# TARGETDROP: A TARGETED REGULARIZATION METHOD FOR CONVOLUTIONAL NEURAL NETWORKS

*Hui Zhu[1,2], Xiaofang Zhao[1,*]*

[1]Institute of Computing Technology, Chinese Academy of Sciences, China
[2]University of Chinese Academy of Sciences, China

## ABSTRACT

Dropout regularization has been widely used in deep learning but performs less effective for convolutional neural networks since the spatially correlated features allow dropped information to still flow through the networks. Some structured forms of dropout have been proposed to address this but prone to result in over or under regularization as features are dropped randomly. In this paper, we propose a targeted regularization method named TargetDrop which incorporates the attention mechanism to drop the discriminative feature units. Specifically, it masks out the target regions of the feature maps corresponding to the target channels. Experimental results compared with the other methods or applied for different networks demonstrate the regularization effect of our method.

***Index Terms***— Dropout, Attention, Targeted Regularization, Convolutional Neural Networks

## 1. INTRODUCTION

Convolutional neural networks are widely used in the field of computer vision and have achieved great success. Many excellent neural architectures have been designed successively such as ResNet [1], DenseNet [2] and SENet [3]. In order to solve the over-fitting problem caused by the increase in the number of parameters for convolutional neural networks, many regularization methods have been proposed, such as weight decay, data augmentation and dropout [4].

However, The effect of dropout for convolutional neural networks is not as significant as that for fully connected networks because the spatially correlated features allow dropped information to still flow through convolutional networks [5]. To address this problem, some structured forms of dropout have been proposed such as SpatialDropout [6], Cutout [7] and DropBlock [5]. But these methods prone to result in over or under regularization as features are dropped randomly.

Some methods attempt to combine structured dropout with attention mechanism such as AttentionDrop [8] and CorrDrop [9]. However, these methods only mask out the units with higher activation values or the regions with less



**Fig. 1**. Masks of naive Dropout [4], Dropblock [5] and our TargetDrop. The red regions denote the regions to be masked.

discriminative information in the spatial dimension. They ignore the instructive information in the channel dimension which is proven to be meaningful in convolutional neural networks [3], even in dropout-based regularization methods[5].

In this paper, we propose a novel regularization method named TargetDrop, which drops the feature units with a clear target. Specifically, we choose the target channels and then drop the target regions in the corresponding feature maps. As is shown in Fig.1, compared with naive Dropout and DropBlock which may lead to unexpected results by dropping randomly, our TargetDrop prone to precisely mask out several effective features of the main object, thus forcing the network to learn more crucial information. Our experimental results demonstrate that TargetDrop can greatly improve the performance of convolutional neural networks and outperforms many other state-of-the-art methods on public datasets CIFAR-10 and CIFAR-100 which we attribute to our method.

Our contributions are summarized as follows:

• We propose a targeted regularization method which incorporates attention mechanism to address the problem for unexpected results caused by dropping randomly.
• We propose the rule of choosing target channels and target regions, and further analyse the regularization effect.
• Our method achieves better regularization effect compared with the other state-of-the-art methods and is applicable to different architectures for image classification tasks.
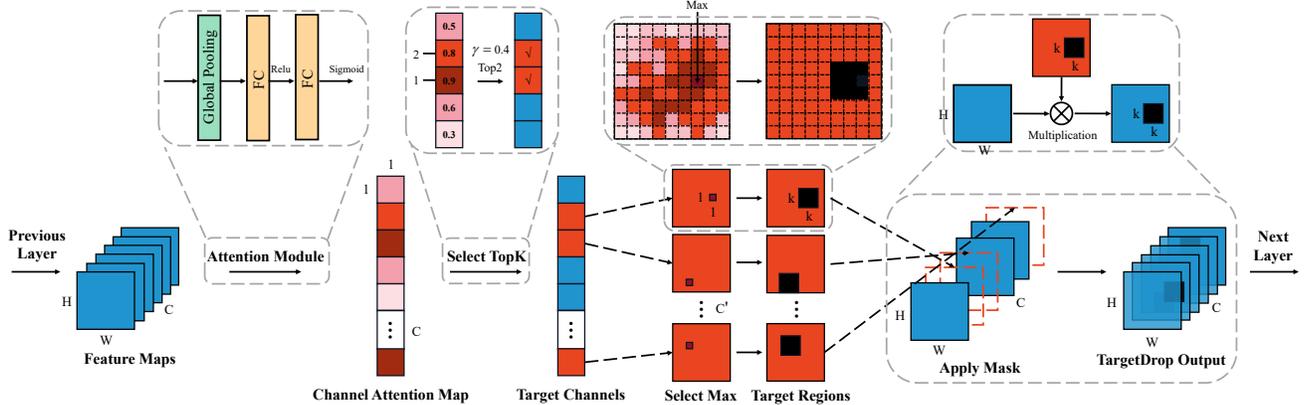
---

