

Gradient Centralization: A New Optimization Technique for Deep Neural Networks

Hongwei Yong^{1,2}, Jianqiang Huang², Xiansheng Hua², and Lei Zhang^{1,2}

¹ Department of Computing, The Hong Kong Polytechnic University
{cshyong, cslzhang}@comp.polyu.edu.hk

² DAMO Academy, Alibaba Group
{jianqiang.jqh, huaxiansheng}@gmail.com

Abstract. Optimization techniques are of great importance to effectively and efficiently train a deep neural network (DNN). It has been shown that using the first and second order statistics (e.g., mean and variance) to perform Z-score standardization on network activations or weight vectors, such as batch normalization (BN) and weight standardization (WS), can improve the training performance. Different from these existing methods that mostly operate on activations or weights, we present a new optimization technique, namely gradient centralization (GC), which operates directly on gradients by centralizing the gradient vectors to have zero mean. GC can be viewed as a projected gradient descent method with a constrained loss function. We show that GC can regularize both the weight space and output feature space so that it can boost the generalization performance of DNNs. Moreover, GC improves the Lipschitzness of the loss function and its gradient so that the training process becomes more efficient and stable. GC is very simple to implement and can be easily embedded into existing gradient based DNN optimizers with only one line of code. It can also be directly used to fine-tune the pre-trained DNNs. Our experiments on various applications, including general image classification, fine-grained image classification, detection and segmentation, demonstrate that GC can consistently improve the performance of DNN learning. The code of GC can be found at <https://github.com/Yonghongwei/Gradient-Centralization>.

Keywords: Deep network optimization, gradient descent

1 Introduction

The broad success of deep learning largely owes to the recent advances on large-scale datasets [43], powerful computing resources (e.g., GPUs and TPUs), sophisticated network architectures [15,16] and optimization algorithms [4,24]. Among these factors, the efficient optimization techniques, such as stochastic gradient descent (SGD) with momentum [38], Adagrad [10] and Adam [24], make it possible to train very deep neural networks (DNNs) with a large-scale dataset, and consequently deliver more powerful and robust DNN models in practice. The generalization performance of the trained DNN models as well as the efficiency of training process depend essentially on the employed optimization techniques.

There are two major goals for a good DNN optimizer: accelerating the training process and improving the model generalization capability. The first goal aims to spend less time and cost to reach a good local minima, while the second goal aims to ensure that the learned DNN model can make accurate predictions on test data. A variety of optimization algorithms [38,10,24,10,24] have been proposed to achieve these goals. SGD [4,5] and its extension SGD with momentum (SGDM) [38] are among the most commonly used ones. They update the parameters along the opposite direction of their gradients in one training step. Most of the current DNN optimization methods are based on SGD and improve SGD to better overcome the gradient vanishing or explosion problems. A few successful techniques have been proposed, such as weight initialization strategies [11,14], efficient active functions (e.g., ReLU [35]), gradient clipping [36,37], adaptive learning rate optimization algorithms [10,24], and so on.

In addition to the above techniques, the sample/feature statistics such as mean and variance can also be used to normalize the network activations or weights to make the training process more stable. The representative methods operating on activations include batch normalization (BN) [19], instance normalization (IN) [47,18], layer normalization (LN) [29] and group normalization (GN) [51]. Among them, BN is the most widely used optimization technique which normalizes the features along the sample dimension in a mini-batch for training. BN smooths the optimization landscape [45] and it can speed up the training process and boost model generalization performance when a proper batch size is used [53,15]. However, BN works not very well when the training batch size is small, which limits its applications to memory consuming tasks, such as object detection [13,42], video classification [21,2], etc.

Another line of statistics based methods operate on weights. The representative ones include weight normalization (WN) [44,17] and weight standardization (WS) [39]. These methods re-parameterize weights to restrict weight vectors during training. For example, WN decouples the length of weight vectors from their direction to accelerate the training of DNNs. WS uses the weight vectors' mean and variance to standardize them to have zero mean and unit variance. Similar to BN, WS can also smooth the loss landscape and speed up training. Nevertheless, such methods operating on weight vectors cannot directly adopt the pre-trained models (e.g., on ImageNet) because their weights may not meet the condition of zero mean and unit variance.

Different from the above techniques which operate on activations or weight vectors, we propose a very simple yet effective DNN optimization technique, namely gradient centralization (GC), which operates on the gradients of weight vectors. As illustrated in Fig. 1(a), GC simply centralizes the gradient vectors to have zero mean. It can be easily embedded into the current gradient based optimization algorithms (e.g., SGDM [38], Adam [24]) using only one line of code. Though simple, GC demonstrates various desired properties, such as accelerating the training process, improving the generalization performance, and the compatibility for fine-tuning pre-trained models. The main contributions of this paper are highlighted as follows:

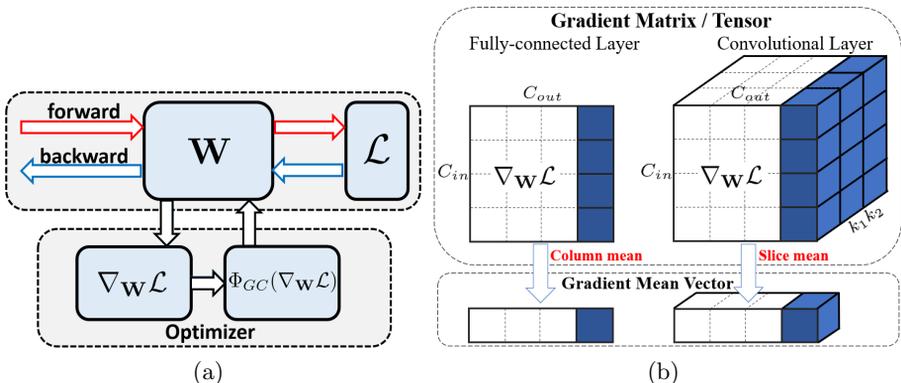


Fig. 1. (a) Sketch map for using gradient centralization (GC). \mathbf{W} is the weight, \mathcal{L} is the loss function, $\nabla_{\mathbf{W}}\mathcal{L}$ is the gradient of weight, and $\Phi_{GC}(\nabla_{\mathbf{W}}\mathcal{L})$ is the centralized gradient. It is very simple to embed GC into existing network optimizers by replacing $\nabla_{\mathbf{W}}\mathcal{L}$ with $\Phi_{GC}(\nabla_{\mathbf{W}}\mathcal{L})$. (b) Illustration of the GC operation on gradient matrix/tensor of weights in the fully-connected layer (left) and convolutional layer (right). GC computes the column/slice mean of gradient matrix/tensor and centralizes each column/slice to have zero mean.

- We propose a new general network optimization technique, namely gradient centralization (GC), which can not only smooth and accelerate the training process of DNN but also improve the model generalization performance.
- We analyze the theoretical properties of GC, and show that GC constrains the loss function by introducing a new constraint on weight vector, which regularizes both the weight space and output feature space so that it can boost model generalization performance. Besides, the constrained loss function has better Lipschitzness than the original one, which makes the training process more stable and efficient.

Finally, we perform comprehensive experiments on various applications, including general image classification, fine-grained image classification, object detection and instance segmentation. The results demonstrate that GC can consistently improve the performance of learned DNN models in different applications. It is a simple, general and effective network optimization method.

2 Related Work

In order to accelerate the training and boost the generalization performance of DNNs, a variety of optimization techniques [19,51,44,39,38,36] have been proposed to operate on activation, weight and gradient. In this section we briefly review the related work from these three aspects.

