

ROTATIONOUT AS A REGULARIZATION METHOD FOR NEURAL NETWORK

Anonymous authors

Paper under double-blind review

ABSTRACT

In this paper, we propose a novel regularization method, RotationOut, for neural networks. Different from Dropout that handles each neuron/channel independently, RotationOut regards its input layer as an entire vector and introduces regularization by randomly rotating the vector. RotationOut can also be used in convolutional layers and recurrent layers with small modifications. We further use a noise analysis method to interpret the difference between RotationOut and Dropout in co-adaptation reduction. Using this method, we also show how to use RotationOut/Dropout together with Batch Normalization. Extensive experiments in vision and language tasks are conducted to show the effectiveness of the proposed method. Codes will be available.

1 INTRODUCTION

Dropout (Srivastava et al., 2014) has proven to be effective for preventing overfitting over many deep learning areas, such as image classification (Shrivastava et al., 2017), natural language processing (Hu et al., 2016) and speech recognition (Amodei et al., 2016). In the years since, a wide range of variants have been proposed for wider scenarios, and most related work focus on the improvement of Dropout structures, i.e., how to drop. For example, drop connect (Wan et al., 2013) drops the weights instead of neurons, evolutionary dropout (Li et al., 2016) computes the adaptive dropping probabilities on-the-fly, max-pooling dropout (Wu & Gu, 2015) drops neurons in the max-pooling kernel so smaller feature values have some probabilities to affect the activations.

These Dropout-like methods process each neuron/channel in one layer independently and introduce randomness by dropping. These architectures are certainly simple and effective. However, randomly dropping independently is not the only method to introduce randomness. Hinton et al. (2012) argues that overfitting can be reduced by preventing co-adaptation between feature detectors. Thus it is helpful to consider other neurons' information when adding noise to one neuron. For example, lateral inhibition noise could be more effective than independent noise.

In this paper, we propose RotationOut as a regularization method for neural networks. RotationOut regards the neurons in one layer as a vector and introduces noise by randomly rotating the vector. Specifically, consider a fully-connected layer with n neurons: $\mathbf{x} \in \mathbb{R}^n$. If applying RotationOut to this layer, the output is $\mathcal{R}\mathbf{x}$ where $\mathcal{R} \in \mathbb{R}^{n \times n}$ is a random rotation matrix. It rotates the input with random angles and directions, bringing noise to the input. The noise added to a neuron comes not only from itself, but also from other neurons. It is the major difference between RotationOut and Dropout-like methods. We further show that RotationOut uses the activations of the other neurons as the noise to one neuron so that the co-adaptation between neurons can be reduced.

RotationOut uses random rotation matrices instead of unrestricted matrices because the directions of feature vectors are important. Random rotation provides noise to the directions directly. Most neural networks use dot product between the feature vector and weight vector as the output. The network actually learns the direction of the weights, especially when there is a normalization layer (e.g. Batch Normalization (Ioffe & Szegedy, 2015) or Weight Normalization (Salimans & Kingma, 2016)) after the weight layer. Random rotation of feature vectors introduces noise into the angle between the feature and the weight, making the learning of weights directions more stable. Sabour et al. (2017) also uses the orientation of feature vectors to represent the instantiation parameters in capsules. Another motivation for rotating feature vectors comes from network dissection. Bau et al. (2017) finds that random rotations of a learned representation can destroy the interpretability which

is axis-aligned. Thus random rotating the feature during training makes the network more robust. Even small rotations can be a strong regularization.

We study how RotationOut helps prevent neural networks from overfitting. Hinton et al. (2012) introduces *co-adaptation* to interpret Dropout but few literature give a clear concept of *co-adaptation*. In this paper, we provide a metric to approximate co-adaptations and derive a general formula for noise analysis. Using the formula, we prove that RotationOut can reduce co-adaptations more effectively than Dropout and show how to combine Dropout and Batch Normalization together.

In our experiments, RotationOut can achieve results on par with or better than Dropout and Dropout-like methods among several deep learning tasks. Applying RotationOut after convolutional layers and fully connected layers improves image classification accuracy of ConvNet on CIFAR100 and ImageNet datasets. On COCO datasets, RotationOut also improves the generalization of object detection models. For LSTM models, RotationOut can achieve competitive results with existing RNN dropout method for speech recognition task on Wall Street Journal (WSJ) corpus.

The main contributions of this paper are as follows: We propose RotationOut as a regularization method for neural networks which is different from existing Dropout-like methods that operate on each neuron independently. RotationOut randomly rotates the feature vector and introduces noise to one neuron with other neurons' information. We present a theoretical analysis method for general formula of noise. Using the method, we answer two questions: 1) how noise-based regularization methods reduce co-adaptations and 2) how to combine noise-based regularization methods with Batch Normalization. Experiments in vision and language tasks are conducted to show the effectiveness of the proposed RotationOut method.

Related Work Dropout is effective for fully connected layers. When applied to convolution layers, it is less effective. Ghiasi et al. (2018) argues that information about the input can still be sent to the next layer even with dropout, which causes the networks to overfit (Ghiasi et al., 2018). SpatialDropout (Tompson et al., 2015) drops the entire channel from the feature map. Shake-shake regularization (Gastaldi, 2017) drops the residual branches. Cutout (DeVries & Taylor, 2017) and Dropblock (Ghiasi et al., 2018) drop a continuous square region from the inputs/feature maps.

Applying standard dropout to recurrent layers also results in poor performance (Zaremba et al., 2014; Labach et al., 2019), since the noise caused by dropout at each time step prevents the network from retaining long-term memory. Gal & Ghahramani (2016); Moon et al. (2015); Merity et al. (2017) generate a dropout mask for each input sequence, and keep it the same at every time step so that memory can be retained.

Batch Normalization (BN) (Ioffe & Szegedy, 2015) accelerates deep network training. It is also a regularization to the network, and discourages the strength of dropout to prevent overfitting (Ioffe & Szegedy, 2015). Many modern ConvNet architectures such as ResNet (He et al., 2016) and DenseNet (Huang et al., 2017) do not apply dropout in convolutions. Li et al. (2019) is the first to argue that it is caused by the a variance shift. In this paper, we use the noise analysis method to further explore this problem.

There is a lot of work studying rotations in networks. Rotations on the images (Lenc & Vedaldi, 2015; Simard et al., 2003) are important data augmentation methods. There are also studies about rotation equivalence. Worrall et al. (2017) uses an enriched feature map explicitly capturing the underlying orientations. Marcos et al. (2017) applies multiple rotated versions of each filter to the input to solve problems requiring different responses with respect to the inputs' rotation. The motivations of these work are different from ours. The most related work is network dissection (Bau et al., 2017). They discuss the impact on the interpretability of random rotations of learned features, showing that rotation in training can be a strong regularization.

2 ROTATIONOUT

In this section, we first introduce the formulation of RotationOut. Next, we use linear models to demonstrate how RotationOut helps for regularization. In the last part, we discuss the implementation of RotationOut in neural networks.

2.1 RANDOM ROTATION MATRIX

A rotation in D dimension is represented by the product between a rotation matrix $\mathcal{R} \in \mathbb{R}^{D \times D}$ and the feature vector $\mathbf{x} \in \mathbb{R}^n$. The complexity for random rotation matrix generation and the matrix multiplication are both $O(D^2)$, which would be less efficient than Dropout with $O(D)$ complexity. We consider a special case that uses Givens rotations (Anderson, 2000) to construct random rotation matrices to reduce the complexity.