**Fig. 2**. The pipeline of TargetDrop. First, channel attention map is generated by processing the output of the previous layer through attention mechanism. Next, $topK$ elements are selected as the target channels according to the drop probability $\gamma$. Then, locate to a pixel with the maximum value in the feature map corresponding to each target channel and generate a mask by dropping the $k \times k$ target region. Finally, the mask is applied to the original feature maps by the multiplication operation.

## 2. RELATED WORK

Since Dropout [4] was proposed to improve the performace of networks by avoiding overfitting the training data, a series of regularization variations have been proposed such as Drop-Connect [10], SpatialDropout [6], DropPath [11], DropBlock [5], AttentionDrop [8], CorrDrop[9] and DropFilterR [12]. In addition, several methods about attention processing are also related as our method incorporates the attention mechanism into the dropout-based regularization. Methods for computing the spatial or channel-wise attention have achieved a certain effect such as Squeeze-and-Excitation (SE) module [3], Convolutional Block Attention Module (CBAM) [13] and Selective Kernel (SK) unit [14]. Our method outperforms the dropout-based regularization counterparts by utilizing the attention mechanism to achieve the targeted dropout.

## 3. METHODS

In this section, we propose our method TargetDrop which mainly contains seeking out the target channels and target regions. The pipeline of TargetDrop is shown in Fig. 2.

### 3.1. Target Channels

Given the output of the previous convolutional layer as $U = [u_1, u_2, \cdots, u_C] \in \mathbb{R}^{H \times W \times C}$, where $H$ and $W$ are the height and width of the feature map respectively, $C$ is the number of channels. As a first step, we are eager to figure out the importance of each channel. We aggregate the spatial information of each feature map into channel-wise vector by using global average pooling which has been proven to be effective [3, 13]. This vector $v \in \mathbb{R}^{1 \times 1 \times C}$ can be regarded as the statistic generated by shrinking through spatial dimensions $H \times W$ and this operation $F_{U \to v}$ can be defined as:

$$v_c = F_{U \to v}(u_c) = \frac{1}{HW} \sum_{i=1}^{H} \sum_{j=1}^{W} u_c(i, j) \qquad (1)$$

where $v_c$ denotes the $c$-th element of $v$. To further capture channel-wise dependencies, the vector is then forwarded to a shared network to produce the channel attention map $M \in \mathbb{R}^{1 \times 1 \times C}$. The shared network is composed of two fully connected (FC) layers and two activation functions. Specifically, a dimensionality-reduction layer with parameters $W_1 \in \mathbb{R}^{\frac{C}{r} \times C}$, a ReLU, a dimensionality-increasing layer with parameters $W_2 \in \mathbb{R}^{C \times \frac{C}{r}}$ and then a Sigmoid function are connected alternately. Here, $r$ is the reduction ratio to adjust the bottleneck. The map indicates the inter-channel relationships, and this operation $F_{v \to M}$ can be defined as:

$$M = F_{v \to M}(v, W) = \sigma(W_2 \delta(W_1 v)) \qquad (2)$$

where $\delta$ and $\sigma$ refer to the ReLU and Sigmoid, respectively.

Then, we sort all the values in $M$ and select the elements (tag "1" means to be selected and "0" if not) with top $K$ values as the target according to the drop probability $\gamma$. Specifically, the channels corresponding to those elements marked as tag "1" in the vector $T \in \mathbb{R}^{1 \times 1 \times C}$ are the target channels. Given the top $K$-th value in $M$ as $M_{topK}$ and this process can be described as:

$$K = \lfloor \gamma C \rfloor, \qquad T_p = \begin{cases} 1 & M_p \geq M_{topK} \\ 0 & \text{otherwise} \end{cases} \qquad (3)$$

where $M_p$ and $T_p$ denote the $p$-th elements of $M$ and $T$. Based on this, we further select the target regions of the original $H \times W$ feature maps corresponding to the target channels which we will elaborate in the following subsection.