Activation: The activation normalization layer has become a common setting in DNN, such as batch normalization (BN) [19] and group normalization (GN) [51]. BN was originally introduced to solve the internal covariate shift by normalizing the activations along the sample dimension. It allows higher learning

rates [3], accelerates the training speed and improves the generalization accuracy [33,45]. However, BN does not perform well when the training batch size is small, and GN is proposed to address this problem by normalizing the activations or feature maps in a divided group for each input sample. In addition, layer normalization (LN) [29] and instance normalization (IN) [47,18] have been proposed for RNN and style transfer learning, respectively.

Weight: Weight normalization (WN) [44] re-parameterizes the weight vectors and decouples the length of a weight vector from its direction. It speeds up the convergence of SGDM algorithm to a certain degree. Weight standardization (WS) [39] adopts the Z-score standardization to re-parameterize the weight vectors. Like BN, WS can also smooth the loss landscape and improve training speed. Besides, binarized DNN [40,9,8] quantifies the weight into binary values, which can improve the generalization capability for certain DNNs. However, a shortcoming of those methods operating on weights is that they cannot be directly used to fine-tune pre-trained models since the pre-trained weight may not meet their constraints. As a consequence, we have to design specific pre-training methods for them in order to fine-tune the model.

Gradient: A commonly used operation on gradient is to compute the momentum of gradient [38]. By using the momentum of gradient, SGDM accelerates SGD in the relevant direction and dampens oscillations. Besides, L_2 regularization based weight decay, which introduces L_2 regularization into the gradient of weight, has long been a standard trick to improve the generalization performance of DNNs [27,54]. To make DNN training more stable and avoid gradient explosion, gradient clipping [36,37,1,23] has been proposed to train a very deep DNNs. In addition, the projected gradient methods [12,28] and Riemannian approach [7,49] project the gradient on a subspace or a Riemannian manifold to regularize the learning of weights.

3 Gradient Centralization

3.1 Motivation

BN [19] is a powerful DNN optimization technique, which uses the first and second order statistics to perform Z-score standardization on activations. It has been shown in [45] that BN reduces the Lipschitz constant of loss function and makes the gradients more Lipschitz smooth so that the optimization landscape becomes smoother. WS [39] can also reduce the Lipschitzness of loss function and smooth the optimization landscape through Z-score standardization on weight vectors. BN and WS operate on activations and weight vectors, respectively, and they implicitly constrict the gradient of weights, which improves the Lipschitz property of loss for optimization.

Apart from operating on activation and weight, can we directly operate on gradient to make the training process more effective and stable? One intuitive idea is that we use Z-score standardization to normalize gradient, like what has been done by BN and WS on activation and weight. Unfortunately, we found that normalizing gradient cannot improve the stability of training. Instead, we

propose to compute the mean of gradient vectors and centralize the gradients to have zero mean. As we will see in the following development, the so called gradient centralization (GC) method can have good Lipschitz property, smooth the DNN training and improve the model generalization performance.

3.2 Notations

We define some basic notations. For fully connected layers (FC layers), the weight matrix is denoted as $\mathbf{W}_{fc} \in \mathbb{R}^{C_{in} \times C_{out}}$, and for convolutional layers (Conv layers) the weight tensor is denoted as $\mathbf{W}_{conv} \in \mathbb{R}^{C_{in} \times C_{out} \times (k_1 k_2)}$, where C_{in} is the number of input channels, C_{out} is the number of output channels, and k_1, k_2 are the kernel size of convolution layers. For the convenience of expression, we unfold the weight tensor of Conv layer into a matrix/tensor and use a unified notation $\mathbf{W} \in \mathbb{R}^{M \times N}$ for weight matrix in FC layer ($\mathbf{W} \in \mathbb{R}^{C_{in} \times C_{out}}$) and Conv layers ($\mathbf{W} \in \mathbb{R}^{(C_{in} k_1 k_2) \times C_{out}}$). Denote by $\mathbf{w}_i \in \mathbb{R}^M$ ($i = 1, 2, \dots, N$) the i -th column vector of weight matrix \mathbf{W} and \mathcal{L} the objective function. $\nabla_{\mathbf{W}} \mathcal{L}$ and $\nabla_{\mathbf{w}_i} \mathcal{L}$ denote the gradient of \mathcal{L} w.r.t. the weight matrix \mathbf{W} and weight vector \mathbf{w}_i , respectively. The size of gradient matrix $\nabla_{\mathbf{W}} \mathcal{L}$ is the same as weight matrix \mathbf{W} . Let \mathbf{X} be the input activations for this layer and $\mathbf{W}^T \mathbf{X}$ be its output activations. $\mathbf{e} = \frac{1}{\sqrt{M}} \mathbf{1}$ denotes an M dimensional unit vector and $\mathbf{I} \in \mathbb{R}^{M \times M}$ denotes an identity matrix.

3.3 Formulation of GC

For a FC layer or a Conv layer, suppose that we have obtained the gradient through backward propagation, then for a weight vector \mathbf{w}_i whose gradient is $\nabla_{\mathbf{w}_i} \mathcal{L}$ ($i = 1, 2, \dots, N$), the GC operator, denoted by Φ_{GC} , is defined as follows:

$$\Phi_{GC}(\nabla_{\mathbf{w}_i} \mathcal{L}) = \nabla_{\mathbf{w}_i} \mathcal{L} - \mu_{\nabla_{\mathbf{w}_i} \mathcal{L}} \quad (1)$$

where $\mu_{\nabla_{\mathbf{w}_i} \mathcal{L}} = \frac{1}{M} \sum_{j=1}^M \nabla_{\mathbf{w}_{i,j}} \mathcal{L}$. The formulation of GC is very simple. As shown in Fig. 1(b), we only need to compute the mean of the column vectors of the weight matrix, and then remove the mean from each column vector. We can also have a matrix formulation of Eq. (1):

$$\Phi_{GC}(\nabla_{\mathbf{W}} \mathcal{L}) = \mathbf{P} \nabla_{\mathbf{W}} \mathcal{L}, \quad \mathbf{P} = \mathbf{I} - \mathbf{e} \mathbf{e}^T \quad (2)$$

The physical meaning of \mathbf{P} will be explained later in Section 4.1. In practical implementation, we can directly remove the mean value from each weight vector to accomplish the GC operation. The computation is very simple and efficient.

3.4 Embedding of GC to SGDM/Adam

GC can be easily embedded into the current DNN optimization algorithms such as SGDM [38,5] and Adam [24]. After obtaining the centralized gradient $\Phi_{GC}(\nabla_{\mathbf{W}} \mathcal{L})$, we can directly use it to update the weight matrix. Algorithm 1 and Algorithm 2 show how to embed GC into the two most popular optimization algorithms, SGDM and Adam, respectively. Moreover, if we want to use weight decay, we can set $\hat{\mathbf{g}}^t = \mathbf{P}(\mathbf{g}^t + \lambda \mathbf{w})$, where λ is the weight decay factor. It only

needs to add one line of code into most existing DNN optimization algorithms to execute GC with negligible additional computational cost. For example, it costs only 0.6 sec extra training time in one epoch on CIFAR100 with ResNet50 model in our experiments (71 sec for one epoch).