Let $D = 2d$ be an even number, and $P = [n_1, n_2, \dots, n_{2d}]$ be a permutation of $\{1, 2, \dots, D\}$. A rotation matrix can be generated by function $\mathbf{M}(\theta, P) = \{r_{ij}\} \in \mathbb{R}^{D \times D}$:

$$r_{ij} = \begin{cases} \cos \theta & \text{if } i = j \\ \sin \theta & \text{if } i = P_l, j = P_{l+d} \\ -\sin \theta & \text{if } i = P_{l+d}, j = P_l \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

Here P_l represents the l^{th} element of P where $1 \leq l \leq d$. See Appendix A.1 for some examples of such rotation matrices. Suppose we sample the angle θ from zero-centered distributions, e.g., truncated Gaussian distribution or uniform distribution and sample the permutation P from \mathcal{P} , the set of all permutations of $\{1, 2, \dots, D\}$, with equal probability. The RotationOut operator \mathcal{R} can be generated using the function $\mathbf{M}(P, \theta)$:

$$P \sim \mathcal{P}, \theta \sim \text{Unif}(-\Theta, \Theta), \mathcal{R} = \frac{1}{\cos \theta} \mathbf{M}(P, \theta). \quad (2)$$

Here $1/\cos \theta$ is a normalization term and \mathcal{R} is not a rotation matrix strictly speaking. The random operator generated from Equation 2 have some good properties. 1) The noise is zero centered: $\mathbb{E}_{\mathcal{R}}[\mathcal{R}\mathbf{x}] = \mathbf{x}$. 2) For any vector \mathbf{x} and any random permutation P , the angle between \mathbf{x} and $\mathcal{R}\mathbf{x}$ is determined by angle θ : $\langle \mathbf{x}, \mathcal{R}\mathbf{x} \rangle = \theta$. 3) For fixed angle θ , there exists $D!/d!$ different rotations. 4) The complexity for random rotation matrix generation and the matrix multiplication are both $O(D)$.

Permutation P draws the rotation direction and angle θ draws the rotation angle. As an analogy, permutation P is similar to the dropout mask widely used in RNN dropout. There exists 2^D different dropout mask ($2^D \ll D!/d!$ for $D > 8$), thus the diversity of random rotation in Equation 1 is sufficient for network training. Angle θ is similar to the percentage of dropped neurons in Dropout, and the distribution of θ controls the regularization strength. (Srivastava et al., 2014) used the multiplier's variance to compare Bernoulli dropout and Gaussian dropout. Following this setting, RotationOut is equivalent to Bernoulli Dropout with the keeping rate p and Gaussian dropout with variance σ^2 if $(1-p)/p = \sigma^2 = \mathbb{E}_{\theta} \tan^2 \theta$.

Reviewing the formulation of the random rotation matrix, it arranges all D dimensions of the input into d pairs randomly, and rotates the two dimension vectors with angle θ in each pair. Suppose u and v are two dimensions/neurons in one pair, the outputs of u and v after RotationOut are

$$\begin{bmatrix} u' \\ v' \end{bmatrix} = \begin{bmatrix} 1 & \tan \theta \\ -\tan \theta & 1 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} u + v \tan \theta \\ v - u \tan \theta \end{bmatrix} \quad (3)$$

The noise of u' comes from v and the noise of v' comes from u since θ is random. Note that the pairs are randomly arranged, thus RotationOut uses all other dimensions/neurons as the noise for one dimension/neuron of the feature vector. With RotationOut, the neurons are trained to work more independently since one neuron has to regard the activation of other neurons as noise. Thus the co-adaptations are reduced.

Consider Gaussian dropout, the outputs are $u' = u + u\epsilon$, $v' = v + v\epsilon$ where $\mathbb{E}\epsilon = 0$, $\mathbb{E}\epsilon^2 = \mathbb{E}_{\theta} \tan^2 \theta$. The difference between Gaussian dropout and RotationOut is the source of noise, i.e., the Gaussian dropout noise for one neuron comes from itself while the RotationOut noise comes from other neurons.

2.2 ROTATIONOUT IN LINEAR MODELS

First we consider a simple case of applying RotationOut to the classical problem of linear regression. Let $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$ be the dataset where $\mathbf{x}_i \in \mathbb{R}^D$, $y_i \in \mathbb{R}$. Linear regression tries to find the weight

$\mathbf{w} \in \mathbb{R}^D$ that minimizes $\sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}_i)^2$. When applied RotationOut to each \mathbf{x}_i , we generate \mathcal{R}_i from Equation 2 for each \mathbf{x}_i . The objective function becomes:

$$\min_{\mathbf{w}} \mathbb{E}_{\mathcal{R}} \left[\sum_{i=1}^N (y_i - \mathbf{w}^\top \mathcal{R}_i \mathbf{x}_i)^2 \right]. \quad (4)$$

Denote $\mathbf{y} = [y_1, y_2, \dots, y_n]^\top \in \mathbb{R}^N$, $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{N \times D}$. To compare RotationOut with Dropout with keep rate p , we suppose $\mathbb{E}_{\theta} \tan^2 \theta = (1-p)/p = \lambda$. Equation 4 reduces to:

$$\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \mathbf{w}^\top \frac{\text{trace}(\mathbf{X}^\top \mathbf{X}) \mathbf{I} - \mathbf{X}^\top \mathbf{X}}{D-1} \mathbf{w}. \quad (5)$$

Details see Appendix A.2. Solutions to Equation 5 (LR with Rotation) and the mirror problem with dropout (Srivastava et al., 2014) are :

$$\begin{aligned} \mathbf{w}_{\text{Rot}} &= \left[\mathbf{X}^\top \mathbf{X} + \lambda \frac{\text{trace}(\mathbf{X}^\top \mathbf{X}) \mathbf{I} - \mathbf{X}^\top \mathbf{X}}{D-1} \right]^{-1} \mathbf{X}^\top \mathbf{y} \\ \mathbf{w}_{\text{Drop}} &= [\mathbf{X}^\top \mathbf{X} + \lambda \text{diag}(\mathbf{X}^\top \mathbf{X})]^{-1} \mathbf{X}^\top \mathbf{y} \end{aligned} \quad (6)$$

Therefore, linear regression with RotationOut and Dropout are equivalent to ridge regression with different regularization terms. Set $\lambda = 1$ (Dropout rate $p = 0.5$) for simplicity. LR with Dropout doubles the diagonal elements of $\mathbf{X}^\top \mathbf{X}$ to make the problem numerical stable. LR with RotationOut is more close to ridge regression:

$$\mathbf{X}^\top \mathbf{X} + \frac{\text{trace}(\mathbf{X}^\top \mathbf{X}) \mathbf{I} - \mathbf{X}^\top \mathbf{X}}{D-1} = \frac{D-2}{D-1} \left[\mathbf{X}^\top \mathbf{X} + \frac{\text{trace}(\mathbf{X}^\top \mathbf{X})}{D-2} \mathbf{I} \right] \quad (7)$$

The condition number of Equation 7 and the LR with RotationOut problem is up bounded by $D-1$. For the Dropout case, if some data dimensions have extremely small variances, both $\mathbf{X}^\top \mathbf{X}$ and $\text{diag}(\mathbf{X}^\top \mathbf{X})$ are ill-conditioned. LR with Dropout problem has unbounded condition number.

Next we consider an m -way classification model of logistic regression. The input is $\mathbf{x} \in \mathbb{R}^D$ and the weights are $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_m] \in \mathbb{R}^{m \times D}$. The probability that the input belongs to the k category is:

$$p_k = \frac{\exp(\mathbf{w}_k \mathbf{x})}{\sum_i \exp(\mathbf{w}_i \mathbf{x})} = \frac{\exp(\|\mathbf{w}_k\| \|\mathbf{x}\| \cos \theta_k)}{\sum_i \exp(\|\mathbf{w}_i\| \|\mathbf{x}\| \cos \theta_i)}. \quad (8)$$

In Equation 8, θ_i denotes the angel between \mathbf{x} and \mathbf{w}_i . Assume that the length of each weights \mathbf{w}_i are very close, the input \mathbf{x} belongs to the k category if \mathbf{x} is most close to \mathbf{w}_k in angle.

Consider a hard sample case that $\theta_i < \theta_j$ are the two smallest weight-data angles. But θ_i and θ_j are very close: $\theta_i \approx \theta_j$, i.e., the data are close to the decision boundary. The model should classify the data correctly but could make mistakes if there is some noise. Applying RotationOut, the angle between the data and the weights can be changed, and the new angles can be $\hat{\theta}_i > \hat{\theta}_j$. To classify the data correctly, there should be a gap between θ_i and θ_j . In other words, the decision boundary changed from $\theta_i < \theta_j$ to $\theta_i < \theta_j - \Theta$ where Θ is a positive constant that depends on the regularization. Thus RotationOut can be regarded as a margin-based hard sample mining.

