

Gradient Computation in Winograd Domain

1 Winograd Convolution

Assume we need to perform convolution on a $m \times m$ input tile \mathbf{I} with a $n \times n$ weight kernel \mathbf{W} . The output tile \mathbf{O} will have a size of $(m - n + 1) \times (m - n + 1)$.

With the Winograd transformation, this convolution operation can be performed by

$$\mathbf{O} = \mathbf{A}^\top [(\mathbf{G}\mathbf{W}\mathbf{G}^\top) \odot (\mathbf{B}^\top \mathbf{I} \mathbf{B})] \mathbf{A} \quad (1)$$

where \mathbf{A} , \mathbf{G} and \mathbf{B} are matrices used for the Winograd transformation. \odot is element-wise matrix multiplication.

As an example, for $m = 4$ and $n = 3$, we have

$$\mathbf{B}^\top = \begin{bmatrix} 1 & 0 & -1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & -1 & 1 & 0 \\ 0 & 1 & 0 & -1 \end{bmatrix} \quad \mathbf{G} = \begin{bmatrix} 1 & 0 & 0 \\ 1/2 & 1/2 & 1/2 \\ 1/2 & -1/2 & 1/2 \\ 0 & 0 & 1 \end{bmatrix} \quad \mathbf{A}^\top = \begin{bmatrix} 1 & 1 & 1 & 0 \\ 0 & 1 & -1 & -1 \end{bmatrix} \quad (2)$$

In this note, we only discuss the situation with a one-channel input tile and a one-channel weight kernel. It can be extended to a regular convolutional layer.

2 Winograd Layer

Deep Neural Networks (DNNs) have much internal redundancy. We can remove weight parameters from DNN models to reduce computation and model sizes. In this case, the weight kernels, e.g., \mathbf{W} , will become sparse.

In the Winograd convolution, the majority of the floating-point multiplication is performed in the element-wise multiplication. Therefore, to fully utilize the redundancy in DNN models, we need higher sparsity in $\mathbf{G}\mathbf{W}\mathbf{G}^\top$. However, the sparsity in \mathbf{W} cannot be directly transformed into the sparsity in $\mathbf{G}\mathbf{W}\mathbf{G}^\top$.

Assume $\mathbf{Q} = \mathbf{G}\mathbf{W}\mathbf{G}^\top$. Here \mathbf{Q} has a size of $m \times m$ and \mathbf{W} has a size of $n \times n$. Since we will have $m > n$, the mapping from \mathbf{W} to \mathbf{Q} is non-invertible.

One solution to increase the sparsity in \mathbf{Q} is to directly consider \mathbf{Q} as the parameters for a convolutional layer. Convolutional layers with weight parameters in the Winograd domain are named as Winograd layers. Then we can remove unimportant parameters directly from \mathbf{Q} .

The computation performed by a Winograd layer is

$$\mathbf{O} = \mathbf{A}^\top [\mathbf{Q} \odot (\mathbf{B}^\top \mathbf{I} \mathbf{B})] \mathbf{A} \quad (3)$$

3 Explosion of Update Steps

With Winograd layers instead of original convolutional layers, we can still use stochastic gradient descent (SGD) algorithm to train the network. However, we found it requires a small learning rate (LR) and, therefore, the training converges extremely slow.

This is due to the explosion of update steps caused by the linear transformation $\mathbf{G}\mathbf{W}\mathbf{G}^\top$. Here we use a simple example to show the explosion problem of update steps.

Assume we have a parameter x and its gradient with respect to the loss (L) is

$$\frac{\partial L}{\partial x} = 1 \quad (4)$$

with $LR = 0.1$, x will be updated by

$$x := x + \Delta x \quad (5)$$

$$\Delta x = -LR \cdot \frac{\partial L}{\partial x} = -0.1 \quad (6)$$

Assume we use a linear transformation

$$y = 0.1x \quad (7)$$

to map x to y . If we are training parameter x , then we need to update y by

$$y := y + \Delta y \quad (8)$$

$$\Delta y = 0.1 \cdot \Delta x = -0.01 \quad (9)$$

However, assuming we are now directly performing training on y , the gradient of y with respect to the loss (L) is

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial x} \cdot \frac{\partial x}{\partial y} = 10 \quad (10)$$

In this case, y will be updated by

$$y := y + \Delta y' \quad (11)$$

$$\Delta y' = -LR \cdot \frac{\partial L}{\partial y} = -1 \quad (12)$$

We will have $\Delta y' = 100 \cdot \Delta y$. It means, if we directly train the parameter y instead of x , the updates on y will be extended by 100x, which is an explosion in update steps.

