

Convolution and Transposed Convolution as Matrix Multiplication

Introduction

In deep learning, [convolution](#) and [transposed convolution](#) are often used in the neural networks. Unlike convolution, transposed convolution is sometimes confusing and not too many people know why it's called transposed convolution.

In this blog post, I would like to discuss how to view convolution and transposed convolution as matrix multiplication, and how to understand the name of transposed convolution.

Convolution and Transposed Convolution as Matrix Multiplication

At each step of the convolution operation, we apply the kernel tensor onto a portion of the elements in the input tensor, compute their element-wise products, and further sum the products together. We repeat the step many times until the kernel convolves over the entire tensor. With this in mind, we could actually formulate the convolution operation using matrix operations equivalently.

Concretely, for any convolution operation in deep learning, $Y = K * X$, it could be formulated as matrix multiplication

where X' is the flatten representation of X which might have already been padded and dilated, W is the sparse matrix representation of kernel K , and Y' is the flatten representation of the output Y .

More specifically, when K and X are 2D matrices, $K \in R^{h_K \times w_K}$, $X \in R^{h_X \times w_X}$, $X' \in R^{h_X w_X}$, $Y \in R^{h_Y \times w_Y}$, $Y' \in R^{h_Y w_Y}$, we should have a sparse matrix $W \in R^{h_Y w_Y \times h_X w_X}$ that does the same transformation as the convolution.

Similarly, for any transposed convolution operation in deep learning, $Z = K \star Y$, it could also be formulated as matrix multiplication

where Z' is the flatten representation of the output Z , Y and Y' are just ones we just obtained from the previous convolution operation $Y = K * X$.

Z must have $Z \in R^{h_Z \times w_Z}$ and $Z' \in R^{h_Z w_Z}$, as if the shape have been reverted back before the convolution operation $Y = K * X$. Notice that the kernels used for the convolution and transposed convolution, K , are exactly the same, and the weight matrices used in the convolution and transposed convolution matrix multiplications, W and W^T , are just transposed to each other. Z , however, usually does not equal to X . Because $Z \neq X$, we cannot call transposed convolution as deconvolution.

Implementing Convolution and Transposed Convolution as Matrix Operation

Let's ignore the channel dimension and the bias term for convolution and transposed convolution for now, and implement convolution and transposed convolution as matrix operation. We also assume the stride is 1 for both convolution and transposed convolution.

```

1  import torch
2  from torch import nn
3
4
5  def corr2d(X, K):
6
7      # Convolution in deep learning is a misnomer.
8      # In fact, it is cross-correlation.
9      # https://d2l.ai/chapter\_convolutional-neural-networks/conv-layer.html
10     # This is equivalent as Conv2D that that input_channel == output_channel == 1 and stride == 1.
11
12     assert X.dim() == 2 and K.dim() == 2
13
14     h, w = K.shape
15     Y = torch.zeros((X.shape[0] - h + 1, X.shape[1] - w + 1))
16     for i in range(Y.shape[0]):
17         for j in range(Y.shape[1]):
18             Y[i, j] = (X[i:i + h, j:j + w] * K).sum()
19
20     return Y
21
22
23  def get_sparse_kernel_matrix(K, h_X, w_X):
24
25     # Assuming no channels and stride == 1.
26     # Convert the kernel matrix to sparse matrix (dense matrix with lots of zeros in fact).
27     # This is a little bit brain-twisting.
28
29     h_K, w_K = K.shape
30
31     h_Y, w_Y = h_X - h_K + 1, w_X - w_K + 1
32
33     W = torch.zeros((h_Y * w_Y, h_X * w_X))
34     for i in range(h_Y):
35         for j in range(w_Y):
36             for ii in range(h_K):
37                 for jj in range(w_K):
38                     W[i * w_Y + j, i * w_X + j + ii * w_X + jj] = K[ii, jj]
39
40     return W
41
42
43  def conv2d_as_matrix_mul(X, K):
44
45     # Assuming no channels and stride == 1.
46     # Convert the kernel matrix to sparse matrix (dense matrix with lots of zeros in fact).
47     # This is a little bit brain-twisting.
48
49     h_K, w_K = K.shape

```

```

50     h_X, w_X = X.shape
51
52     h_Y, w_Y = h_X - h_K + 1, w_X - w_K + 1
53
54     W = get_sparse_kernel_matrix(K=K, h_X=h_X, w_X=w_X)
55
56     Y = torch.matmul(W, X.reshape(-1)).reshape(h_Y, w_Y)
57
58     return Y
59
60
61 def conv_transposed_2d_as_matrix_mul(X, K):
62
63     # Assuming no channels and stride == 1.
64     # Convert the kernel matrix to sparse matrix (dense matrix with lots of zeros in fact).
65     # This is a little bit brain-twisting.
66
67     h_K, w_K = K.shape
68     h_X, w_X = X.shape
69
70     h_Y, w_Y = h_X + h_K - 1, w_X + w_K - 1
71
72     # It's like the kernel were applied on the output tensor.
73     W = get_sparse_kernel_matrix(K=K, h_X=h_Y, w_X=w_Y)
74
75     # Weight matrix tranposed.
76     Y = torch.matmul(W.T, X.reshape(-1)).reshape(h_Y, w_Y)
77
78     return Y
79
80
81 def main():
82
83     X = torch.arange(30).reshape(5, 6).float()
84     K = torch.arange(8).reshape(2, 4).float()
85     print("X:")
86     print(X)
87     print("K:")
88     print(K)
89     print("Cross-Correlation:")
90     Y = corr2d(X=X, K=K)
91     print(Y)
92
93     conv = nn.Conv2d(in_channels=1,
94                      out_channels=1,
95                      kernel_size=K.shape,
96                      padding=0,
97                      stride=1,
98                      bias=False)
99     conv.weight.data = K.unsqueeze(0).unsqueeze(0)

