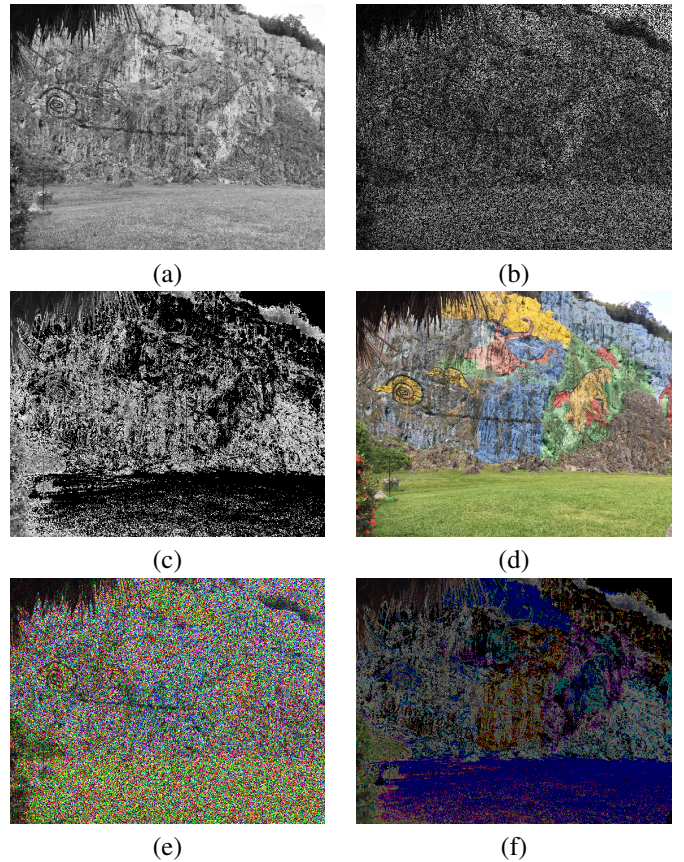


Fig. 1: Simulation using grayscale (a)-(c) and colored images (d)-(f): (a) original grayscale image and its outcomes after (b) Dropout and (c) MaxDropout transformations, respectively, and (d) original colored image and its outcomes after (e) Dropout and (f) MaxDropout transformations, respectively. In all cases, the drop rout ate is 50%.

The rationale behind the proposed approach can be better visualized in a tensor-like data. Considering the colored image showed in Figure 1d, one can observe its outcome after Dropout and MaxDropout transformations in Figures 1e

and 1f, respectively. Regarding standard Dropout, the image looks like a colored uniform noise, while MaxDropout could remove entire regions composed of bright pixels (i.e., pixels with high activation values, as expected).

For the sake of clarification purposes, Algorithm 1 implements the proposed MaxDropout²: the main loop in Lines 1–9 is in charge of the training procedure, and the inner loop in Lines 2–8 is executed for each hidden layer. Line 3 computes a random value uniformly distributed that is going to work as the dropout rate r . The output of each layer produces an $x \times y \times z$ tensor, where x and y stand for the image’s size, and z denotes the number of feature maps produced for each convolutional kernel. Line 4 creates a copy of the original tensor and uses an L_2 normalization to produce an output between 0 and 1.

Algorithm 1 Pseudocode for MaxDropout training algorithm.

```

1: while training do
2:   for each layer do
3:      $rate \leftarrow U(0, r)$ 
4:      $normTensor \leftarrow L2Normalize(Tensor)$ 
5:      $max \leftarrow Max(normTensor)$ 
6:      $keptIdx \leftarrow IdxOf(normTensor, (1 - rate) * max)$ 
7:      $returnTensor \leftarrow Tensor * KeptIdx$ 
8:   end for
9: end while

```

Later, Line 5 finds the biggest value in the normalized tensor, once it may not be equal to one³. Line 6 creates another tensor with the same shape as the input one and assigns 1 where $(1 - rate) \times max$ at a certain tensor position is greater than a given threshold; otherwise it sets such a position to 0. Finally, Line 7 creates the tensor to be used in the training phase, where each position of the original tensor is multiplied by the value in the respective position of the tensor created in the line before. Therefore, such a procedure guarantees that only values smaller than the threshold employed in Line 3 go further on.

IV. EXPERIMENTS

In this section, we describe the methodology employed to validate the robustness of the proposed approach. The hardware used for the paper is an Intel Xeon Bronze[®] 3104 CPU with 6 cores (12 threads), 1.70GHz, 96GB RAM with 2666Mhz, and a GPU Nvidia Tesla P4 with 8GB. Since most of the regularization methods aim to improve image classification tasks, we decided to follow the same protocol and approaches for a fair comparison.