Algorithm 1 SGDM with Gradient Centralization

Input: Weight vector \mathbf{w}^0 , step size α , momentum factor β , \mathbf{m}^0 Training step: 1: for $t = 1, \dots, T$ do 2: $\mathbf{g}^t = \nabla_{\mathbf{w}^t} \mathcal{L}$	3: $\hat{\mathbf{g}}^t = \Phi_{GC}(\mathbf{g}^t)$ 4: $\mathbf{m}^t = \beta \mathbf{m}^{t-1} + (1 - \beta) \hat{\mathbf{g}}^t$ 5: $\mathbf{w}^{t+1} = \mathbf{w}^t - \alpha \mathbf{m}^t$ 6: end for
---	---

Algorithm 2 Adam with Gradient Centralization

Input: Weight vector \mathbf{w}^0 , step size α , β_1 , β_2 , ϵ , $\mathbf{m}^0, \mathbf{v}^0$ Training step: 1: for $t = 1, \dots, T$ do 2: $\mathbf{g}^t = \nabla_{\mathbf{w}^t} \mathcal{L}$ 3: $\hat{\mathbf{g}}^t = \Phi_{GC}(\mathbf{g}^t)$	4: $\mathbf{m}^t = \beta_1 \mathbf{m}^{t-1} + (1 - \beta_1) \hat{\mathbf{g}}^t$ 5: $\mathbf{v}^t = \beta_2 \mathbf{v}^{t-1} + (1 - \beta_2) \hat{\mathbf{g}}^t \odot \hat{\mathbf{g}}^t$ 6: $\hat{\mathbf{m}}^t = \mathbf{m}^t / (1 - (\beta_1)^t)$ 7: $\hat{\mathbf{v}}^t = \mathbf{v}^t / (1 - (\beta_2)^t)$ 8: $\mathbf{w}^{t+1} = \mathbf{w}^t - \alpha \frac{\hat{\mathbf{m}}^t}{\sqrt{\hat{\mathbf{v}}^t + \epsilon}}$ 9: end for
---	---

4 Properties of GC

As we will see in the section of experimental result, GC can accelerate the training process and improve the generalization performance of DNNs. In this section, we perform theoretical analysis to explain why GC works.

4.1 Improving Generalization Performance

One important advantage of GC is that it can improve the generalization performance of DNNs. We explain this advantage from two aspects: weight space regularization and output feature space regularization.

Weight space regularization: Let's first explain the physical meaning of \mathbf{P} in Eq.(2). Actually, it is easy to prove that:

$$\mathbf{P}^2 = \mathbf{P} = \mathbf{P}^T, \quad \mathbf{e}^T \mathbf{P} \nabla_{\mathbf{w}} \mathcal{L} = 0. \quad (3)$$

The above equations show that \mathbf{P} is the projection matrix for the hyperplane with normal vector \mathbf{e} in weight space, and $\mathbf{P} \nabla_{\mathbf{w}} \mathcal{L}$ is the projected gradient.

The property of projected gradient has been investigated in some previous works [12,28,7,49], which indicate that projecting the gradient of weight will constrict the weight space in a hyperplane or a Riemannian manifold. Similarly, the role of GC can also be viewed from the perspective of projected gradient descent. We give a geometric illustration of SGD with GC in Fig. 2. As shown in Fig. 2, in the t -th step of SGD with GC, the gradient is first projected on the hyperplane determined by $\mathbf{e}^T(\mathbf{w} - \mathbf{w}^t) = 0$, where \mathbf{w}^t is the weight vector in the t -th iteration, and then the weight is updated along the direction of projected gradient $-\mathbf{P} \nabla_{\mathbf{w}^t} \mathcal{L}$. From $\mathbf{e}^T(\mathbf{w} - \mathbf{w}^t) = 0$, we have $\mathbf{e}^T \mathbf{w}^{t+1} = \mathbf{e}^T \mathbf{w}^t =$

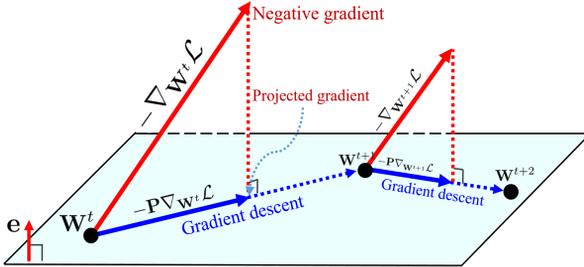


Fig. 2. The geometrical interpretation of GC. The gradient is projected on a hyperplane $\mathbf{e}^T(\mathbf{w} - \mathbf{w}^t) = 0$, where the projected gradient is used to update the weight.

... = $\mathbf{e}^T \mathbf{w}^0$, i.e., $\mathbf{e}^T \mathbf{w}$ is a constant during training. Mathematically, the latent objective function w.r.t. one weight vector \mathbf{w} can be written as follows:

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}), \quad \text{s.t.} \quad \mathbf{e}^T(\mathbf{w} - \mathbf{w}^0) = 0 \quad (4)$$

Clearly, this is a constrained optimization problem on weight vector \mathbf{w} . It regularizes the solution space of \mathbf{w} , reducing the possibility of over-fitting on training data. As a result, GC can improve the generalization capability of trained DNN models, especially when the number of training samples is limited.

It is noted that WS [39] uses a constraint $\mathbf{e}^T \mathbf{w} = 0$ for weight optimization. It reparameterizes weights to meet this constraint. However, this constraint largely limits its practical applications because the initialized weight may not satisfy this constraint. For example, a pretrained DNN on ImageNet usually cannot meet $\mathbf{e}^T \mathbf{w}^0 = 0$ for its initialized weight vectors. If we use WS to fine-tune this DNN, the advantages of pretrained models will disappear. Therefore, we have to retrain the DNN on ImageNet with WS before we fine-tune it. This is very cumbersome. Fortunately the weight constraint of GC in Eq. (4) fits any initialization of weight, e.g., ImageNet pretrained initialization, because it involves the initialized weight \mathbf{w}^0 into the constraint so that $\mathbf{e}^T(\mathbf{w}^0 - \mathbf{w}^0) = 0$ is always true. This greatly extends the applications of GC.

Output feature space regularization: For SGD based algorithms, we have $\mathbf{w}^{t+1} = \mathbf{w}^t - \alpha^t \mathbf{P} \nabla_{\mathbf{w}^t} \mathcal{L}$. It can be derived that $\mathbf{w}^t = \mathbf{w}^0 - \mathbf{P} \sum_{i=0}^{t-1} \alpha^{(i)} \nabla_{\mathbf{w}^{(i)}} \mathcal{L}$. For any input feature vector \mathbf{x} , we have the following theorem:

Theorem 4.1: *Suppose that SGD (or SGDM) with GC is used to update the weight vector \mathbf{w} , for any input feature vectors \mathbf{x} and $\mathbf{x} + \gamma \mathbf{1}$, we have*

$$(\mathbf{w}^t)^T \mathbf{x} - (\mathbf{w}^t)^T (\mathbf{x} + \gamma \mathbf{1}) = \gamma \mathbf{1}^T \mathbf{w}^0 \quad (5)$$

where \mathbf{w}^0 is the initial weight vector and γ is a scalar.

Please find the proof in the **Appendix**. Theorem 4.1 indicates that a constant intensity change (i.e., $\gamma \mathbf{1}$) of an input feature causes a change of output activation; interestingly, this change is only related to γ and $\mathbf{1}^T \mathbf{w}^0$ but not the current weight vector \mathbf{w}^t . $\mathbf{1}^T \mathbf{w}^0$ is the scaled mean of the initial weight vector \mathbf{w}^0 . In particular, if the mean of \mathbf{w}^0 is close to zero, then the output activation is not sensitive to the intensity change of input features, and the output feature space becomes more robust to training sample variations.

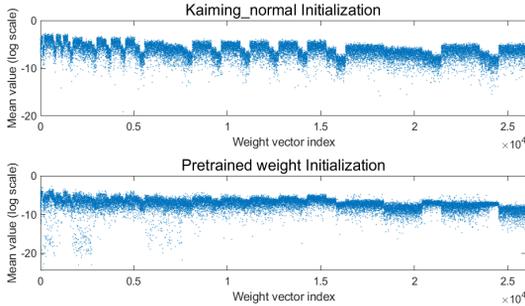


Fig. 3. The absolute value (log scale) of the mean of weight vectors for convolution layers in ResNet50. The x -axis is the weight vector index. We plot the mean value of different convolution layers from left to right with the order from shallow to deep layers. Kaiming normal initialization [14] (top) and ImageNet pre-trained weight initialization (bottom) are employed here. We can see that the mean values are usually very small (less than e^{-7}) for most of the weight vectors.