```

```

100     Z1 = conv(X.unsqueeze(0).unsqueeze(0)).squeeze(0).squeeze(0).detach()
101     print("Convolution:")
102     print(Z1)
103     assert torch.equal(Y, Z1)
104
105     print("Convolution as Matrix Multiplication:")
106     Z2 = conv2d_as_matrix_mul(X=X, K=K)
107     print(Z2)
108     assert torch.equal(Y, Z2)
109
110     conv_transposed = nn.ConvTranspose2d(in_channels=1,
111                                           out_channels=1,
112                                           kernel_size=K.shape,
113                                           padding=0,
114                                           stride=1,
115                                           bias=False)
116     conv_transposed.weight.data = K.unsqueeze(0).unsqueeze(0)
117     Z3 = conv_transposed(Y.unsqueeze(0).unsqueeze(0)).squeeze(0).squeeze(0).detach()
118     print("Transposed Convolution:")
119     print(Z3)
120     # The shape will "go back".
121     assert Z3.shape == X.shape
122
123     print("Transposed Convolution as Matrix Multiplication:")
124     Z4 = conv_transposed_2d_as_matrix_mul(X=Y, K=K)
125     print(Z4)
126     assert torch.equal(Z3, Z4)
127     assert Z4.shape == X.shape
128
129     return
130
131
132 if __name__ == "__main__":
133
134     main()

```

We could see that outputs from the matrix multiplication implementations match the expectation.

```

1  $ python conv_as_gemm.py
2  X:
3  tensor([[ 0.,  1.,  2.,  3.,  4.,  5.],
4         [ 6.,  7.,  8.,  9., 10., 11.],
5         [12., 13., 14., 15., 16., 17.],
6         [18., 19., 20., 21., 22., 23.],
7         [24., 25., 26., 27., 28., 29.]])
8  K:
9  tensor([[0., 1., 2., 3.],
10         [4., 5., 6., 7.]])
11 Cross-Correlation:

```

```

12  tensor([[184., 212., 240.],
13         [352., 380., 408.],
14         [520., 548., 576.],
15         [688., 716., 744.]])
16  Convolution:
17  tensor([[184., 212., 240.],
18         [352., 380., 408.],
19         [520., 548., 576.],
20         [688., 716., 744.]])
21  Convolution as Matrix Multiplication:
22  tensor([[184., 212., 240.],
23         [352., 380., 408.],
24         [520., 548., 576.],
25         [688., 716., 744.]])
26  Transposed Convolution:
27  tensor([[ 0., 184., 580., 1216., 1116., 720.],
28         [ 736., 2120., 4208., 5984., 4880., 2904.],
29         [1408., 3800., 7232., 10016., 7904., 4584.],
30         [2080., 5480., 10256., 14048., 10928., 6264.],
31         [2752., 6304., 10684., 12832., 9476., 5208.]])
32  Transposed Convolution as Matrix Multiplication:
33  tensor([[ 0., 184., 580., 1216., 1116., 720.],
34         [ 736., 2120., 4208., 5984., 4880., 2904.],
35         [1408., 3800., 7232., 10016., 7904., 4584.],
36         [2080., 5480., 10256., 14048., 10928., 6264.],
37         [2752., 6304., 10684., 12832., 9476., 5208.]])

```

For multi-channel kernel, input tensor, and output tensor, the derivation of the weight matrix for matrix multiplication is more complicated, but it does not change the nature that the convolution and transposed convolution can be viewed as matrix multiplications.

Backward Propagation

Considering the convolution as matrix multiplications.

Given the flattened input vector $x \in R^{m_1}$, weight matrix $W \in R^{m_2 \times m_1}$, and the flattened output vector $y \in R^{m_2}$, the matrix multiplication is

Because the gradients matrix with respect to the input x is

where $\nabla_x y \in R^{m_1 \times m_2}$.

In the backward propagation, the input is $\frac{\partial L}{\partial y} \in R^{m_2}$, and the output of the backward propagation is

$$\begin{aligned}
 \frac{\partial L}{\partial x} &= \nabla_x y \frac{\partial L}{\partial y} \\
 &= W^T \frac{\partial L}{\partial y}
 \end{aligned}$$

Therefore, in the forward propagation and the backward propagation of the convolution operation, the weight matrices are W and W^T , respectively.

Similarly, considering the transposed convolution as matrix multiplications. In the forward propagation and the backward propagation of the transposed convolution operation, the weight matrices are W^T and W , respectively.

What does it imply? If the convolution and transposed convolution were implemented as matrix multiplication. To compute the gradient matrix W_1 for transposed convolution operation $Z = K \star Y$, if we happen to have already computed the gradient matrix W_2 for $X = K \star Z$, we must have $W_1 = W_2^T$. However, I think in practice this property is less useful because usually the kernel used for different layers are different.

Because the weight matrices used in the forward propagation and the backward propagation of a transposed convolution are just the transpose of the weight matrices used in the forward propagation and the backward propagation of a convolution which has the same kernel parameters, that's probably why transposed convolution is called transposed convolution.

Miscellaneous

We could see that without non-linear activation function, a sequence of convolution operations is just a linear function. Therefore, having non-linear activation is also important for convolutional neural networks.

References