A. Neural Network Structure

Regarding the neural network structure, we evaluated the proposed approach in two different practices. For the former experiments, regularization layers were added to a neural network that does not drop any transformation between layers.

²The pseudocode uses Keras syntax.

³Depending on the floating-point precision, the maximum value can be extremely close but not equal to one.

Concerning the latter experiments, the standard Dropout [6] layers were changed by the MaxDropout one to compare results.

For the first experiment, ResNet18 [17] was chosen because such an architecture has been used in several works for comparison purposes when coming to new regularizer techniques. ResNet18 is compounded by a sequence of convolutional residual blocks, followed by the well-known BatchNormalization [5]. As such, a MaxDropout layer was added between these blocks, changing the basic structure during training but keeping it to inference purposes.

In the second experiment, a slightly different approach has been performed. Here, a neural network that already has the Dropout regularization in its composition was considered for direct comparison among methods. The WideResNet [18] uses Dropout layers in its blocks with outstanding results on image classification tasks, thus becoming a good choice.

B. Training Protocol

In this work, we considered a direct comparison with other regularization algorithms. To be consistent with the literature, we provided the error rate instead of the accuracy itself [15], [18], [16]. Nonetheless, to ensure that the only difference between the proposed approach and the baselines used for comparison purposes concerns the MaxDropout layer, we strictly followed the protocols according to the original works.

To compare MaxDropout with other regularizers, we followed the protocol proposed by DeVries and Taylor [15], in which five runs were repeated, and the mean and the standard deviation are used for comparison purposes. For the experiment, the images from the datasets were normalized per-channel using mean and standard deviation.

During the training procedure, the images were shifted four pixels in every direction and then cropped into 32x32 pixels. Besides, the images were horizontally mirrored with a 50% probability. In such a case, two comparisons were provided. In the first case, besides the data augmentation already described, only the MaxDropout was included in the ResNet18 structure, directly comparing to the other methods. Regarding the second case, the Cutout data augmentation was included, providing a direct comparison of the results, showing that the proposed approach can work nicely.

As previously mentioned, to evaluate the MaxDropout against the standard Dropout, we choose the Wide Residual Network [18], and the same training protocol and parameters were employed to make sure the only difference concerns the type of neuron dropping.

C. Datasets

In this work, two well-established datasets in the literature were employed, i.e., CIFAR-10 [19] and its enhanced version CIFAR-100 [19]. Using such datasets allows us to compare the proposed approach toward important baseline methods, such as the standard Dropout [6] and CutOut [15]. Figure 2 portrays random samples extracted from the datasets mentioned above.

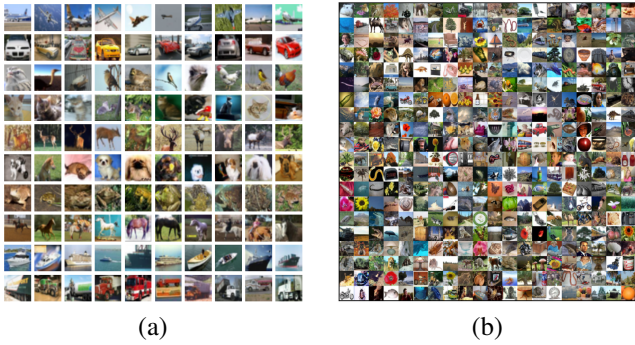


Fig. 2: Random training samples from: (a) CIFAR-10 and (b) CIFAR-100 datasets.

CIFAR-10 dataset comprises 10 classes equally distributed in 60,000 colored image samples, with a dimension of 32x32 pixels. The entire dataset is partitioned into 50,000 training images and 10,000 test images. On the other hand, CIFAR-100 dataset holds similar aspects of its smaller version, but now with 100 classes equally distributed in 60,000 colored image samples, with 600 images samples per class. Nonetheless, the higher number of classes and the low number of samples per class make image classification significantly hard in this case.

V. EXPERIMENTAL RESULTS

This section is divided into four main parts. First, we provided a convergence study during training for all experiments. Later, we compared the results of MaxDropout with other methods showing that, when combined with other regularizers, MaxDropout can lead to even better performance than their original versions. Finally, in the last part, we make a direct comparison between the proposed approach and standard Dropout by replacing the equivalent layer with the MaxDropout in the Wide-ResNet.

A. Training Evolution

Figures 3 and 4 depict the mean accuracies concerning the test set considering the 5 runs during training phase. Since we are dealing with regularizers, it makes sense to analyze their behavior during training and, for each epoch, compute their accuracy over the test set. One can notice that the proposed approach can improve the results even when the model is near to overfit.