The linear transformation $\mathbf{Q} = \mathbf{G}\mathbf{W}\mathbf{G}^\top$ will have a similar effect. The direct training of \mathbf{Q} will suffer from the explosion of the update steps, which requires LR to be extremely small. Also, it is difficult to find appropriate LR for each parameter in \mathbf{Q} .

4 Gradients in Winograd Domain

For parameter y , one solution for the explosion of update steps is to calculate Δy by

$$\Delta y = \Delta(0.1x) = 0.1 \cdot \Delta x = 0.1 \cdot (-LR) \cdot \frac{\partial L}{\partial x} \quad (13)$$

In this case, the effective gradient of y , $(\frac{\partial L}{\partial y})^*$, can be calculated by

$$(\frac{\partial L}{\partial y})^* = \frac{\Delta y}{-LR} = \frac{0.1 \cdot \Delta x}{-LR} = \frac{0.1 \cdot (-LR) \cdot \frac{\partial L}{\partial x}}{-LR} = 0.1 \cdot \frac{\partial L}{\partial x} \quad (14)$$

Similarly, for the Winograd layer parameters \mathbf{Q} , we need to have

$$\Delta \mathbf{Q} = \Delta(\mathbf{G}\mathbf{W}\mathbf{G}^\top) = \mathbf{G}(\Delta \mathbf{W})\mathbf{G}^\top \quad (15)$$

Therefore, we can calculate the effective gradients for \mathbf{Q} by

$$(\frac{\partial L}{\partial \mathbf{Q}})^* = \frac{\Delta \mathbf{Q}}{-LR} = \mathbf{G} \frac{\Delta \mathbf{W}}{-LR} \mathbf{G}^\top = \mathbf{G} \frac{\partial L}{\partial \mathbf{W}} \mathbf{G}^\top \quad (16)$$

Since $\mathbf{Q} = \mathbf{G}\mathbf{W}\mathbf{G}^\top$, we will have

$$\frac{\partial L}{\partial \mathbf{W}} = \mathbf{G}^\top \frac{\partial L}{\partial \mathbf{Q}} \mathbf{G} \quad (17)$$

where $\frac{\partial L}{\partial \mathbf{Q}}$ is the original gradients of \mathbf{Q} calculated through back-propagation.

Therefore, based on Equation 16 and 17, we will have

$$(\frac{\partial L}{\partial \mathbf{Q}})^* = \mathbf{G} \frac{\partial L}{\partial \mathbf{W}} \mathbf{G}^\top = \mathbf{G}(\mathbf{G}^\top \frac{\partial L}{\partial \mathbf{Q}} \mathbf{G})\mathbf{G}^\top = (\mathbf{G}\mathbf{G}^\top) \frac{\partial L}{\partial \mathbf{Q}} (\mathbf{G}\mathbf{G}^\top) \quad (18)$$

Then \mathbf{Q} will be updated by

$$\mathbf{Q} := \mathbf{Q} - LR \cdot (\frac{\partial L}{\partial \mathbf{Q}})^* \quad (19)$$

5 Regularization

In this section, we will discuss how to perform regularization on \mathbf{Q} .

The loss function of DNN models usually includes two parts: the divergence between the DNN output and the ground-truth output, L^D , and the regularization term, L^R . The regularization term is used to prevent over-fitting. In this case, we have

$$L = L^D + L^R \quad (20)$$

Assume we are using the L2 regularization, then for the original weight kernel \mathbf{W} ,

$$L^R = \frac{\lambda}{2} \|\mathbf{W}\|_2 \quad (21)$$

where λ is the regularization strength which is a scalar value.