Here we provide an intuitive understanding of how Dropout with low keep rates leads to lower performance. Randomly zeroing units, Dropout method also rotates the feature vector. A lower keep rate results in a bigger rotation angle: $\cos^2 \theta = \frac{(\sum_i p_i x_i^2)^2}{\sum_i x_i^2 \sum_i p_i^2 x_i^2} \approx \frac{(\mathbb{E} p_i)^2}{\mathbb{E} p_i^2} = p$. Consider the last hidden layer in neural networks, it is similar to logistic regression on the features. If one feature \mathbf{x} is most close to \mathbf{w}_k , it belongs to the k^{th} . A lower keep rate Dropout would rotate the feature with a bigger angle, and the Dropout output can be most close to another weight with higher probability, which may hurts the training.

2.3 ROTATIONOUT IN NEURAL NETWORKS

Consider a neural network with L hidden layers. Let \mathbf{x}^l , \mathbf{y}^l , and \mathbf{W}^l denote the vector of inputs, the vector of output before activation, and the weights for the layer l . Let \mathcal{R} be generated from

Equation 2 and a be the activation function, for example Rectified Linear Unit (ReLU). The MLP feed-forward operation with RotationOut in training time can be:

$$\tilde{\mathbf{x}}^l = \mathcal{R}(\mathbf{x}^l - \mathbb{E}[\mathbf{x}^l]) + \mathbb{E}[\mathbf{x}^l], \quad \mathbf{y}^l = \mathbf{W}^l \tilde{\mathbf{x}}^l, \quad \mathbf{x}^{l+1} = a(\mathbf{y}^l). \quad (9)$$

We rotate the zero-centered features and then add the expectation back. The reasons will be explained later. Here we give an intuitive understanding. If features are not zero-centered, we do not know the exact regularization strength. Suppose all features elements are in one interval, say $1 < x < 2$. The angle between any two feature vectors is a sharp angle. In this case a rotation angle of $\pi/4$ would be too big. It is the same for Dropout. The regularization strength is influenced by the mean value of features which we may not know. At test time, the RotationOut operation is removed.

Consider 2D case for example, the input for 2D convolutional layers are three dimensional: number of channels C , width H and height W :

$$\mathbf{X} = \{\mathbf{x}_{hw}\}, \mathbf{x}_{hw} \in \mathbb{R}^C, 1 \leq h \leq H, 1 \leq w \leq W \quad (10)$$

We regard each \mathbf{x}_{hw} as a feature vector with semantic information for each position (h, w) , and apply rotation to each position. As Ghiasi et al. (2018) argued, the convolutional feature maps are spatially correlated, so information can still flow through convolutional layers if features are dropped out randomly. Similarly, if we rotate feature vectors in different positions with random directions, random directions offset each other and result in no rotation. So we rotate all feature vectors with the same directions but different angles. The operation on convolutional feature maps can be:

$$\begin{aligned} P &\sim \mathcal{P}, \quad \theta_{11}, \dots, \theta_{HW} \sim \text{Unif}(0, \Theta), \\ \forall h, w \quad \mathcal{R}_{hw} &= \mathbf{M}(P, \theta_{hw}), \\ \tilde{\mathbf{x}}_{hw} &= \mathcal{R}_{hw}(\mathbf{x}_{hw} - \mathbb{E}\mathbf{x}_{hw}) + \mathbb{E}\mathbf{x}_{hw}. \end{aligned} \quad (11)$$

The operation for general convolutional networks are very similar. Also note that RotationOut can combined with DropBlock (Ghiasi et al., 2018) easily: only rotating features in a continuous block. Experiments show that the combination can get extra performance gain. As mentioned in Section 3.1, the rotation directions defined by P is similar to the dropout mask in RNN drops. RotationOut can also be used in recurrent networks following Equation 11.

3 NOISE ANALYSIS

In this section, we first study the general formula of adding noise. Using the formula, we show how introducing randomness/noise helps reduce co-adaptations and why RotationOut is more efficient than the vanilla dropout. Next, we explain the variance between Dropout-and batch normalization (Li et al., 2019) using the formula and propose some solutions.

3.1 RANDOMNESS AND CO-ADAPTATIONS

Strictly speaking, the co-adaptations describe the dependence between neurons. The mutual information between two neurons may be the best metric to define co-adaptations. To compute mutual information, we need the exact distributions of neurons, which are generally unknown. So we consider the correlation coefficient to evaluate co-adaptations, which only need the first and second moment. Moreover, if we assume the distributions of neurons are Gaussian, correlation coefficient and mutual information are equivalent in co-adaptations evaluation.

Suppose $\mathbf{x} \in \mathbb{R}^D$ is the activations of one hidden layer. Let $\mathbb{E}[\mathbf{x}] = \mathbf{c} \in \mathbb{R}^D$, $\text{Var}[\mathbf{x}] = \Sigma \in \mathbb{R}^{D \times D}$. The ideal situation is that $\Sigma = \text{diag}\Sigma$, i.e., the neurons are mutually independent. We define the co-adaptations as the distance between Σ and $\Sigma = \text{diag}(\Sigma)$.

$$\text{co}(\mathbf{x}) = \frac{\|\Sigma - \text{diag}(\Sigma)\|_1}{\|\text{diag}(\Sigma)\|_1} = \frac{\|\Sigma - \text{diag}(\Sigma)\|_1}{\text{trace}(\Sigma)} \quad (12)$$

Here $\text{trace}(\Sigma)$ is a normalization term that defines the regularization strength. We use L_1 distance but not L_2 distance because the variance is already a second moment of \mathbf{x} . Let $\tilde{\mathbf{x}}$ be the out of \mathbf{x} with arbitrary noise (e.g. Dropout or RotationOut). We assume that the noise should follow two

assumptions: 1) zero-center: $\mathbb{E}[\tilde{\mathbf{x}}|\mathbf{x}] = \mathbf{x}$; 2) non-trivial: $\text{Var}[\tilde{\mathbf{x}}|\mathbf{x}] \neq \mathbf{O}$ (avoid that $\tilde{\mathbf{x}}$ always equals to \mathbf{x}). Consider the law of total variance, we have:

$$\text{Var}[\tilde{\mathbf{x}}] = \mathbb{E}[\text{Var}[\tilde{\mathbf{x}}|\mathbf{x}]] + \text{Var}[\mathbb{E}[\tilde{\mathbf{x}}|\mathbf{x}]] = \mathbb{E}[\text{Var}[\tilde{\mathbf{x}}|\mathbf{x}]] + \text{Var}[\mathbf{x}] \quad (13)$$

Let $\tilde{\mathbf{x}}_{\text{Drop}}$ be the out of \mathbf{x} after Dropout with drop rate p , and $\tilde{\mathbf{x}}_{\text{Rot}}$ be the out of \mathbf{x} after RotationOut with $\mathbb{E}_\theta \tan^2 \theta = (1-p)/p$, we have Lemma 1 (proof see Appendix A.3):

Lemma 1. $\text{Var}[\tilde{\mathbf{x}}_{\text{Drop}}|\mathbf{x}] = \frac{1-p}{p} \text{diag}(\mathbf{x}\mathbf{x}^T)$, $\text{Var}[\tilde{\mathbf{x}}_{\text{Rot}}|\mathbf{x}] = \frac{1-p}{p(D-1)}(\mathbf{x}^T\mathbf{x}\mathbf{I} - \mathbf{x}\mathbf{x}^T)$.

