# COSE474 2024-2 HW 2

## 2023320119 김동현

```
pip install d2l==1.0.3
```

⮒  **숨겨진 출력 표시**

## 7.2 Convolution for Images

```python
import torch
from torch import nn
from d2l import torch as d2l


def corr2d(X, K):
    h, w = K.shape
    Y = torch.zeros((X.shape[0] - h + 1), X.shape[1] - w + 1)
    for i in range(Y.shape[0]):
        for j in range(Y.shape[1]):
            Y[i, j] = (X[i : i + h, j : j + w] * K).sum()
    return Y

X = torch.tensor([[0.0, 1.0, 2.0], [3.0, 4.0, 5.0], [6.0, 7.0, 8.0]])
K = torch.Tensor([[0.0, 1.0], [2.0, 3.0]])
corr2d(X, K)
```

⮒  tensor([[19., 25.],
        [37., 43.]])

```python
class Conv2D(nn.Module):
    def __init__(self, kernel_size):
        super().__init__()
        self.weight = nn.Parameter(torch.rand(kernel_size))
        self.bias = nn.Parameter(torch.zeros(1))

    def forward(self, x):
        return corr2d(x, self.weight) + self.bias

X = torch.ones((6, 8))
X[:, 2:6] = 0
X
```

⮒  tensor([[1., 1., 0., 0., 0., 0., 1., 1.],
        [1., 1., 0., 0., 0., 0., 1., 1.],
        [1., 1., 0., 0., 0., 0., 1., 1.],
        [1., 1., 0., 0., 0., 0., 1., 1.],
        [1., 1., 0., 0., 0., 0., 1., 1.],
        [1., 1., 0., 0., 0., 0., 1., 1.]])

```python
K = torch.tensor([[1.0, -1.0]])
```

```python
Y = corr2d(X, K)
Y
```

⮒  tensor([[ 0.,  1.,  0.,  0.,  0., -1.,  0.],
        [ 0.,  1.,  0.,  0.,  0., -1.,  0.],
        [ 0.,  1.,  0.,  0.,  0., -1.,  0.],
        [ 0.,  1.,  0.,  0.,  0., -1.,  0.],
        [ 0.,  1.,  0.,  0.,  0., -1.,  0.],
        [ 0.,  1.,  0.,  0.,  0., -1.,  0.]])

```python
corr2d(X.t(), K)
```

⮒  tensor([[0., 0., 0., 0., 0.],
        [0., 0., 0., 0., 0.],
        [0., 0., 0., 0., 0.],
        [0., 0., 0., 0., 0.],
        [0., 0., 0., 0., 0.],
        [0., 0., 0., 0., 0.],
        [0., 0., 0., 0., 0.],
        [0., 0., 0., 0., 0.]])

```
conv2d = nn.LazyConv2d(1, kernel_size=(1, 2), bias=False)
X = X.reshape((1, 1, 6, 8))
Y = Y.reshape((1, 1, 6, 7))
lr = 3e-2

for i in range(10):
    Y_hat = conv2d(X)
    l = (Y_hat - Y) ** 2
    conv2d.zero_grad()
    l.sum().backward()

    conv2d.weight.data[:] -= lr * conv2d.weight.grad
    if (i + 1) % 2 == 0:
        print(f'epoch {i + 1}, loss {l.sum():.3f}')
```

```
epoch 2, loss 9.499
epoch 4, loss 2.401
epoch 6, loss 0.733
epoch 8, loss 0.258
epoch 10, loss 0.099
```

```
conv2d.weight.data.reshape((1, 2))
```

```
tensor([[ 1.0188, -0.9554]])
```