Indeed, the mean of \mathbf{w}^0 is very close to zero by the commonly used weight initialization strategies, such as Xavier initialization [11], Kaiming initialization [14] and even ImageNet pre-trained weight initialization. Fig. 3 shows the absolute value (log scale) of the mean of weight vectors for Conv layers in ResNet50 with Kaiming normal initialization and ImageNet pre-trained weight initialization. We can see that the mean values of most weight vectors are very small and close to zero (less than e^{-7}). This ensures that if we train the DNN model with GC, the output features will not be sensitive to the variation of the intensity of input features. This property regularizes the output feature space and boosts the generalization performance of DNN training.

4.2 Accelerating Training Process

Optimization landscape smoothing: It has been shown in [45,39] that both BN and WS smooth the optimization landscape. Although BN and WS operate on activations and weights, they implicitly constrict the gradient of weights, making the gradient of weight more predictive and stable for fast training. Specifically, BN and WS use the gradient magnitude $\|\nabla f(\mathbf{x})\|_2$ to capture the Lipschitzness of function $f(\mathbf{x})$. For the loss and its gradients, $f(\mathbf{x})$ will be \mathcal{L} and $\nabla_{\mathbf{w}}\mathcal{L}$, respectively, and \mathbf{x} will be \mathbf{w} . The upper bounds of $\|\nabla_{\mathbf{w}}\mathcal{L}\|_2$ and $\|\nabla_{\mathbf{w}}^2\mathcal{L}\|_2$ ($\nabla_{\mathbf{w}}^2\mathcal{L}$ is the Hessian matrix of \mathbf{w}) have been given in [45,39] to illustrate the optimization landscape smoothing property of BN and WS. Similar conclusion can be made for our proposed GC by comparing the Lipschitzness of original loss function $\mathcal{L}(\mathbf{w})$ with the constrained loss function in Eq. (4) and the Lipschitzness of their gradients. We have the following theorem:

Theorem 4.2: *Suppose $\nabla_{\mathbf{w}}\mathcal{L}$ is the gradient of loss function \mathcal{L} w.r.t. weight vector \mathbf{w} . With the $\Phi_{GC}(\nabla_{\mathbf{w}}\mathcal{L})$ defined in Eq.(2), we have the following conclusion for the loss function and its gradient, respectively:*

$$\begin{cases} \|\Phi_{GC}(\nabla_{\mathbf{w}}\mathcal{L})\|_2 \leq \|\nabla_{\mathbf{w}}\mathcal{L}\|_2, \\ \|\nabla_{\mathbf{w}}\Phi_{GC}(\nabla_{\mathbf{w}}\mathcal{L})\|_2 \leq \|\nabla_{\mathbf{w}}^2\mathcal{L}\|_2. \end{cases} \quad (6)$$

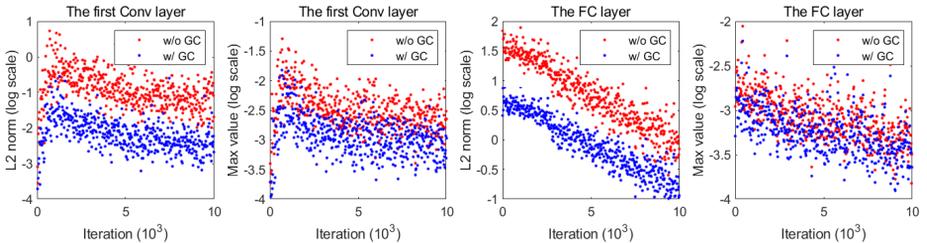


Fig. 4. The L_2 norm (log scale) and max value (log scale) of gradient matrix or tensor vs. iterations. ResNet50 trained on CIFAR100 is used as the DNN model here. The left two sub-figures show the results on the first Conv layer and the right two show the FC layer. The red points represent the results of training without GC and the blue points represent the results with GC. We can see that GC largely reduces the L_2 norm and max value of gradient.

The proof of Theorem 4.2 can be found in the **Appendix**. Theorem 4.2 shows that for the loss function \mathcal{L} and its gradient $\nabla_{\mathbf{w}}\mathcal{L}$, the constrained loss function in Eq. (4) by GC leads to a better Lipschitzness than the original loss function so that the optimization landscape becomes smoother. This means that GC has similar advantages to BN and WS on accelerating training. A good Lipschitzness on gradient implies that the gradients used in training are more predictive and well-behaved so that the optimization landscape can be smoother for faster and more effective training.

Gradient explosion suppression: Another benefit of GC for DNN training is that GC can avoid gradient explosion and make training more stable. This property is similar to gradient clipping [36,37,23,1]. Too large gradients will make the weights change abruptly during training so that the loss may severely oscillate and hard to converge. It has been shown that gradient clipping can suppress large gradient so that the training can be more stable and faster [36,37]. There are two popular gradient clipping approaches: element-wise value clipping [36,23] and norm clipping [37,1], which apply thresholding to element-wise value and gradient norm to gradient matrix, respectively. In order to investigate the influence of GC on clipping gradient, in Fig. 4 we plot the max value and L_2 norm of gradient matrix of the first convolutional layer and the fully-connected layer in ResNet50 (trained on CIFAR100) with and without GC. It can be seen that both the max value and the L_2 norm of the gradient matrix become smaller by using GC in training. This is in accordance to our conclusion in Theorem 4.2 that GC can make training process smoother and faster.

5 Experimental Results

5.1 Setup of Experiments

Extensive experiments are performed to validate the effectiveness of GC. To make the results as comprehensive and clear as possible, we arrange the experiments as follows:

- We start from experiments on the Mini-ImageNet dataset [48] to demonstrate that GC can accelerate the DNN training process and improve the model

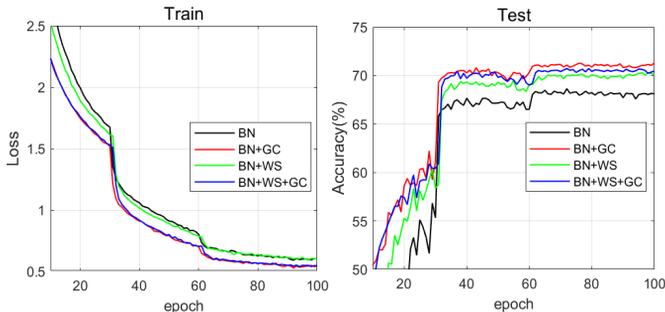


Fig. 5. Training loss (left) and testing accuracy (right) curves vs. training epoch on the Mini-ImageNet. The ResNet50 is used as the DNN model. The compared optimization techniques include BN, BN+GC, BN+WS and BN+WS+GC.

generalization performance. We also evaluate the combinations of GC with BN and WS to show that GC can improve them for DNN optimization.

- We then use the CIFAR100 dataset [26] to evaluate GC with various DNN optimizers (e.g., SGDM, Adam, Adagrad), various DNN architectures (e.g., ResNet, DenseNet, VGG), and different hyper-parameters.
- We then perform experiments on ImageNet [43] to demonstrate that GC also works well on large scale image classification, and show that GC can also work well with normalization methods other than BN, such as GN.
- We consequently perform experiments on four fine-grained image classification datasets (FGVC Aircraft [34], Stanford Cars [25], Stanford Dogs [22] and CUB-200-2011 [50]) to show that GC can be directly adopted to fine-tune the pre-trained DNN models and improve them.
- At last, we perform experiments on the COCO dataset [31] to show that GC also works well on other tasks such as objection detection and segmentation.