Note that $\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \Sigma + \mathbf{c}\mathbf{c}^T$, $\mathbb{E}[\mathbf{x}^T\mathbf{x}] = \text{trace}(\Sigma) + \mathbf{c}^T\mathbf{c}$, we have:

$$\begin{aligned} \text{Var}[\tilde{\mathbf{x}}_{\text{Drop}}] &= \Sigma + \frac{1-p}{p} \text{diag}(\Sigma + \mathbf{c}\mathbf{c}^T) \\ \text{Var}[\tilde{\mathbf{x}}_{\text{Rot}}] &= \Sigma + \frac{1-p}{p} \frac{\text{trace}(\Sigma)\mathbf{I} - \Sigma + \mathbf{c}^T\mathbf{c}\mathbf{I} - \mathbf{c}\mathbf{c}^T}{D-1} \end{aligned} \quad (14)$$

We can compute the co-adaptations of $\tilde{\mathbf{x}}$ (Assume $\mathbf{c} = \mathbf{0}$):

$$\begin{aligned} \text{co}(\tilde{\mathbf{x}}_{\text{Drop}}) &= \frac{\|\tilde{\mathbf{x}}_{\text{Drop}} - \text{diag}(\tilde{\mathbf{x}}_{\text{Drop}})\|_1}{\text{trace}(\tilde{\mathbf{x}}_{\text{Drop}})} = \frac{\|\Sigma - \text{trace}(\Sigma)\|_1}{\frac{1}{p}\text{trace}(\Sigma)} = p \text{co}(\mathbf{x}) \\ \text{co}(\tilde{\mathbf{x}}_{\text{Rot}}) &= \frac{(1 - \frac{1-p}{p(D-1)})\|\Sigma - \text{trace}(\Sigma)\|_1}{\text{trace}(\Sigma) + \frac{1-p}{p} \frac{D\text{trace}(\Sigma) - \text{trace}(\Sigma)}{D-1}} = (p - \frac{1-p}{D-1}) \text{co}(\mathbf{x}) \end{aligned} \quad (15)$$

Under zero-center assumption, Dropout with keep rate p reduces co-adaptation by p times, and the equivalent RotationOut reduces co-adaptation by $p - \frac{1-p}{D-1}$ times.

We take a close look at the correlation coefficient to see what makes the difference. Let \mathbf{x}_i be the i^{th} element of \mathbf{x} . Recall Equation 13, we have:

$$|\text{cor}(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j)| = \frac{|\text{cov}(\mathbf{x}_i, \mathbf{x}_j) + \mathbb{E}[\text{cov}(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j|\mathbf{x})]|}{\sqrt{(\text{Var}[\mathbf{x}_i] + \mathbb{E}[\text{Var}[\tilde{\mathbf{x}}_i|\mathbf{x}]]) (\text{Var}[\mathbf{x}_j] + \mathbb{E}[\text{Var}[\tilde{\mathbf{x}}_j|\mathbf{x}]])}} \quad (16)$$

For Dropout-and other dropout-like methods, they add noise to different neurons independently, so $\text{cov}(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j|\mathbf{x}) = 0$. The only term to reduce correlation coefficients in Equation 16 is $\mathbb{E}[\text{Var}[\tilde{\mathbf{x}}_j|\mathbf{x}]]$. Under out non-trivial noise assumption, $\text{Var}[\tilde{\mathbf{x}}_j|\mathbf{x}]$ is always positive. Thus non-trivial noise can always reduce co-adaptations. For RotationOut, there is another term to reduce correlation coefficients: $\mathbb{E}[\text{cov}(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j|\mathbf{x})] = -\frac{1-p}{p(D-1)}\text{cov}(\mathbf{x}_i, \mathbf{x}_j)$ and typically $0 < \frac{1-p}{p(D-1)} < 1$. In addition to increasing the uncertainty of each neuron as Dropout does, RotationOut can also reduce the correlation between two neurons. In other words, *inhibition noise*.

Here we explain the first reason why we need a zero-center assumption and rotate the zero-centered features in Section 2.3. Equation 14 and 16 show that the non-zero mean value can further reduce the co-adaptations. If we do not know the exact mean value, we do not know the exact regularization strength. Suppose the neurons $x \sim \mathcal{N}(0, 1)$ follow a normal distribution, and we apply Dropout on the ReLU activations $y = \text{ReLU}(x)$. With a keep rate 0.9, Dropout reduces the co-adaptations by 0.86 times, while Dropout reduces the co-adaptations by 0.61 times with a keep rate 0.7, which is a non-linear mapping and influenced by the mean value. We rotate/drop the zero-centered features so that the regularization strength is independent with the mean value.

BN also introduces noise to the neurons by using the batch mean and variance. The noise to different neurons/channels are independent, so the effect of BN's noise is similar to Dropout. It is widely believed that the noise causes BN performance to decrease with small batch size (Wu & He, 2018; Luo et al., 2018). However, Dropout usually decrease the performance when the keep rate is very low. We study the effect of BN's noise and argue that BN is not a linear operation. The non-linearity increases when the batch size decreases, which is also one reason for the small batch size BN's performance drop. Details see Appendix A.4.

3.2 DROPOUT-BEFORE BATCH NORMALIZATION

Dropout changes the variance of a specific neuron when transferring the network from training to inference. However, BN requires a consistent statistical variance. The variance inconsistency

(variance shift) in training and inference leads to unstable numerical behaviors and more erroneous predictions when applying Dropout-before BN.

We can easily understand this using Equation 13. If a Dropout layer is applied right before a BN layer. In training time, the BN layer records the diagonal element of $\text{Var}[\tilde{\mathbf{x}}]$ as the running variance and uses them in inference. However, the actual variance in inference should be the diagonal element of $\text{Var}[\mathbf{x}]$ which is small than the recorded running variance (train variance). Li et al. (2019) argues:

- P1 Instead of using Dropout, a more variance-stable form *Uout* can be used to mitigate the problem: $\tilde{\mathbf{x}}_i = \mathbf{x}_i(1 + r_i)$ where $r_i \sim \text{Unif}[-\beta, \beta]$.
- P2 Instead of applying Dropout-a (Figure 1), applying Dropout-b can mitigate the problem.
- P3 In Dropout-b, let r be the ratio between train variance and test variance. Expanding the input dimension of weight layer D can mitigate the problem: $D \rightarrow \infty, r \rightarrow 1$.

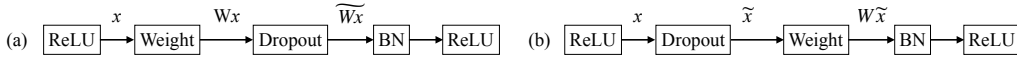


Figure 1: Two types of Dropout. The weight layer can be convolutional or fully connected layer.

We revisit these propositions and discuss how to mitigate the problem. For Proposition 1, *Uout* is unlikely to mitigate the problem. The *Uout* noise to different neurons are independent, so the variance shift is the only term to reduce co-adaptations in Equation 16. Though *Uout* is variance-stable, it provides less regularization, which is equivalent to Dropout with a higher keep rate.

Proposition 2 and 3 discuss the positions to insert Dropout. Let \mathbf{x} be the output from ReLU layer with $\mathbb{E}[\mathbf{x}] = \mathbf{c}$, $\text{Var}[\mathbf{x}] = \Sigma$ and \mathbf{y} be the input of BN layer. The weight layer in Dropout-a and b are the same with weight $\mathbf{W} \in \mathbb{R}^{n \times D}$. During test time, the inputs to BN layers in Dropout-a and b are the same $\mathbf{y} = \mathbf{W}\mathbf{x}$ with variance $\text{Var}[\mathbf{y}] = \mathbf{W}\Sigma\mathbf{W}^T$. During training time, the inputs are different. In Dropout-a, the formulation is $\tilde{\mathbf{y}}_a = \tilde{\mathbf{W}}\mathbf{x}$ where $\mathbb{E}[\tilde{\mathbf{W}}\mathbf{x}|\mathbf{W}\mathbf{x}] = \mathbf{W}\mathbf{x}$. In Dropout-b, the formulation is $\tilde{\mathbf{y}}_b = \mathbf{W}\tilde{\mathbf{x}}$ where $\mathbb{E}[\tilde{\mathbf{x}}|\mathbf{x}] = \mathbf{x}$. So the training variance for the two types are different. Recall Lemma 1, we have:

$$\begin{aligned} \text{Var}[\tilde{\mathbf{y}}_a] &= \mathbb{E} \left[\text{Var}[\tilde{\mathbf{W}}\mathbf{x}|\mathbf{W}\mathbf{x}] \right] + \text{Var}[\mathbf{W}\mathbf{x}] = \mathbf{W}\Sigma\mathbf{W}^T + \frac{1-p}{p} \text{diag}(\mathbf{W}(\Sigma + \mathbf{c}\mathbf{c}^T)\mathbf{W}^T) \\ \text{Var}[\tilde{\mathbf{y}}_b] &= \mathbf{W}(\mathbb{E}[\text{Var}[\tilde{\mathbf{x}}|\mathbf{x}]] + \text{Var}[\mathbf{x}])\mathbf{W}^T = \mathbf{W}\Sigma\mathbf{W}^T + \frac{1-p}{p} \mathbf{W} \text{diag}(\Sigma + \mathbf{c}\mathbf{c}^T)\mathbf{W}^T \end{aligned} \quad (17)$$