## 3.2. Target Regions

For each feature map corresponding to a target channel, we hope to further seek out a region with much discriminative information in convolution operation. Utilizing spatial attention mechanism like a convolution operation with the kernel size of $7 \times 7$ is not necessary and may lead to considerable additional computation. Considering the continuity of image pixel values [8], we can simply locate to a pixel with maximum value and the other top values distributed in the surrounding continuous regions are most likely to be certain crucial features of the main object. Hence the location $(a, b)$ with the maximum value will be selected and the $k \times k$ region centered around it will be dropped. $h_1$, $h_2$, $w_1$ and $w_2$ represent the boundaries of the target region and the TargetDrop mask $\boldsymbol{S} = [\boldsymbol{s_1}, \boldsymbol{s_2}, \cdots, \boldsymbol{s_C}] \in \mathbb{R}^{H \times W \times C}$ can be described as:

$$h_1 = a - \lfloor \tfrac{k}{2} \rfloor, h_2 = a + \lfloor \tfrac{k}{2} \rfloor \qquad (4)$$
$$w_1 = b - \lfloor \tfrac{k}{2} \rfloor, w_2 = b + \lfloor \tfrac{k}{2} \rfloor$$

$$s_q(m, n) = \begin{cases} 0 & T_q = 1 \wedge h_1 \le m \le h_2 \wedge w_1 \le n \le w_2 \\ 1 & \text{otherwise} \end{cases} \qquad (5)$$

where $s_q$ and $T_q$ denote the $q$-th elements of $\boldsymbol{s}$ and $\boldsymbol{T}$. Given the final output as $\widetilde{\boldsymbol{U}} = [\widetilde{\boldsymbol{u}}_1, \widetilde{\boldsymbol{u}}_2, \cdots, \widetilde{\boldsymbol{u}}_C] \in \mathbb{R}^{H \times W \times C}$. Finally, we apply the mask and normalize the features:

$$\widetilde{\boldsymbol{u_z}} = \boldsymbol{u}_z \odot \boldsymbol{s}_z \times \frac{numel(\boldsymbol{s}_z)}{sum(\boldsymbol{s}_z)} \qquad (6)$$

where $\boldsymbol{u}_z$ and $\boldsymbol{s}_z$ denote the $z$-th elements of $\boldsymbol{u}$ and $\boldsymbol{s}$, $numel(\boldsymbol{s}_z)$ counts the number of units in $\boldsymbol{s}_z$, $sum(\boldsymbol{s}_z)$ counts the number of units where the value is "1" and $\odot$ represents the point-wise multiplication operation.

## 3.3. TargetDrop

---
**Algorithm 1** TargetDrop
___

**Input:** Output of the previous layer: $U$, Drop probability: $\gamma$,
    Size of the dropped block: $k \times k$, Phase of run: $phase$
**Output:** Final feature maps: $\widetilde{U}$.
  1: **if** $phase == interface$ **then**
  2:     **return** $U$
  3: **end if**
  4: Generate the channel attention map $\mathcal{M}$;
  5: Select the top $\mathcal{K}$-th value $\mathcal{M}_{top\mathcal{K}}$ in $\mathcal{M}$ according to $\gamma$ and generate the target channels: $\mathcal{T}$;
  6: Locate to the maximum value $u(a, b)$, produce a mask with the center being it and the width, height being $k$: $\mathcal{S}$;
  7: Apply the mask $\mathcal{S}$ by the multiplication operation and normalize the features to generate the output: $\widetilde{U}$;
  8: **return** $\widetilde{U}$

---

The pseudocode of our method is described in Algorithm 1.

**Table 1**. Comparison against the results of different state-of-the-art dropout-based regularization methods for classification accuracy of ResNet-18 on CIFAR-10 and CIFAR-100.

| Methods | C10 Error(%) | C100 Error(%) |
|---|---|---|
| No Regularization | 4.72 | 22.46 |
| Dropout [4] | 5.14 | 23.82 |
| DropBlock [5] | 4.59 | 21.95 |
| AttentionDrop [8] | 4.51 | 21.53 |
| TargetDrop (Ours) | **4.41** | **21.37** |
| Cutout [7] | 3.99 | 21.96 |
| Cutout + TargetDrop (Ours) | **3.67** | **21.25** |

**Table 2**. The regularization effect of TargetDrop on CIFAR-10 with different convolutional neural networks.

| Networks | Params | C10 Test Error(%) | | |
|---|---|---|---|---|
| | (Mil.) | Baseline | Dropout [4] | TargetDrop (Ours) |
| ResNet-20 [1] | 0.27 | 8.21 | 7.80 | **7.61** |
| VGG-16 [15] | 14.73 | 6.17 | 6.43 | **5.89** |
| WRN-28-10 [16] | 36.48 | 4.02 | 4.04 | **3.68** |

## 4. EXPERIMENTS

In this section, we introduce the implementation of experiments and report the performance of our method. We compare TargetDrop with the other state-of-the-art dropout-based methods on CIFAR-10 and CIFAR-100 [17] and apply it for different architectures. We further analyse the selection of hyper-parameters and visualize the class activation map.