GC can be applied to either Conv layer or FC layer, or both of them. In all of our following experiments, if not specified, we always apply GC to both Conv and FC layers. Except for Section 5.3 where we embed GC into different DNN optimizers for test, in all other sections we embed GC into SGDM for experiments, and the momentum is set to 0.9. All experiments are conducted under the Pytorch 1.3 framework and the GPUs are NVIDIA Tesla P100.

We would like to stress that no additional hyper-parameter is introduced in our GC method. Only one line of code is needed to embed GC into the existing optimizers, while keeping all other settings remain unchanged. We compare the performances of DNN models trained with and without GC to validate the effectiveness of GC.

5.2 Results on Mini-Imagenet

Mini-ImageNet [48] is a subset of the ImageNet dataset [43], which was originally proposed for few shot learning. We use the train/test splits provided by [41,20]. It consists of 100 classes and each class has 500 images for training and 100 images for testing. The image resolution is 84×84 . We resize the images into

Table 1. Testing accuracies of different DNN models on CIFAR100

Model	R18	R101	X29	V11	D121
w/o GC	76.87±0.26	78.82± 0.42	79.70±0.30	70.94± 0.34	79.31±0.33
w/ GC	78.82±0.31	80.21±0.31	80.53±0.33	71.69±0.37	79.68±0.40

Table 2. Testing accuracies of different optimizers on CIFAR100

Algorithm	SGDM	Adam	Adagrad	SGDW	AdamW
w/o GC	78.23±0.42	71.64±0.56	70.34 ±0.31	74.02±0.27	74.12±0.42
w/ GC	79.14±0.33	72.80±0.62	71.58±0.37	76.82±0.29	75.07±0.37

224 × 224, which is the standard ImageNet training input size. The DNN we used here is ResNet50, which is trained on 4 GPUs with batch size 128. Other settings are the same as training ImageNet. We repeat the experiments for 10 times and report the average results over the 10 runs.

BN, WS and GC operate on activations, weights and gradients, respectively, and they can be used together to train DNNs. Actually, it is necessary to normalize the feature space by methods such as BN; otherwise, the model is hard to be well trained. Therefore, we evaluate four combinations here: BN, BN+GC, BN+WS and BN+WS+GC. The optimizer is SGDM with momentum 0.9. Fig. 5 presents the training loss and testing accuracy curves of these four combinations. Compared with BN, the training loss of BN+GC decreases much faster and the testing accuracy increases more rapidly. For both BN and BN+WS, GC can further speed up their training speed. Moreover, we can see that BN+GC achieves the highest testing accuracy, validating that GC can accelerate training and enhance the generalization performance simultaneously.

5.3 Experiments on CIFAR100

CIFAR100 [26] consists of 50K training images and 10K testing images from 100 classes. The size of input image is 32 × 32. Since the image resolution is small, we found that applying GC to the Conv layer is good enough on this dataset. All DNN models are trained for 200 epochs using one GPU with batch size 128. The experiments are repeated for 10 times and the results are reported in mean ± std format.

Different networks: We testify GC on different DNN architectures, including ResNet18 (R18), ResNet101 (R101) [15], ResNeXt29 4x64d (X29) [52], VGG11 (V11) [46] and DenseNet121 (D121) [16]. SGDM is used as the network optimizer. The weight decay is set to 0.0005. The initial learning rate is 0.1 and it is multiplied by 0.1 for every 60 epochs. Table 1 shows the testing accuracies of these DNNs. It can be seen that the performance of all DNNs is improved by GC, which verifies that GC is a general optimization technique for different DNN architectures.

Different optimizers: We embed GC into different DNN optimizers, including SGDM [38], Adagrad [10], Adam [24], SGDW and AdamW [32], to test their performance. SGDW and AdamW optimizers directly apply weight decay on weight, instead of using L_2 weight decay regularization. Weight decay is set to 0.001 for SGDW and AdamW, and 0.0005 for other optimizers. The initial learning rate is set to 0.1, 0.01 and 0.001 for SGDM/SGDW, Adagrad,

Table 3. Testing accuracies of different weight decay on CIFAR100 with ResNet50.

Weight decay	0	$1e^{-4}$	$2e^{-4}$	$5e^{-4}$	$1e^{-3}$
w/o GC	71.62±0.31	73.91±0.35	75.57±0.33	78.23±0.42	77.43±0.30
w/ GC	72.83±0.29	76.56±0.31	77.62±0.37	79.14±0.33	78.10±0.36

Table 4. Testing accuracies of different learning rates on CIFAR100 with ResNet50 for SGDM and Adam.

Algorithm	SGDM	SGDM	SGDM	Adam	Adam	Adam
Learning rate	0.05	0.1	0.2	0.0005	0.001	0.0015
w/o GC	76.81±0.27	78.23±0.42	76.53±0.32	73.88±0.46	71.64±0.56	70.63±0.44
w/ GC	78.12±0.33	79.14±0.33	77.71±0.35	74.32±0.55	72.80±0.62	71.22±0.49

Adam/AdamW, respectively, and the learning rate is multiplied by 0.1 for every 60 epochs. The other hyper-parameters are set by their default settings on Pytorch. The DNN model used here is ResNet50. Table 2 shows the testing accuracies. It can be seen that GC boosts the generalization performance of all the five optimizers. It is also found that adaptive learning rate based algorithms Adagrad and Adam have poor generalization performance on CIFAR100, while GC can improve their performance by $> 0.9\%$.

Different hyper-parameter settings: In order to illustrate that GC can achieve consistent improvement with different hyper-parameters, we present the results of GC with different settings of weight decay and learning rates on the CIFAR100 dataset. ResNet50 is used as the backbone. Table 3 shows the testing accuracies with different settings of weight decay, including 0, $1e^{-4}$, $2e^{-4}$, $5e^{-4}$ and $1e^{-3}$. The optimizer is SGDM with learning rate 0.1. It can be seen that the performance of weight decay is consistently improved by GC. Table 4 shows the testing accuracies with different learning rates for SGDM and Adam. For SGDM, the learning rates are 0.05, 0.1 and 0.2, and for Adam, the learning rates are 0.0005, 0.001 and 0.0015. The weight decay is set to $5e^{-4}$. Other settings are the same as those in the manuscript. We can see that GC consistently improves the performance.

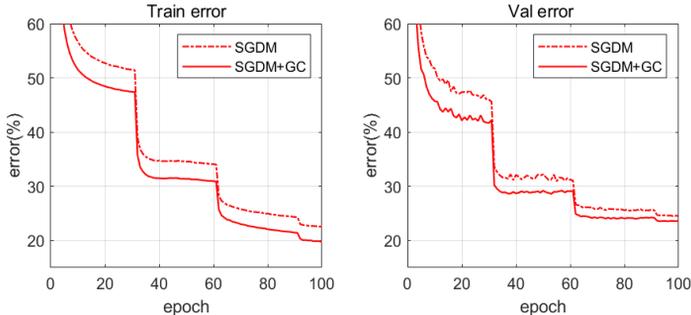
5.4 Results on ImageNet

We then evaluate GC on the large-scale ImageNet dataset [43] which consists of 1.28 million images for training and 50K images for validation from 1000 categories. We use the common training settings and embed GC to SGDM on Conv layer. The ResNet50 and ResNet101 are used as the backbone networks. For the former, we use 4 GPUs with batch size 64 per GPU, and for the latter, we use 8 GPUs with batch size 32 per GPU.