Let \mathbf{w} be i^{th} row of \mathbf{W} and assume \mathbf{w}_i is uniformly distributed on the unit ball. Since the length of \mathbf{w} expands the training and testing variance with the same proportion, it does not affect the ratio between training and testing variance, and we can assume the length of \mathbf{w} is fixed. The i^{th} element of actual testing variance is $\mathbf{w}\Sigma\mathbf{w}^T$. For Dropout-a, the i^{th} element of running variance (i.e., the training variance) is $\text{Var}[\tilde{\mathbf{y}}_a|\mathbf{w}]_i = \mathbf{w}\Sigma\mathbf{w}^T + \frac{1-p}{p} \mathbf{w}(\Sigma + \mathbf{c}\mathbf{c}^T)\mathbf{w}^T$. For Dropout-b, the i^{th} element of running variance is $\text{Var}[\tilde{\mathbf{y}}_b|\mathbf{w}]_i = \mathbf{w}\Sigma\mathbf{w}^T + \frac{1-p}{p} \mathbf{w} \text{diag}(\Sigma + \mathbf{c}\mathbf{c}^T)\mathbf{w}^T$. Dropout-a and b have the same expected variance shift:

$$\mathbb{E}_{\mathbf{w}}[\text{Var}[\tilde{\mathbf{y}}_a|\mathbf{w}]_i - \text{Var}[\mathbf{y}|\mathbf{w}]_i] = \mathbb{E}_{\mathbf{w}}[\text{Var}[\tilde{\mathbf{y}}_b|\mathbf{w}]_i - \text{Var}[\mathbf{y}|\mathbf{w}]_i] = \frac{1-p}{p}(\text{trace}(\Sigma) + \mathbf{c}^T\mathbf{c}) \quad (18)$$

Though the expected variance shift is the same, the variance of the shift is different. Let $r(\mathbf{w})$ be the ratio between the training variance and the testing variance: $r(\mathbf{w}) = \text{Var}[\tilde{\mathbf{y}}|\mathbf{w}]_i / \text{Var}[\mathbf{y}|\mathbf{w}]_i$. We have the following observation:

Observation. If $\mathbf{c} > 0$ which is the case that the activation function is ReLU. The ratio in Dropout-b is more centered: $\text{Var}_{\mathbf{w}}[r_b(\mathbf{w})] < \text{Var}_{\mathbf{w}}[r_a(\mathbf{w})] = o(\frac{1}{D})$. Sample n weights to make the weight layer \mathbf{W} , the maximum ratio in Dropout-a is bigger than the maximum ratio in Dropout-b with high probability: $\max_{1 \leq k \leq n} r_b(\mathbf{w}_i) < \max_{1 \leq k \leq n} r_a(\mathbf{w}_i)$.

According to this observation, Proposition 2 and 3 are basically right but might not be precise. Dropout-b does help mitigate the problem but there might be other reasons. The expected variance

shift is the same in Dropout-a and b: $D \rightarrow \infty, r \rightarrow 1$. Dropout-b has more stable variance shift among different dimensions. Dropout-a is more likely to have very big training/testing variance ratio, leading to more serious unstable numerical behaviour.

Consider zero-centered Dropout-a in Equation 17: $\text{Var}[\tilde{\mathbf{y}}_a] = \mathbf{W}\Sigma\mathbf{W}^T + \frac{1-p}{p}\text{diag}(\mathbf{W}\Sigma\mathbf{W}^T)$. The ratio is fixed to be $1/p$ for any weights, i.e. $\text{Var}_{\mathbf{w}}[r_a(\mathbf{w})] = 0$. It leads to fewer unstable numerical behaviour since there is no extreme variance shift ratio, and we can modify BN layer’s validation mode (reduce the running variance by $1/p$ times). Zero-centered Dropout-a can be one solution to mitigate the variance shift problem.

4 EXPERIMENTS

In this section, we evaluate the performance of RotationOut for image classification, object detection, and speech recognition. First, we conduct detailed ablation studies with CIFAR100 dataset. Next, we compare RotationOut with other regularization techniques using more data and higher resolution. We test on two tasks: image classification on ILSVRC dataset and object detection on COCO dataset. Finally, we show that RotationOut can also help training LSTMs for speech recognition task.

4.1 ABLATION STUDY ON CIFAR100

The CIFAR100 dataset consists of 60,000 colour images of size 32×32 pixels and 100 classes. The official version of the dataset is split into a training set with 50,000 images and a test set with 10,000 images. We conduct image classification experiments on the dataset.

Experiment settings. Our focus is on the regularization abilities, so the experiment settings for different regularization techniques are the same. The network inputs are 32×32 and normalized using per-channel mean and standard deviation. The data augmentation methods are as follows: first zero-pad the images with 4 pixels on each side to obtain a 40×40 pixel image, then randomly crop a 32×32 pixel image, and finally mirror the images horizontally with 50% probability. For all of these experiments, we use the same optimizer: training for 200 epochs with batches of 128 images using SGD, momentum of 0.9, and weight decay of $1e-5$. The learning rate is initially set to 0.1, but is scheduled to decrease by a factor of 5x after each of the 60th, 120th, and 160th epochs. Each experiment is repeated 5 times and we report the top 1 accuracy as “mean \pm standard deviation”.

How RotationOut regularizes network. We first compare the results with and without RotationOut. Here we study how the hyper-parameters of the angle distribution influence the results and compare with the equivalent to Bernoulli Dropout. We experiment on three different distributions: 1) uniform distribution $\theta \sim \text{Unif}(0, \Theta)$, 2) uniform distribution of tangent $\theta \sim \mathcal{N}(0, \Theta^2)$, and 3) normal distribution of tangent $\tan \theta \sim \mathcal{N}(0, \sigma^2)$. As discussed earlier, the regularization strength is controlled by $\mathbb{E} \tan \theta^2 = 1/p - 1$. We compare RotationOut with the corresponding Dropout. We

Table 1: Top 1 accuracy of Dropout and corresponding RotationOut on CIFAR100

(a) Standard Dropout		(b) $\theta \sim \text{Unif}(0, \Theta)$		(c) $\tan \theta \sim \text{Unif}(0, \Theta)$		(d) $\tan \theta \sim \mathcal{N}(0, \sigma^2)$	
keep rate	top-1(%)	Θ	top-1(%)	Θ	top-1(%)	σ	top-1(%)
0	77.67 ± 0.23	0	77.67 ± 0.23	0	77.67 ± 0.23	0	77.67 ± 0.23
0.9	78.38 ± 0.26	0.542	78.41 ± 0.23	0.577	78.53 ± 0.36	0.333	78.58 ± 0.22
0.8	78.67 ± 0.37	0.759	78.81 ± 0.27	0.866	78.91 ± 0.29	0.500	79.12 ± 0.32
0.7	78.80 ± 0.32	0.921	78.61 ± 0.25	1.134	78.50 ± 0.32	0.655	78.60 ± 0.37
0.6	78.43 ± 0.33	1.053	28.24 ± 0.30	1.414	78.24 ± 0.37	0.816	78.11 ± 0.29

use ResNet50 as the backbone network. Since the input size 32 is much smaller than the standard input size, we only stride at the beginning of Res2, Res3 and Res4. The tensor size before global pooling is $2048 \times 4 \times 4$. We apply RotationOut to convolutional layers in Res3 and Res4. The results are shown in Table 1. The rotation angles play similar roles as the dropout rate to control the regularization strength. In a proper range, the model can have better generalization ability. But

the performance would be affected if we rotate the features too much. RotationOut with normal distribution of tangent performs best with a 1.45% improvement compared with the baseline and a 0.32% improvement compared with the best DropOut performance.