## 4.1. Datasets

We use CIFAR-10 and CIFAR-100 [17] for image classification as basic datasets in our experiments. For preprocessing, we normalize the images using channel means and standard deviations and apply a standard data augmentation scheme: zero-padding the image with 4 pixels on each side, then cropping it to $32 \times 32$ and flipping it horizontally at random.

## 4.2. Training Method

Networks using the official PyTorch implementation are trained on the full training dataset until convergence and we report the highest test accuracy following common practice. The hyper-parameters in our experiments are as follows: the batch size is 128, the optimizer is SGD with Nesterov's momentum of 0.9, the initial learning rate is 0.1 and is decayed by the factor of 2e-1 at 0.4, 0.6, 0.8 ratio of total epochs. For Cutout [7], which is used in a few experiments, the cutout size is $16 \times 16$ for CIFAR-10 and $8 \times 8$ for CIFAR-100.
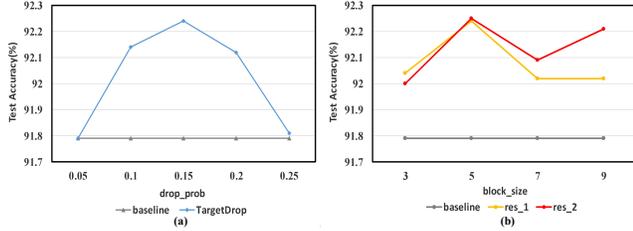
**Fig. 3**. (a) and (b) are the test accuracy on CIFAR-10 with different drop probabilities and block sizes, respectively.

### 4.3. Experiments for TargetDrop and Results

The experiments for TargetDrop we conducted mainly contain two parts: comparing the regularization effect with the other state-of-the-art dropout-based methods for ResNet-18 and show the performance for different architectures. We demonstrate the experiments and results in detail from these two aspects in the following two paragraphs. Specifically, in this part, the reduction ratio $r$ is 16, the drop probability and block size for TargetDrop are 0.15 and 5, respectively.

**Comparison against the results of other methods.** We compare the regularization effect with the other state-of-the-art dropout-based methods on ResNet-18. We apply our Target-Drop in the same way with the other methods that adding the regulation to the outputs of first two groups for a fair comparison. Specially, several methods we reproduced can not reach the reported results, so we refer to the data in the original paper [8] directly for these. As is shown in Table 1, the results of our method outperform the other methods on CIFAR-10 and CIFAR-100. Moreover, combined with Cutout [7], our method can achieve better regularization effect.

**Regularization on different architectures.** We further conduct experiments on CIFAR-10 with several classical convolutional neural networks to demonstrate that our method is applicable to different architectures. As is shown in Table 2, our method TargetDrop can improve the performances of different architectures. We can notice that Dropout [4] is not applicable for convolutional neural networks which is mentioned above and TargetDrop is an effective dropout-based regularization method for the networks on different scales.

The number of parameters added in our method is presented as follows. The additional parameters come from the channel attention mechanism, and the additional computation includes the simple selection of the maximum pixel besides this. The number of parameters in the training process only increase by about 0.02% and the amount of computation increase similarly. While in the test process, TargetDrop is closed like other methods, so the complexity will not change.

### 4.4. Analysis of the Hyper-parameters Selection

In this subsection, we further analyse the selection of hyper-parameters mentioned above: the drop probability $\gamma$ and the
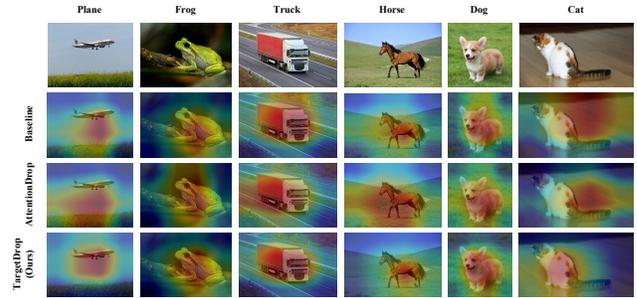


**Fig. 4**. Class activation mapping(CAM) [18] for ResNet-18 model trained with no regularization, AttentionDrop [8] and our method TargetDrop.

block size $k$. To analyse the former, we constrain the block size to 5 and TargetDrop is applied to the output of the first group (the size of the feature map is 32×32). To analyse the latter, we constrain the drop probability to 0.15 and Target-Drop is applied to the outputs of the first two groups (the size of the feature map is 32×32 for res_1 and 16×16 for res_2). As is shown in Fig. 3, our method is suitable for more channels and insensitive to different hyper-parameters to a certain extent which may be due to the targeted dropout. The drop probability of 0.15 and the block size of 5 are slightly better.