Assume with the Winograd layers, we still use the same loss function as for the original convolutional layers. Then for Equation 18, we have

$$\begin{aligned} \left(\frac{\partial L}{\partial \mathbf{Q}}\right)^* &= (\mathbf{G}\mathbf{G}^\top) \frac{\partial L}{\partial \mathbf{Q}} (\mathbf{G}\mathbf{G}^\top) \\ &= (\mathbf{G}\mathbf{G}^\top) \frac{\partial L^D}{\partial \mathbf{Q}} (\mathbf{G}\mathbf{G}^\top) + (\mathbf{G}\mathbf{G}^\top) \frac{\partial L^R}{\partial \mathbf{Q}} (\mathbf{G}\mathbf{G}^\top) \end{aligned} \quad (22)$$

$\frac{\partial L^D}{\partial \mathbf{Q}}$ can be calculated with the back-propagation. But

$$\frac{\partial L^R}{\partial \mathbf{Q}} = \frac{\partial}{\partial \mathbf{Q}} \left(\frac{\lambda}{2} \|\mathbf{W}\|_2 \right) \quad (23)$$

is intractable since the linear transformation, $\mathbf{Q} = \mathbf{G}\mathbf{W}\mathbf{G}^\top$, is non-invertible.

However, we can simplify the computation since

$$\begin{aligned} (\mathbf{G}\mathbf{G}^\top) \frac{\partial L^R}{\partial \mathbf{Q}} (\mathbf{G}\mathbf{G}^\top) &= \mathbf{G} (\mathbf{G}^\top \frac{\partial L^R}{\partial \mathbf{Q}} \mathbf{G}) \mathbf{G}^\top \\ &= \mathbf{G} \frac{\partial L^R}{\partial \mathbf{W}} \mathbf{G}^\top \\ &= \mathbf{G} (\lambda \cdot \mathbf{W}) \mathbf{G}^\top \\ &= \lambda \cdot (\mathbf{G}\mathbf{W}\mathbf{G}^\top) \\ &= \lambda \cdot \mathbf{Q} \end{aligned} \quad (24)$$

Then, Equation 22 can be simplified to

$$\begin{aligned} \left(\frac{\partial L}{\partial \mathbf{Q}}\right)^* &= (\mathbf{G}\mathbf{G}^\top) \frac{\partial L^D}{\partial \mathbf{Q}} (\mathbf{G}\mathbf{G}^\top) + (\mathbf{G}\mathbf{G}^\top) \frac{\partial L^R}{\partial \mathbf{Q}} (\mathbf{G}\mathbf{G}^\top) \\ &= (\mathbf{G}\mathbf{G}^\top) \frac{\partial L^D}{\partial \mathbf{Q}} (\mathbf{G}\mathbf{G}^\top) + \lambda \cdot \mathbf{Q} \end{aligned} \quad (25)$$

The term, $\lambda \cdot \mathbf{Q}$, means, effectively, we are directly applying an L2 regularization

$$(L^R)^* = \frac{\lambda}{2} \|\mathbf{Q}\|_2 \quad (26)$$

on \mathbf{Q} , the parameters in the Winograd domain.

In conclusion, the update we need to perform in each step of SGD is

$$\mathbf{Q} := \mathbf{Q} - LR \cdot [(\mathbf{G}\mathbf{G}^\top) \frac{\partial L^D}{\partial \mathbf{Q}} (\mathbf{G}\mathbf{G}^\top) + \lambda \cdot \mathbf{Q}] \quad (27)$$

6 Evaluation and Results

We perform the training on two networks on the MNIST dataset: LeNet-5 and LeNet-5-3x3. LeNet-5 includes 2 convolution layers with the kernel size of 5 and 2 fully-connected layers. LeNet-5-3x3 is designed by ourselves. It includes 4 convolution layers with the kernel size of 3. The size of the output tiles in the Winograd convolution is fixed to 4×4 , which means $m - n + 1 = 4$. All networks are trained from scratch. Table 1 shows the results.

Table 1: Training Accuracy

	Accuracy (%)	
	Original	Winograd
LeNet-5	99.44%	99.44%
LeNet-5-3x3	99.26%	99.24%