4.2 EXPERIMENTS WITH MORE DATA AND HIGHER RESOLUTION

ImageNet Classification. The ILSVRC 2012 classification dataset contains 1.2 million training images and 50,000 validation images with 1,000 categories. We following the training and test schema as in (Szegedy et al., 2015; He et al., 2016) but train the model for 240 epochs. The learning rate is decayed by the factor of 0.1 at 120, 190 and 230 epochs. We apply RotationOut with normal distribution of tangent $\mathbb{E} \tan \theta^2 = 1/4$ to convolutional layers in Res3 and Res4 as well as the last fully connected layer. As mentioned earlier, RotationOut is easily combined with DropBlock idea. We rotate features in a continuous block size of 7×7 in Res3 and 3×3 in Res4.

Table 2 shows the results of some state of the art methods and our results. Our results are average over 5 runs. Results of other methods are from Ghiasi et al. (2018). Our result is significantly better than Dropout and SpatialDropout. By using the DropBlock idea, RotationOut can get competitive results compared with state of the art methods and get a 2.07% improvement compared with the baseline.

Table 2: Comparison with state of the art: Top 1 accuracy of ResNet50 on ImageNet Validation

Model	top-1(%)	top-5(%)
ResNet-50	76.51 ± 0.07	93.20 ± 0.05
ResNet-50 + dropout(kp=0.7)	76.80 ± 0.04	93.41 ± 0.04
ResNet-50 + DropPath(kp=0.9)	77.10 ± 0.08	93.50 ± 0.05
ResNet-50 + SpatialDropout(kp=0.9)	77.41 ± 0.04	93.74 ± 0.02
ResNet-50 + Cutout	76.52 ± 0.07	93.21 ± 0.04
ResNet-50 + DropBlock(kp=0.9)	78.13 ± 0.05	94.02 ± 0.02
ResNet-50 + RotationOut	77.87 ± 0.34	93.94 ± 0.17
ResNet-50 + RotationOut (Block)	78.58 ± 0.48	94.27 ± 0.24

COCO Object Detection. Our proposed method can also be used in other vision tasks, for example Object Detection on MS COCO (Lin et al., 2014). In this task, we use RetinaNet (Lin et al., 2017) as the detection method and apply RotationOut to the ResNet backbone. We use the same hyper parameters as in ImageNet classification. We follow the implementation details in (Ghiasi et al., 2018): resize images between scales [512, 768] and then crop the image to max dimension 640. The model are initialized with ImageNet pretraining and trained for 35 epochs with learning decay at 20 and 28 epochs. We set $\alpha = 0.25$ and $\gamma = 1.5$ for focal loss, a weight decay of 0.0001, a momentum of 0.9 and a batch size of 64. The model is trained on COCO train2017 and evaluated on COCO val2017. We compare our result with DropBlock as table 3 shows.

Table 3: Object detection in COCO using RetinaNet and ResNet-50 FPN backbone

Model	Initialization	AP	AP50	AP75
RetinaNet	ImageNet	36.5	55.0	39.1
RetinaNet, no DropBlock	Random	36.8	54.6	39.4
RetinaNet, Dropout, keep_prob = 0.9	Random	37.9	56.1	40.6
RetinaNet, keep_prob = 0.9, block_size = 5	Random	38.4	56.4	41.2
RetinaNet, RotationOut	ImageNet	38.2	56.2	41.0
RetinaNet, RotationOut (Block)	ImageNet	38.7	56.6	41.4

Due to limited computing resources, we finetune the model from ImageNet classification models while DropBlock method trained the model from scratch. We think it is fair to compare DropBlock

method since the initialization does not help increase the results as showed in the first two rows. Our RotationOut can still have additional 0.3 AP based on the DropBlock result.

4.3 EXPERIMENT IN SPEECH RECOGNITION

We show that our RotationOut can also help train LSTMs. We conduct an Auto2Text experiment on the WSJ (Wall Street Journal) dataset (Paul & Baker, 1992). The dataset is a database with 80 hours of transcribed speech. The inputs are variable length speech $\mathbf{X} \in \mathbb{R}^{T \times L}$ where T is the length and L is the feature dimension for one time step. The labels are character-based words. We use a four-layer bidirectional LSTM network to design a CTC (Connectionist temporal classification) Graves et al. (2006) model. The input dimension, hidden dimension and output dimension of the four-layer bidirectional LSTM network are 40, 512, 137 respectively. We use Adam optimizer with learning rate 1e-3, weight decay 1e-5 and batch size 32, and train the model for 80 epochs and reduce the learning rate by 5x at epoch 40. We report the edit distance between our prediction and ground truth on the “eval92” test set. Table 4 shows the performance of different regularization methods.

Table 4: Auto2Text experiment on the WSJ

Method	Distance
No regularization	9.1
Standard Dropout(kp=0.9)	8.6
Weight Drop(kp=0.8)	7.8
Variational Weight Drop(kp=0.8)	7.5
Locked Drop(kp=0.7)	7.3
Locked Drop(kp=0.8)+Variational Weight Drop(kp=0.8)	6.7
RotationOut	6.8
RotationOut +Variational Weight Drop(kp=0.9)	6.4

5 CONCLUSION

In this work, we introduce RotationOut as an alternative for dropout for neural network. RotationOut adds continuous noise to data/features and keep the semantics. We further establish an analysis of noise to show how co-adaptations are reduced in neural network and why dropout is more effective than dropout. The analysis of noise provide a clear view of the dropout-BN architecture, indicating how to use dropout with batch normalization. Our experiments show that applying RotationOut in neural network helps training and increase the accuracy. For further work, we would reduce the running time of RotationOut. Though the proposed method has the same complexity as Dropout, it is still much slower than Dropout. Another direction is the theoretical analysis of co-adaptations. As discussed earlier, the proposed correlation analysis is not optimal. It cannot explain the difference between standard Dropout and Gaussian dropout. And it can not explain some methods such as Shake-shake regularization. Further work on co-adaptation analysis can help better understand noise-based regularization methods.

REFERENCES

- Dario Amodei, Sundaram Ananthanarayanan, Rishita Anubhai, Jingliang Bai, Eric Battenberg, Carl Case, Jared Casper, Bryan Catanzaro, Qiang Cheng, Guoliang Chen, et al. Deep speech 2: End-to-end speech recognition in english and mandarin. In *International conference on machine learning*, pp. 173–182, 2016.
- Edward Anderson. Discontinuous plane rotations and the symmetric eigenvalue problem. 2000.
- David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6541–6549, 2017.

- Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.
- Yarin Gal and Zoubin Ghahramani. A theoretically grounded application of dropout in recurrent neural networks. In *Advances in neural information processing systems*, pp. 1019–1027, 2016.
- Xavier Gastaldi. Shake-shake regularization. *arXiv preprint arXiv:1705.07485*, 2017.
- Golnaz Ghiasi, Tsung-Yi Lin, and Quoc V Le. Dropblock: A regularization method for convolutional networks. In *Advances in Neural Information Processing Systems*, pp. 10727–10737, 2018.
- Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd international conference on Machine learning*, pp. 369–376. ACM, 2006.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012.
- Ronghang Hu, Huazhe Xu, Marcus Rohrbach, Jiashi Feng, Kate Saenko, and Trevor Darrell. Natural language object retrieval. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4555–4564, 2016.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708, 2017.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- Alex Labach, Hojjat Salehinejad, and Shahrokh Valaee. Survey of dropout methods for deep neural networks. *arXiv preprint arXiv:1904.13310*, 2019.
- Karel Lenc and Andrea Vedaldi. Understanding image representations by measuring their equivariance and equivalence. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 991–999, 2015.
- Xiang Li, Shuo Chen, Xiaolin Hu, and Jian Yang. Understanding the disharmony between dropout and batch normalization by variance shift. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2682–2690, 2019.
- Zhe Li, Boqing Gong, and Tianbao Yang. Improved dropout for shallow and deep learning. In *Advances in Neural Information Processing Systems*, pp. 2523–2531, 2016.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pp. 740–755. Springer, 2014.
- Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pp. 2980–2988, 2017.
- Ping Luo, Xinjiang Wang, Wenqi Shao, and Zhanglin Peng. Towards understanding regularization in batch normalization. In *arXiv preprint arXiv:1809.00846*, pp. 17, 2018.
- Diego Marcos, Michele Volpi, Nikos Komodakis, and Devis Tuia. Rotation equivariant vector field networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 5048–5057, 2017.