### 4.5. Activation Visualization

In this subsection, we utilize the class activation mapping (CAM) [18] to visualize the activation units of ResNet-18 [1] on several images as shown in Fig. 4. We can notice that the activation map generated by model regularized with our method TargetDrop demonstrates strong competence in capturing the extensive and relevant features towards the main object. Compared with the other methods, the model regularized with TargetDrop tends to precisely focus on those discriminative regions for image classification which we attribute to targeting and masking out certain effective features corresponding to the crucial channels.

## 5. CONCLUSION

In this paper, we propose the novel regularization method TargetDrop for convolutional neural networks, which addresses the problem for unexpected results caused by the untargeted methods to some extent by considering the importance of the channels and regions of feature maps. Extensive experiments demonstrate the outstanding performance of TargetDrop by comparing it with the other methods and applying it for different architectures. Furthermore, we analyse the selection of hyper-parameters and visualize the activation map to prove the rationality of our method. In addition to image classification tasks, we believe that TargetDrop is suitable for more datasets and tasks in the field of computer vision.

## 6. REFERENCES

[1] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, 2016, pp. 770–778.

[2] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger, "Densely connected convolutional networks," in *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, 2017, pp. 2261–2269.

[3] Jie Hu, Li Shen, and Gang Sun, "Squeeze-and-excitation networks," in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, 2018, pp. 7132–7141.

[4] Nitish Srivastava, Geoffrey E. Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.

[5] Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V. Le, "Dropblock: A regularization method for convolutional networks," in *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada.*, 2018, pp. 10750–10760.

[6] Jonathan Tompson, Ross Goroshin, Arjun Jain, Yann LeCun, and Christoph Bregler, "Efficient object localization using convolutional networks," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*. 2015, pp. 648–656, IEEE Computer Society.

[7] Terrance Devries and Graham W. Taylor, "Improved regularization of convolutional neural networks with cutout," *CoRR*, vol. abs/1708.04552, 2017.

[8] Zhihao Ouyang, Yan Feng, Zihao He, Tianbo Hao, Tao Dai, and Shu-Tao Xia, "Attentiondrop for convolutional neural networks," in *IEEE International Conference on Multimedia and Expo, ICME 2019, Shanghai, China, July 8-12, 2019*. 2019, pp. 1342–1347, IEEE.

[9] Yuyuan Zeng, Tao Dai, and Shu-Tao Xia, "Corrdrop: Correlation based dropout for convolutional neural networks," in *2020 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2020, Barcelona, Spain, May 4-8, 2020*. 2020, pp. 3742–3746, IEEE.

[10] Li Wan, Matthew D. Zeiler, Sixin Zhang, Yann Le-Cun, and Rob Fergus, "Regularization of neural networks using dropconnect," in *Proceedings of the 30th International Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16-21 June 2013*. 2013, vol. 28 of *JMLR Workshop and Conference Proceedings*, pp. 1058–1066, JMLR.org.

[11] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V. Le, "Learning transferable architectures for scalable image recognition," in *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*. 2018, pp. 8697–8710, IEEE Computer Society.

[12] Hengyue Pan, Xin Niu, Rongchun Li, Siqi Shen, and Yong Dou, "Dropfilterr: A novel regularization method for learning convolutional neural networks," *Neural Process. Lett.*, vol. 51, no. 2, pp. 1285–1298, 2020.

[13] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon, "CBAM: convolutional block attention module," in *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part VII*, Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss, Eds. 2018, vol. 11211 of *Lecture Notes in Computer Science*, pp. 3–19, Springer.

[14] Xiang Li, Wenhai Wang, Xiaolin Hu, and Jian Yang, "Selective kernel networks," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*. 2019, pp. 510–519, Computer Vision Foundation / IEEE.

[15] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.

[16] 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.

[17] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," *Computer Science Department, University of Toronto, Tech. Rep*, vol. 1, 01 2009.

[18] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 2921–2929.