B. Comparison Against Other Regularizers

As aforementioned, we considered a comparison against some baselines over five runs and exposed their mean accuracies and standard deviation in Table I. Such results evidence the robustness of the proposed approach against two other well-known regularizers, i.e., Cutout, and the RandomErasing.

From Table I, one can notice that when MaxDropout is incorporated within ResNet18 blocks, it allows the model to accomplish relevant and better results. Regarding the CIFAR-10 dataset, the model that uses MaxDropout achieved a reduction of around 0.5% in the error rate when compared to

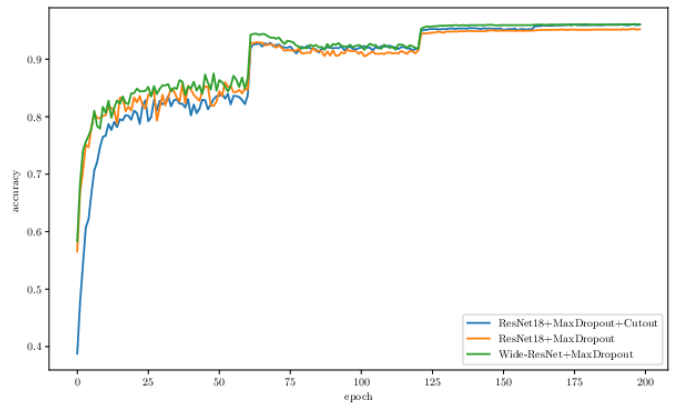


Fig. 3: Convergence over CIFAR-10 test set.

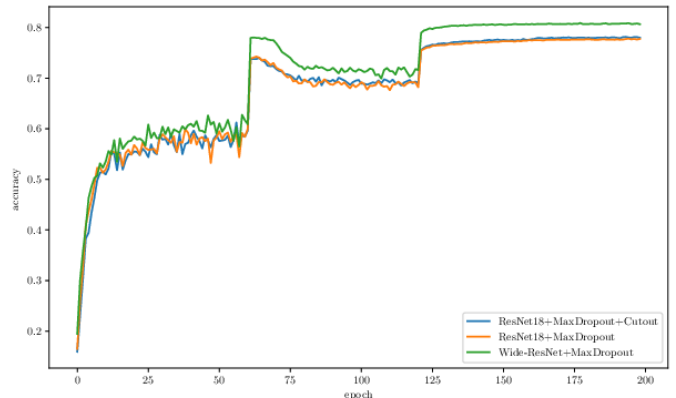


Fig. 4: Convergence over CIFAR-100 test set.

Approach	CIFAR-100	CIFAR-10
ResNet18 [17], [16]	24.50 ± 0.19	5.17 ± 0.18
ResNet18+RandomErasing [16]	24.03 ± 0.19	4.31 ± 0.07
ResNet18+Cutout [15]	21.96 ± 0.24	3.99 ± 0.13
ResNet18+MaxDropout	21.93 ± 0.07	4.66 ± 0.14

TABLE I: Results of MaxDropout and other regularizers

ResNet18. However, concerning the CIFAR-100 dataset, the model achieved over 2% less error than the same baseline, besides being statistically similar to Cutout.

C. Working Along with Other Regularizers

Since MaxDropout works inside the neural network by changing the hidden layers' values, it permits the concomitant functionality with other methods that change information from the input, such as Cutout. Table II portrays the results of each stand-alone approach and their combination. From these results, one can notice a slight improvement in performance considering the CIFAR-100 dataset, but it ends up as a relevant gain on CIFAR-10 dataset, reaching the best results so far.

D. MaxDropout x Dropout

One interesting point such a work stands for concerns the following question: Is the MaxDropout comparable to the standard Dropout [6]? To answer this question, we compared

Regularizer	CIFAR-100	CIFAR-10
Cutout [15]	21.96 \pm 0.24	3.99 \pm 0.13
MaxDropout	21.93 \pm 0.07	4.66 \pm 0.14
MaxDropout + Cutout	21.82 \pm 0.13	3.76 \pm 0.08

TABLE II: Results of the MaxDropout combined with Cutout.

the proposed approach against standard Dropout by replacing it with MaxDropout on the Wide Residual Network (WRN).

From Table III, one can observe the model using MaxDropout works slightly better than standard Dropout, leading to dropping in the error rate regarding CIFAR-100 and CIFAR-10 datasets by 0.04 and 0.05%, respectively. Although it may not look an impressive improvement, we showed that the proposed approach has a margin to improve the overall results, mainly when the threshold of the MaDropout is taken into account (i.e., ablation studies)⁴.

Model	CIFAR-100	CIFAR-10
WRN [18]	19.25	4.00
WRN + Dropout [18]	18.85	3.89
WRN + MaxDropout	18.81	3.84

TABLE III: Results of Dropout and MaxDropout over the WRN.