- Stephen Merity, Nitish Shirish Keskar, and Richard Socher. Regularizing and optimizing lstm language models. *arXiv preprint arXiv:1708.02182*, 2017.
- Taesup Moon, Heeyoul Choi, Hoshik Lee, and Inchul Song. Rnndrop: A novel dropout for rnns in asr. In *2015 IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*, pp. 65–70. IEEE, 2015.
- Douglas B Paul and Janet M Baker. The design for the wall street journal-based csr corpus. In *Proceedings of the workshop on Speech and Natural Language*, pp. 357–362. Association for Computational Linguistics, 1992.
- Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic routing between capsules. In *Advances in neural information processing systems*, pp. 3856–3866, 2017.
- Tim Salimans and Durk P Kingma. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In *Advances in Neural Information Processing Systems*, pp. 901–909, 2016.
- Ashish Shrivastava, Tomas Pfister, Oncel Tuzel, Joshua Susskind, Wenda Wang, and Russell Webb. Learning from simulated and unsupervised images through adversarial training. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2107–2116, 2017.
- Patrice Y Simard, David Steinkraus, John C Platt, et al. Best practices for convolutional neural networks applied to visual document analysis. In *Icdar*, volume 3, 2003.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- 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*, pp. 1–9, 2015.
- Jonathan Tompson, Ross Goroshin, Arjun Jain, Yann LeCun, and Christoph Bregler. Efficient object localization using convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 648–656, 2015.
- Li Wan, Matthew Zeiler, Sixin Zhang, Yann Le Cun, and Rob Fergus. Regularization of neural networks using dropout. In *International conference on machine learning*, pp. 1058–1066, 2013.
- Daniel E Worrall, Stephan J Garbin, Daniyar Turmukhambetov, and Gabriel J Brostow. Harmonic networks: Deep translation and rotation equivariance. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5028–5037, 2017.
- Haibing Wu and Xiaodong Gu. Towards dropout training for convolutional neural networks. *Neural Networks*, 71:1–10, 2015.
- Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 3–19, 2018.
- Wojciech Zaremba, Ilya Sutskever, and Oriol Vinyals. Recurrent neural network regularization. *arXiv preprint arXiv:1409.2329*, 2014.

A APPENDIX

A.1 RANDOM ROTATION MATRIX

One example of such a matrix that rotates the (1, 3) dimensions and (2, 4) dimensions can be:

$$\mathbf{M}(\theta, [3, 2, 1, 4]) = \begin{bmatrix} \cos \theta & 0 & -\sin \theta & 0 \\ 0 & \cos \theta & 0 & \sin \theta \\ \sin \theta & 0 & \cos \theta & 0 \\ 0 & -\sin \theta & 0 & \cos \theta \end{bmatrix}. \quad (19)$$

In Section 2, we mentioned the complexity of RotationOut is $O(D)$. It is because we can avoid matrix multiplications to get $\mathcal{R}\mathbf{x}$. For example, let the \mathcal{R} be the operator generated by Equation 19, we have:

$$\mathcal{R}\mathbf{x} = \mathbf{x} + \tan \theta \begin{bmatrix} 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \end{bmatrix} \mathbf{x}. \quad (20)$$

The sparse matrix in Equation 20 is similar to a combine of permutation matrix, and we do not need matrix multiplications to get the output. The output can be get by slicing and an element wise multiplication: $\mathbf{x}[3, 4, 1, 2] * [-1, 1, 1, -1]$.

A.2 MARGINALIZING LINEAR REGRESSION

Recall that $\mathbb{E}_{\mathcal{R}} = \mathbf{I}$, the marginalizing linear regression expression:

$$\begin{aligned} & \mathbb{E}_{\mathcal{R}} \left[\sum_{i=1}^N (y_i - \mathbf{w}^T \mathcal{R}_i \mathbf{x}_i)^2 \right] \\ &= \mathbb{E}_{\mathcal{R}} \left[\sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i + \mathbf{w}^T (\mathbf{I} - \mathcal{R}_i) \mathbf{x}_i)^2 \right] \\ &= \left[\sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \right] + \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i) \mathbb{E}_{\mathcal{R}} \left[\mathbf{w}^T (\mathbf{I} - \mathcal{R}_i) \mathbf{x}_i \right] + \mathbf{w}^T \mathbb{E}_{\mathcal{R}} \left[(\mathbf{I} - \mathcal{R}_i) \mathbf{x}_i \mathbf{x}_i^T (\mathbf{I} - \mathcal{R}_i)^T \right] \mathbf{w} \\ &= \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \mathbf{w}^T \sum_{i=1}^N \text{Var}_{\mathcal{R}}[(\mathbf{I} - \mathcal{R}_i) \mathbf{x}_i] \mathbf{w} \\ &= \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \mathbf{w}^T \sum_{i=1}^N \text{Var}_{\mathcal{R}}[\mathcal{R}_i \mathbf{x}_i] \mathbf{w} \end{aligned} \quad (21)$$

From Lemma one, we have $\text{Var}_{\mathcal{R}}[\mathcal{R}_i \mathbf{x}_i] = \frac{1-p}{pD-p} (\mathbf{x}_i^T \mathbf{x}_i \mathbf{I} - \mathbf{x}_i \mathbf{x}_i^T)$. Write the second term of Equation 21 in the matrix form, we can get Equation 5.

A.3 PROOF OF LEMMA 1

The Dropout form is trivial. We consider the RotationOut equation. Denote \mathbf{x}^i as the i^{th} term of \mathbf{x} . The probability distribution of each element of $\tilde{\mathbf{x}}_{\text{Rot}}$ is:

$$\forall j \neq i, \mathbb{P}(\tilde{\mathbf{x}}_{\text{Rot}}^i = \mathbf{x}^i + \tan \theta \mathbf{x}^j) = \mathbb{P}(\tilde{\mathbf{x}}_{\text{Rot}}^i = \mathbf{x}^i - \tan \theta \mathbf{x}^j) = \frac{1}{2(D-1)} \quad (22)$$

The joint distribution of each two elements of $\tilde{\mathbf{x}}_{\text{Rot}}$ is:

$$\begin{aligned} & \forall m \neq i, n \neq j, m \neq n : \mathbb{P}(\tilde{\mathbf{x}}_{\text{Rot}}^i = \mathbf{x}^i + \tan \theta \mathbf{x}^m, \tilde{\mathbf{x}}_{\text{Rot}}^j = \mathbf{x}^j + \tan \theta \mathbf{x}^n) \\ &= \mathbb{P}(\tilde{\mathbf{x}}_{\text{Rot}}^i = \mathbf{x}^i + \tan \theta \mathbf{x}^m, \tilde{\mathbf{x}}_{\text{Rot}}^j = \mathbf{x}^j - \tan \theta \mathbf{x}^n) \\ &= \mathbb{P}(\tilde{\mathbf{x}}_{\text{Rot}}^i = \mathbf{x}^i - \tan \theta \mathbf{x}^m, \tilde{\mathbf{x}}_{\text{Rot}}^j = \mathbf{x}^j - \tan \theta \mathbf{x}^n) \\ &= \mathbb{P}(\tilde{\mathbf{x}}_{\text{Rot}}^i = \mathbf{x}^i - \tan \theta \mathbf{x}^m, \tilde{\mathbf{x}}_{\text{Rot}}^j = \mathbf{x}^j + \tan \theta \mathbf{x}^n) \\ &= \frac{1}{4(D-1)(D-3)} \\ & \forall i \neq j : \mathbb{P}(\tilde{\mathbf{x}}_{\text{Rot}}^i = \mathbf{x}^i + \tan \theta \mathbf{x}^j, \tilde{\mathbf{x}}_{\text{Rot}}^j = \mathbf{x}^j - \tan \theta \mathbf{x}^i) \\ & \quad \mathbb{P}(\tilde{\mathbf{x}}_{\text{Rot}}^i = \mathbf{x}^i - \tan \theta \mathbf{x}^j, \tilde{\mathbf{x}}_{\text{Rot}}^j = \mathbf{x}^j + \tan \theta \mathbf{x}^i) \\ &= \frac{1}{2(D-1)} \end{aligned} \quad (23)$$