We evaluate four models here: ResNet50 with BN (R50BN), ResNet50 with GN (R50GN), ResNet101 with BN (R101BN) and ResNet101 with GN (R101GN). Table 5 shows the final Top-1 errors of these four DNN models trained with and without GC. We can see that GC can improve the performance by 0.5% \sim 1.2% on ImageNet. Fig. 6 plots the training and validation error curves of ResNet50 (GN is used for feature normalization). We can see that GC can largely speed up the training with GN.

Table 5. Top-1 error rates on ImageNet w/o GC and w/ GC.

Datasets	R50BN	R50GN	R101BN	R101GN
w/o GC	23.71	24.50	22.37	23.34
w/ GC	23.21	23.53	21.82	22.14

**Fig. 6.** Training error (left) and validation error (right) curves vs. training epoch on ImageNet. The DNN model is ResNet50 with GN.

5.5 Results on Fine-grained Image Classification

In order to show that GC can also work with the pre-trained models, we conduct experiments on four challenging fine-grained image classification datasets, including FGVC Aircraft [34], Stanford Cars [25], Stanford Dogs [22] and CUB-200-2011 [50]. The detailed statistics of these four datasets are summarized in Table 6. We use the official pre-trained ResNet50 provided by Pytorch as the baseline DNN for all these four datasets. The original images are resized into 512×512 and we crop the center region with 448×448 as input for both training and testing. The models are pre-trained on ImageNet. We use SGDM with momentum of 0.9 to fine-tune ResNet50 for 100 epochs on 4 GPUs with batch size 256. The initial learning rate is 0.1 for the last FC layer and 0.01 for all pre-trained Conv layers. The learning rate is multiplied by 0.1 at the 50th and 80th epochs. Please note that our goal is to validate the effectiveness of GC but not to push state-of-the-art results, so we only use simple training tricks. We repeat the experiments for 10 times and report the result in mean \pm std format.

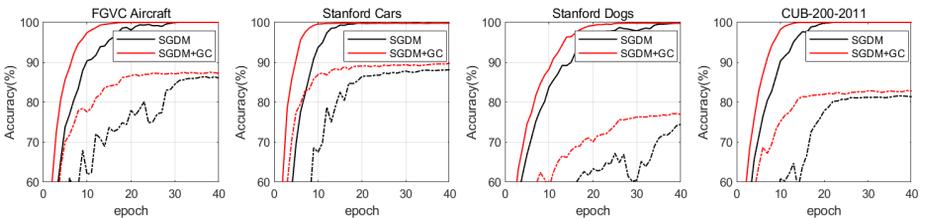
**Fig. 7.** Training accuracy (solid line) and testing accuracy (dotted line) curves vs. training epoch on four fine-grained image classification datasets.

Fig. 7 shows the training and testing accuracies of SGDM and SGDM+GC for the first 40 epochs on the four fine-grained image classification datasets. Table 7 shows the final testing accuracies. We can see that both the training and testing

Table 6. The statistics of fine-grained datasets used in this paper.

Datasets	#Category	#Training	#Testing
FGVC Aircraft	100	6,667	3,333
Stanford Cars	196	8,144	8,041
Stanford Dogs	120	12,000	8,580
CUB-200-2011	200	5,994	5,794

Table 7. Testing accuracies on the four fine-grained image classification datasets.

Datasets	FGVC Aircraft	Stanford Cars	Stanford Dogs	CUB-200-2011
w/o GC	86.62±0.31	88.66±0.22	76.16±0.25	82.07±0.26
w/ GC	87.77±0.27	90.03±0.26	78.23±0.24	83.40±0.30

accuracies of SGDM are improved by GC. For the final classification accuracy, GC improves SGDM by 1.1% ~ 2.1% on these four datasets. This sufficiently demonstrates the effectiveness of GC on fine-tuning pre-trained models.

5.6 Objection Detection and Segmentation

Finally, we evaluate GC on object detection and segmentation tasks to show that GC can also be applied to more tasks beyond image classification. The models are pre-trained on ImageNet. The training batch size for object detection and segmentation is usually very small (e.g., 1 or 2) because of the high resolution of input image. Therefore, the BN layer is usually frozen [15] and the benefits from BN cannot be enjoyed during training. One alternative is to use GN instead. The models are trained on COCO *train2017* dataset (118K images) and evaluated on COCO *val2017* dataset (40K images) [31]. COCO dataset can be used for multiple tasks, including image classification, object detection, semantic segmentation and instance segmentation.

We use the MMDetection [6] toolbox, which contains comprehensive models on object detection and segmentation tasks, as the detection framework. The official implementations and settings are used for all experiments. All the pre-trained models are provided from their official websites, and we fine-tune them on COCO *train2017* set with 8 GPUs and 2 images per GPU. The optimizers are SGDM and SGDM+GC. The backbone networks include ResNet50 (R50), ResNet101 (R101), ResNeXt101-32x4d (X101-32x4d), ResNeXt101-64x4d (X101-32x4d). The Feature Pyramid Network (FPN) [30] is also used. The learning rate schedule is 1X for both Faster R-CNN [42] and Mask R-CNN [13], except R50 with GN and R101 with GN, which use 2X learning rate schedule.

Tabel 8 shows the Average Precision (AP) results of Faster R-CNN. We can see that all the backbone networks trained with GC can achieve a performance gain about 0.3% ~ 0.6% on object detection. Tabel 9 presents the Average Precision for bounding box (AP^b) and instance segmentation (AP^m). It can be seen that the AP^b increases by 0.5% ~ 0.9% for object detection task and the AP^m increases by 0.3% ~ 0.7% for instance segmentation task. Moreover, we find that if 4conv1fc bounding box head, like R50 (4c1f), is used, the performance can increase more by GC. And GC can also boost the performance of GN (see R101GN) and improve the performance of WS (see R50GN+WS). Overall, we can see that GC boosts the generalization performance of all evaluated models.

Table 8. Detection results on COCO by using Faster-RCNN and FPN with various backbone models.

Method	Backbone	AP	AP _{.5}	AP _{.75}	Backbone	AP	AP _{.5}	AP _{.75}
w/o GC	R50	36.4	58.4	39.1	X101-32x4d	40.1	62.0	43.8
w/ GC	R50	37.0	59.0	40.2	X101-32x4d	40.7	62.7	43.9
w/o GC	R101	38.5	60.3	41.6	X101-64x4d	41.3	63.3	45.2
w/ GC	R101	38.9	60.8	42.2	X101-64x4d	41.6	63.8	45.4

Table 9. Detection and segmentation results on COCO by using Mask-RCNN and FPN with various backbone models.

Method	Backbone	AP ^b	AP ^b _{.5}	AP ^b _{.75}	AP ^m	AP ^m _{.5}	AP ^m _{.75}	Backbone	AP ^b	AP ^b _{.5}	AP ^b _{.75}	AP ^m	AP ^m _{.5}	AP ^m _{.75}
w/o GC	R50	37.4	59.0	40.6	34.1	55.5	36.1	R50 (4cif)	37.5	58.2	41.0	33.9	55.0	36.1
w/ GC	R50	37.9	59.6	41.2	34.7	56.1	37.0	R50 (4cif)	38.4	59.5	41.8	34.6	55.9	36.7
w/o GC	R101	39.4	60.9	43.3	35.9	57.7	38.4	R101GN	41.1	61.7	44.9	36.9	58.7	39.3
w/ GC	R101	40.0	61.5	43.7	36.2	58.1	38.7	R101GN	41.7	62.3	45.3	37.4	59.3	40.3
w/o GC	X101-32x4d	41.1	62.8	45.0	37.1	59.4	39.8	R50GN+WS	40.0	60.7	43.6	36.1	57.8	38.6
w/ GC	X101-32x4d	41.6	63.1	45.5	37.4	59.8	39.9	R50GN+WS	40.6	61.3	43.9	36.6	58.2	39.1
w/o GC	X101-64x4d	42.1	63.8	46.3	38.0	60.6	40.9							
w/ GC	X101-64x4d	42.8	64.5	46.8	38.4	61.0	41.1							

This demonstrates that it is a simple yet effective optimization technique, which is general to many tasks beyond image classification.