VI. DISCUSSION

Unfortunately, the approaches employed for comparison purposes did not release their training evolution for a direct comparison in Section V-A. Nevertheless, it is possible to observe that all models performed very well for the image classification task. In Table I, MaxDropout shows a result as good as Cutout for CIFAR-100 dataset, demonstrating it performs as expected when improving baseline models' results. However, it did not perform as well for CIFAR-10 dataset, but it still improves the baseline model results by almost 0.5%.

Results from Table II show that the MaxDropout supports the improvement when another regularizer is used along with. Although Cutout has been used to demonstrate the proposed approach's effectiveness, one can consider other similar regularizers. The most interesting results can be found in Table III, where MaxDropout is directly compared to the standard Dropout. It shows relevant gains over the baseline model, and it performs a little better than Dropout using the same drop rate, indicating that it may be the case to find out the best drop rates for MaxDropout, which can be data or model-dependent.

VII. CONCLUSIONS AND FUTURE WORKS

In this paper, we introduced MaxDropout, an improved version of the original Dropout method. Experiments show that it can be incorporated into existing models, working along

⁴We did not show the standard deviation since the original original study did not present such an information as well.

with other regularizers, such as Cutout, and can replace the standard Dropout with some accuracy improvement.

With relevant results, we intend to conduct a more in-depth investigation to figure out the best drop rates depending on the model and the training data. Moreover, the next step is to re-implement MaxDropout and make it available in other frameworks, like TensorFlow and MXNet, and test in other tasks, such as object detection and image segmentation.

Nonetheless, we showed that MaxDropout works very well for image classification tasks. For future works, we intended to perform evaluations in other different tasks such as natural language processing and automatic speech recognition.

ACKNOWLEDGMENTS

The authors are grateful to São Paulo Research Foundation - FAPESP (#2013/07375-0, #2014/12236-1, #2017/25908-6, and #2019/07825-1, Brazilian National Council for Scientific and Technological Development - CNPq (#307066/2017-7, and #427968/2018-6), and Petrobras (#2017/00285-6).

REFERENCES

- [1] N. Doulamis, "Adaptable deep learning structures for object labeling/tracking under dynamic visual environments," *Multimedia Tools and Applications*, vol. 77, no. 8, pp. 9651–9689, 2018.
- [2] S. Cao and R. Nevatia, "Exploring deep learning based solutions in fine grained activity recognition in the wild," in *2016 23rd International Conference on Pattern Recognition (ICPR)*. IEEE, 2016, pp. 384–389.
- [3] A. Toshev and C. Szegedy, "Deeppose: Human pose estimation via deep neural networks," *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1653–1660, 2014.
- [4] X. Chen and A. L. Yuille, "Articulated pose estimation by a graphical model with image dependent pairwise relations," in *Advances in neural information processing systems*, 2014, pp. 1736–1744.
- [5] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.
- [6] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *The journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [7] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising," *IEEE Transactions on Image Processing*, vol. 26, no. 7, pp. 3142–3155, 2017.
- [8] M. Simon, E. Rodner, and J. Denzler, "Imagenet pre-trained models with batch normalization," *arXiv preprint arXiv:1612.01452*, 2016.
- [9] J. Wang and X. Hu, "Gated recurrent convolution neural network for OCR," in *Advances in Neural Information Processing Systems*, 2017, pp. 335–344.
- [10] M. Tan and Q. V. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," *arXiv preprint arXiv:1905.11946*, 2019.
- [11] S. Wang and C. Manning, "Fast dropout training," in *international conference on machine learning*, 2013, pp. 118–126.
- [12] D. P. Kingma, T. Salimans, and M. Welling, "Variational dropout and the local reparameterization trick," in *Advances in Neural Information Processing Systems*, 2015, pp. 2575–2583.
- [13] Y. Gal, J. Hron, and A. Kendall, "Concrete dropout," in *Advances in Neural Information Processing Systems*, 2017, pp. 3581–3590.
- [14] D. Molchanov, A. Ashukha, and D. Vetrov, "Variational dropout sparsifies deep neural networks," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 2498–2507.
- [15] T. DeVries and G. W. Taylor, "Improved regularization of convolutional neural networks with cutout," *arXiv preprint arXiv:1708.04552*, 2017.
- [16] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang, "Random erasing data augmentation," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2020.

- [17] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," in *European conference on computer vision*. Springer, 2016, pp. 630–645.
- [18] S. Zagoruyko and N. Komodakis, "Wide residual networks," *arXiv preprint arXiv:1605.07146*, 2016.
- [19] A. Krizhevsky, "Learning multiple layers of features from tiny images," University of Toronto, Tech. Rep., 2009.