So we have:

$$\mathbb{E}\tilde{\mathbf{x}}_{\text{Rot}}^i = 0, \mathbb{E}(\tilde{\mathbf{x}}_{\text{Rot}}^i)^2 = \frac{\mathbb{E}_\theta \tan^2 \theta}{D-1} \sum_{j \neq i} (\mathbf{x}^j)^2, \mathbb{E}\tilde{\mathbf{x}}_{\text{Rot}}^i \tilde{\mathbf{x}}_{\text{Rot}}^j = -\frac{\mathbb{E}_\theta \tan^2 \theta}{D-1} \mathbf{x}^i \mathbf{x}^j \quad (24)$$

A.4 RETHINKING SMALL BATCH SIZE BATCH NORMALIZATION

Let $\{x_i\}_{i=1}^D$ be one dimension of the data where D is the dataset size. During mini-batch training, one batch of B data $\{x_{b_k}\}_{k=1}^B$ is sampled and the BN operation can be formulate as:

$$\begin{aligned} \mu_{\mathcal{B}} &= \frac{1}{B} \sum_{k=1}^B x_{b_k}, \\ \sigma_{\mathcal{B}}^2 &= \frac{1}{B} \sum_{k=1}^B (x_{b_k} - \mu_{\mathcal{B}})^2, \\ \hat{x}_{b_k} &= \frac{x_{b_k} - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}, \\ \hat{y}_{b_k} &= \gamma \cdot \hat{x}_{b_k} + \beta. \end{aligned} \quad (25)$$

The batch normalization operation records a running mean $\mu_{\mathcal{B}}$ and running variance $\sigma_{\mathcal{B}}^2$ to be used in testing:

$$\begin{aligned} E[x] &= \mathbb{E}_{\mathcal{B}} \mu_{\mathcal{B}}, \\ Var[x] &= \frac{B}{B-1} \mathbb{E}_{\mathcal{B}} \sigma_{\mathcal{B}}^2, \\ \hat{x}_{b_k} &= \frac{x_{b_k} - E[x]}{\sqrt{Var[x] + \epsilon}}, \end{aligned} \quad (26)$$

We want to check whether the test mode formula can be a good estimation of the training mode formula. Suppose we have a batch of data $\{x_k\}_{k=1}^B$. Denote:

$$\mu_{\mathcal{B}} = \frac{1}{B} \sum_{k=1}^B x_k, \sigma_{\mathcal{B}}^2 = \frac{1}{B} \sum_{k=1}^B (x_k - \mu_{\mathcal{B}})^2, \mu_{\mathcal{B}-1} = \frac{1}{B-1} \sum_{k=2}^B x_k, \sigma_{\mathcal{B}-1}^2 = \frac{1}{B-1} \sum_{k=2}^B (x_k - \mu_{\mathcal{B}-1})^2$$

We have:

$$\frac{x_1 - \mu_{\mathcal{B}}}{\sigma_{\mathcal{B}}} = \sqrt{\frac{B-1}{B}} \frac{x_1 - \mu_{\mathcal{B}-1}}{\sqrt{\sigma_{\mathcal{B}-1}^2 + \frac{1}{B}(x_1 - \mu_{\mathcal{B}-1})^2}} \quad (27)$$

Note that $\mu_{\mathcal{B}-1}$ and $\sigma_{\mathcal{B}-1}^2$ are independent from x_1 . So the expected output of any x is:

$$\mathbb{E}[\text{Normalize}(x)] = \mathbb{E}_{\mu_{\mathcal{B}-1}, \sigma_{\mathcal{B}-1}^2} \sqrt{\frac{B-1}{B}} \frac{x - \mu_{\mathcal{B}-1}}{\sqrt{\sigma_{\mathcal{B}-1}^2 + \frac{1}{B}(x - \mu_{\mathcal{B}-1})^2}} \quad (28)$$

Let the function in 28 be $f(x, B)$. Easy to know that it is not a linear function (but BN assumes it should be $y = x$!). Suppose the data follows normal distribution, we can plot $f(x, B)$ by Monte Carlo sampling: Figure 2 shows that BN is not a linear operation. The non-linearity increases when the batch size decreases. It is another important reason for the small batch size BN's performance drop.

To validate our conclusion, we propose cross normalization. For each data in the batch, cross normalization uses the sample mean and variance except itself to compute its normalization mean and

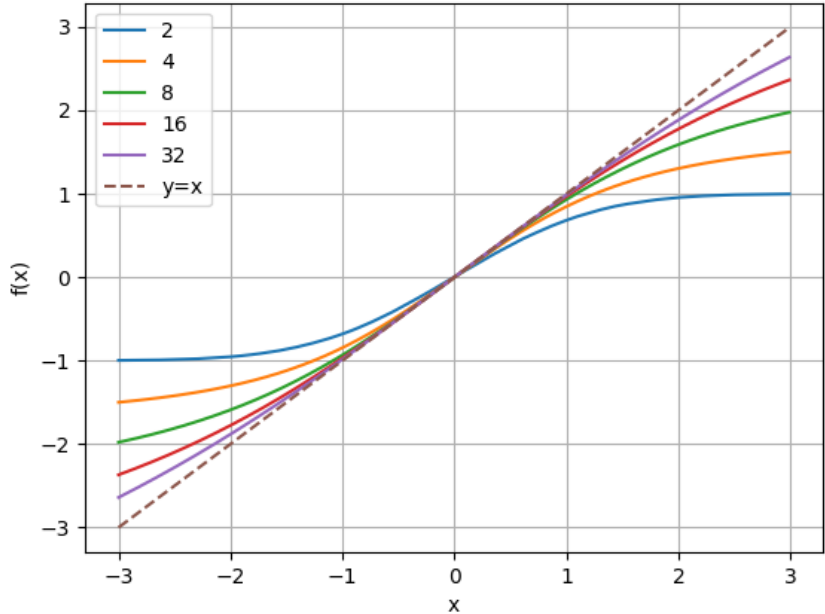


Figure 2: Values of $f(x, B)$ against different values of B .

variance:

$$\begin{aligned}
 \mu_i &= \frac{1}{B} \sum_{k \neq i}^B x_{b_k}, \\
 \sigma_i^2 &= \frac{1}{B} \sum_{k \neq i}^B (x_{b_k} - \mu_i)^2, \\
 \hat{x}_{b_i} &= \frac{x_{b_i} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}, \\
 \hat{y}_{b_i} &= \gamma \cdot \hat{x}_{b_i} + \beta.
 \end{aligned}
 \tag{29}$$

In this case, the expectation of operation on any data is strictly linear, but it uses less data.

We do not intend to propose a better alternative for BN (cross normalization is even more noisy than BN) but want to check whether the non-linearity is an important issue for BN when batch size is small. If cross normalization can outperform BN in batch size case, then the non-linearity is definitely an important issue. On CIFAR100 dataset, following the settings in our ablation study, ResNet50 with cross normalization has lower test loss when the batch size is 8 and 16. But the test accuracy is the small since cross normalization leads to higher variance.