6 Conclusions

How to efficiently and effectively optimize a DNN is one of the key issues in deep learning research. Previous methods such as batch normalization (BN) and weight standardization (WS) mostly operate on network activations or weights to improve DNN training. We proposed a different approach which operates directly on gradients. Specifically, we removed the mean from the gradient vectors and centralized them to have zero mean. The so-called Gradient Centralization (GC) method demonstrated several desired properties of deep network optimization. We showed that GC actually improves the loss function with a constraint on weight vectors, which regularizes both weight space and output feature space. We also showed that this constrained loss function has better Lipschitzness than the original one so that it has a smoother optimization landscape. Comprehensive experiments were performed and the results demonstrated that GC can be well applied to different tasks with different optimizers and network architectures, improving their training efficiency and generalization performance.

Appendix

A1. Proof of Theorem 4.1

Proof. First we show below a simple property of \mathbf{P} :

$$\mathbf{1}^T \mathbf{P} = \mathbf{1}^T (\mathbf{I} - \mathbf{e} \mathbf{e}^T) = \mathbf{1}^T - \frac{1}{M} \mathbf{1}^T \mathbf{1} \mathbf{1}^T = \mathbf{0}^T,$$

where M is the dimension of \mathbf{e} .

For each SGD step with GC, we have:

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \alpha^t \mathbf{P} \nabla_{\mathbf{w}^t} \mathcal{L}.$$

It can be easily derived that:

$$\mathbf{w}^t = \mathbf{w}^0 - \mathbf{P} \sum_{i=0}^{t-1} \alpha^{(i)} \nabla_{\mathbf{w}^{(i)}} \mathcal{L},$$

where t is the number of iterations. Then for the output activations of \mathbf{x} and $\mathbf{x} + \gamma \mathbf{1}$, there is

$$\begin{aligned} (\mathbf{w}^t)^T \mathbf{x} - (\mathbf{w}^t)^T (\mathbf{x} + \gamma \mathbf{1}) &= \gamma \mathbf{1}^T \mathbf{w}^t \\ &= \gamma \mathbf{1}^T (\mathbf{w}^0 - \mathbf{P} \sum_{i=0}^{t-1} \alpha^{(i)} \nabla_{\mathbf{w}^{(i)}} \mathcal{L}) \\ &= \gamma \mathbf{1}^T \mathbf{w}^0 - \gamma \mathbf{1}^T \mathbf{P} \sum_{i=0}^{t-1} \alpha^{(i)} \nabla_{\mathbf{w}^{(i)}} \mathcal{L} \\ &= \gamma \mathbf{1}^T \mathbf{w}^0. \end{aligned} \tag{7}$$

Therefore,

$$(\mathbf{w}^t)^T \mathbf{x} - (\mathbf{w}^t)^T (\mathbf{x} + \gamma \mathbf{1}) = \gamma \mathbf{1}^T \mathbf{w}^0. \tag{8}$$

For SGD with momentum, the conclusion is the same, because we can obtain a term $\gamma \mathbf{1}^T \mathbf{P} \sum_{i=0}^{t-1} \alpha^{(i)} \mathbf{m}^i$ in the third row of Eq.(7), where \mathbf{m}^i is the momentum in the i th iteration, and this term is also equal to zero.

The proof is completed. ■

A2. Proof of Theorem 4.2

Proof. Because \mathbf{e} is a unit vector, there is $\mathbf{e}^T \mathbf{e} = 1$. We can easily prove that:

$$\begin{aligned} \mathbf{P}^T \mathbf{P} &= (\mathbf{I} - \mathbf{e} \mathbf{e}^T)^T (\mathbf{I} - \mathbf{e} \mathbf{e}^T) \\ &= \mathbf{I} - 2\mathbf{e} \mathbf{e}^T + \mathbf{e} \mathbf{e}^T \mathbf{e} \mathbf{e}^T \\ &= \mathbf{I} - \mathbf{e} \mathbf{e}^T \\ &= \mathbf{P}. \end{aligned} \tag{9}$$

Then for $\Phi_{GC}(\nabla_{\mathbf{w}} \mathcal{L})$, we have:

$$\begin{aligned} \|\Phi_{GC}(\nabla_{\mathbf{w}} \mathcal{L})\|_2^2 &= \Phi_{GC}(\nabla_{\mathbf{w}} \mathcal{L})^T \Phi_{GC}(\nabla_{\mathbf{w}} \mathcal{L}) \\ &= (\mathbf{P} \nabla_{\mathbf{w}} \mathcal{L})^T (\mathbf{P} \nabla_{\mathbf{w}} \mathcal{L}) \\ &= \nabla_{\mathbf{w}} \mathcal{L}^T \mathbf{P}^T \mathbf{P} \nabla_{\mathbf{w}} \mathcal{L} \\ &= \nabla_{\mathbf{w}} \mathcal{L}^T \mathbf{P} \nabla_{\mathbf{w}} \mathcal{L} \\ &= \nabla_{\mathbf{w}} \mathcal{L}^T (\mathbf{I} - \mathbf{e} \mathbf{e}^T) \nabla_{\mathbf{w}} \mathcal{L} \\ &= \nabla_{\mathbf{w}} \mathcal{L}^T \nabla_{\mathbf{w}} \mathcal{L} - \nabla_{\mathbf{w}} \mathcal{L}^T \mathbf{e} \mathbf{e}^T \nabla_{\mathbf{w}} \mathcal{L} \\ &= \|\nabla_{\mathbf{w}} \mathcal{L}\|_2^2 - \|\mathbf{e}^T \nabla_{\mathbf{w}} \mathcal{L}\|_2^2 \\ &\leq \|\nabla_{\mathbf{w}} \mathcal{L}\|_2^2. \end{aligned} \tag{10}$$

For $\nabla_{\mathbf{w}}\Phi_{GC}(\nabla_{\mathbf{w}}\mathcal{L})$, we also have

$$\begin{aligned}
 \|\nabla\Phi_{GC}(\nabla_{\mathbf{w}}\mathcal{L})\|_2^2 &= \|\mathbf{P}\nabla_{\mathbf{w}}^2\mathcal{L}\|_2^2 \\
 &= \nabla_{\mathbf{w}}^2\mathcal{L}^T\mathbf{P}^T\mathbf{P}\nabla_{\mathbf{w}}^2\mathcal{L} \\
 &= \nabla_{\mathbf{w}}^2\mathcal{L}^T\mathbf{P}\nabla_{\mathbf{w}}^2\mathcal{L} \\
 &= \|\nabla_{\mathbf{w}}^2\mathcal{L}\|_2^2 - \|\mathbf{e}^T\nabla_{\mathbf{w}}^2\mathcal{L}\|_2^2 \\
 &\leq \|\nabla_{\mathbf{w}}^2\mathcal{L}\|_2^2.
 \end{aligned} \tag{11}$$

The proof is completed. ■

References

1. Abadi, M., Chu, A., Goodfellow, I., McMahan, H.B., Mironov, I., Talwar, K., Zhang, L.: Deep learning with differential privacy. In: Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security. pp. 308–318. ACM (2016)
2. Abu-El-Haija, S., Kothari, N., Lee, J., Natsev, P., Toderici, G., Varadarajan, B., Vijayanarasimhan, S.: Youtube-8m: A large-scale video classification benchmark. arXiv preprint arXiv:1609.08675 (2016)
3. Bjorck, J., Gomes, C., Selman, B., Weinberger, K.Q.: Understanding batch normalization pp. 7694–7705 (2018)
4. Bottou, L.: Stochastic gradient learning in neural networks. Proceedings of NeuroNimes **91**(8), 12 (1991)
5. Bottou, L.: Large-scale machine learning with stochastic gradient descent. In: Proceedings of COMPSTAT'2010, pp. 177–186. Springer (2010)
6. Chen, K., Wang, J., Pang, J., Cao, Y., Xiong, Y., Li, X., Sun, S., Feng, W., Liu, Z., Xu, J., et al.: Mmdetection: Open mmlab detection toolbox and benchmark. arXiv preprint arXiv:1906.07155 (2019)
7. Cho, M., Lee, J.: Riemannian approach to batch normalization. In: Advances in Neural Information Processing Systems. pp. 5225–5235 (2017)
8. Courbariaux, M., Bengio, Y., David, J.P.: Binaryconnect: Training deep neural networks with binary weights during propagations. In: Advances in neural information processing systems. pp. 3123–3131 (2015)
9. Courbariaux, M., Hubara, I., Soudry, D., El-Yaniv, R., Bengio, Y.: Binarized neural networks: Training deep neural networks with weights and activations constrained to+ 1 or-1. arXiv preprint arXiv:1602.02830 (2016)
10. Duchi, J., Hazan, E., Singer, Y.: Adaptive subgradient methods for online learning and stochastic optimization. Journal of Machine Learning Research **12**(Jul), 2121–2159 (2011)
11. Glorot, X., Bengio, Y.: Understanding the difficulty of training deep feedforward neural networks. In: Proceedings of the thirteenth international conference on artificial intelligence and statistics. pp. 249–256 (2010)
12. Gupta, H., Jin, K.H., Nguyen, H.Q., McCann, M.T., Unser, M.: Cnn-based projected gradient descent for consistent ct image reconstruction. IEEE transactions on medical imaging **37**(6), 1440–1453 (2018)
13. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 2961–2969 (2017)
14. He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: Proceedings of the IEEE international conference on computer vision. pp. 1026–1034 (2015)
15. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
16. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4700–4708 (2017)
17. Huang, L., Liu, X., Liu, Y., Lang, B., Tao, D.: Centered weight normalization in accelerating training of deep neural networks. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 2803–2811 (2017)

18. Huang, X., Belongie, S.: Arbitrary style transfer in real-time with adaptive instance normalization. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1501–1510 (2017)
19. Ioffe, S., Szegedy, C.: Batch normalization: Accelerating deep network training by reducing internal covariate shift. arXiv preprint arXiv:1502.03167 (2015)
20. Iscen, A., Tolias, G., Avrithis, Y., Chum, O.: Label propagation for deep semi-supervised learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5070–5079 (2019)
21. Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R., Fei-Fei, L.: Large-scale video classification with convolutional neural networks. In: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. pp. 1725–1732 (2014)
22. Khosla, A., Jayadevaprakash, N., Yao, B., Li, F.F.: Novel dataset for fgvc: Stanford dogs. In: San Diego: CVPR Workshop on FGVC. vol. 1 (2011)
23. Kim, J., Kwon Lee, J., Mu Lee, K.: Accurate image super-resolution using very deep convolutional networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1646–1654 (2016)
24. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
25. Krause, J., Stark, M., Deng, J., Fei-Fei, L.: 3d object representations for fine-grained categorization. In: Proceedings of the IEEE International Conference on Computer Vision Workshops. pp. 554–561 (2013)
26. Krizhevsky, A., Hinton, G., et al.: Learning multiple layers of features from tiny images. Tech. rep., Citeseer (2009)
27. Krogh, A., Hertz, J.A.: A simple weight decay can improve generalization. In: Advances in neural information processing systems. pp. 950–957 (1992)
28. Larsson, M., Arnab, A., Kahl, F., Zheng, S., Torr, P.: A projected gradient descent method for crf inference allowing end-to-end training of arbitrary pairwise potentials. In: International Workshop on Energy Minimization Methods in Computer Vision and Pattern Recognition. pp. 564–579. Springer (2017)
29. Lei Ba, J., Kiros, J.R., Hinton, G.E.: Layer normalization. arXiv preprint arXiv:1607.06450 (2016)
30. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2117–2125 (2017)
31. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference on computer vision. pp. 740–755. Springer (2014)
32. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101 (2017)
33. Luo, P., Wang, X., Shao, W., Peng, Z.: Towards understanding regularization in batch normalization (2018)
34. Maji, S., Rahtu, E., Kannala, J., Blaschko, M., Vedaldi, A.: Fine-grained visual classification of aircraft. arXiv preprint arXiv:1306.5151 (2013)
35. Nair, V., Hinton, G.E.: Rectified linear units improve restricted boltzmann machines. In: Proceedings of the 27th international conference on machine learning (ICML-10). pp. 807–814 (2010)
36. Pascanu, R., Mikolov, T., Bengio, Y.: Understanding the exploding gradient problem. CoRR, abs/1211.5063 2 (2012)
37. Pascanu, R., Mikolov, T., Bengio, Y.: On the difficulty of training recurrent neural networks. In: International conference on machine learning. pp. 1310–1318 (2013)

38. Qian, N.: On the momentum term in gradient descent learning algorithms. *Neural networks* **12**(1), 145–151 (1999)
39. Qiao, S., Wang, H., Liu, C., Shen, W., Yuille, A.: Weight standardization. *arXiv preprint arXiv:1903.10520* (2019)
40. Rastegari, M., Ordonez, V., Redmon, J., Farhadi, A.: Xnor-net: Imagenet classification using binary convolutional neural networks. In: *European Conference on Computer Vision*. pp. 525–542. Springer (2016)
41. Ravi, S., Larochelle, H.: Optimization as a model for few-shot learning (2016)
42. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: *Advances in neural information processing systems*. pp. 91–99 (2015)
43. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. *International journal of computer vision* **115**(3), 211–252 (2015)
44. Salimans, T., Kingma, D.P.: Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In: *Advances in Neural Information Processing Systems*. pp. 901–909 (2016)
45. Santurkar, S., Tsipras, D., Ilyas, A., Madry, A.: How does batch normalization help optimization? (no, it is not about internal covariate shift) pp. 2483–2493 (2018)
46. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014)
47. Ulyanov, D., Vedaldi, A., Lempitsky, V.: Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022* (2016)
48. Vinyals, O., Blundell, C., Lillicrap, T., Wierstra, D., et al.: Matching networks for one shot learning. In: *Advances in neural information processing systems*. pp. 3630–3638 (2016)
49. Vorontsov, E., Trabelsi, C., Kadoury, S., Pal, C.: On orthogonality and learning recurrent networks with long term dependencies. In: *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. pp. 3570–3578. *JMLR. org* (2017)
50. Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: The caltech-ucsd birds-200-2011 dataset (2011)
51. Wu, Y., He, K.: Group normalization. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. pp. 3–19 (2018)
52. Xie, S., Girshick, R., Dollár, P., Tu, Z., He, K.: Aggregated residual transformations for deep neural networks. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 1492–1500 (2017)
53. Zhang, C., Bengio, S., Hardt, M., Recht, B., Vinyals, O.: Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530* (2016)
54. Zhang, G., Wang, C., Xu, B., Grosse, R.: Three mechanisms of weight decay regularization. *arXiv preprint arXiv:1810.12281* (2018)