## 7.3 Padding and Stride

```
import torch
from torch import nn


def comp_conv2d(conv2d, X):
    X = X.reshape((1, 1) + X.shape)
    Y = conv2d(X)
    return Y.reshape(Y.shape[2:])

conv2d = nn.LazyConv2d(1, kernel_size=3, padding=1)
X = torch.rand(size=(8, 8))
comp_conv2d(conv2d, X).shape
```

```
torch.Size([8, 8])
```

```
conv2d = nn.LazyConv2d(1, kernel_size=(5, 3), padding=(2, 1))
comp_conv2d(conv2d, X).shape
```

```
torch.Size([8, 8])
```

```
conv2d = nn.LazyConv2d(1, kernel_size=3, padding=1, stride=2)
comp_conv2d(conv2d, X).shape
```

```
torch.Size([4, 4])
```

```
conv2d = nn.LazyConv2d(1, kernel_size=(3, 5), padding=(0, 1), stride=(3, 4))
comp_conv2d(conv2d, X).shape
```

```
torch.Size([2, 2])
```

## 7.4 Multiple Input and Multiple Output Channels

```
import torch
from d2l import torch as d2l


def corr2d_multi_in(X, K):
    return sum(d2l.corr2d(x, k) for x, k in zip(X, K))

X = torch.tensor([[[0.0, 1.0, 2.0], [3.0, 4.0, 5.0], [6.0, 7.0, 8.0]],
                  [[1.0, 2.0, 3.0], [4.0, 5.0, 6.0], [7.0, 8.0, 9.0]]])
K = torch.Tensor([[[0.0, 1.0], [2.0, 3.0]], [[1.0, 2.0], [3.0, 4.0]]])

corr2d_multi_in(X, K)
```

```
tensor([[ 56.,  72.],
        [104., 120.]])
```

```python
def corr2d_multi_in_out(X, K):
    return torch.stack([corr2d_multi_in(X, k) for k in K], 0)
```

```python
K = torch.stack((K, K + 1, K + 2), 0)
K.shape
```

➔   torch.Size([3, 2, 2, 2])

```python
corr2d_multi_in_out(X, K)
```

➔   tensor([[[ 56.,  72.],
            [104., 120.]],

           [[ 76., 100.],
            [148., 172.]],

           [[ 96., 128.],
            [192., 224.]]])

```python
def corr2d_multi_in_out_1x1(X, K):
    c_i, h, w = X.shape
    c_o = K.shape[0]
    X = X.reshape((c_i, h * w))
    K = K.reshape((c_o, c_i))
    Y = torch.matmul(K, X)
    return Y.reshape((c_o, h, w))
```

```python
X = torch.normal(0, 1, (3, 3, 3))
K = torch.normal(0, 1, (2, 3, 1, 1))
Y1 = corr2d_multi_in_out_1x1(X, K)
Y2 = corr2d_multi_in_out(X, K)
assert float(torch.abs(Y1-Y2).sum()) < 1e-6
```

## ⌄ 7.5 Pooling

```python
import torch
from torch import nn
from d2l import torch as d2l
```

```python
def pool2d(X, pool_size, mode="max"):
    p_h, p_w = pool_size
    Y = torch.zeros((X.shape[0] - p_h + 1, X.shape[1] - p_w + 1))
    for i in range(Y.shape[0]):
        for j in range(Y.shape[1]):
            if mode == 'max':
                Y[i, j] = X[i: i + p_h, j : j + p_w].max()
            elif mode == 'avg':
                Y[i, j] = X[i: i + p_h, j : j + p_w].mean()
    return Y
```

```python
X = torch.tensor([[0.0, 1.0, 2.0], [3.0, 4.0, 5.0], [6.0, 7.0, 8.0]])
pool2d(X, (2, 2))
```

➔   tensor([[4., 5.],
           [7., 8.]])

```python
pool2d(X, (2, 2), 'avg')
```

➔   tensor([[2., 3.],
           [5., 6.]])

```python
X = torch.arange(16, dtype=torch.float32).reshape((1, 1, 4, 4))
X
```

➔   tensor([[[[ 0.,  1.,  2.,  3.],
             [ 4.,  5.,  6.,  7.],
             [ 8.,  9., 10., 11.],
             [12., 13., 14., 15.]]]])

```python
pool2d = nn.MaxPool2d(3)
pool2d(X)
```

➔   tensor([[[[10.]]]])

```
pool2d = nn.MaxPool2d(3, padding=1,stride=2)
pool2d(X)
```

⤷ 
```
tensor([[[[ 5.,  7.],
          [13., 15.]]]])
```

```
pool2d = nn.MaxPool2d((2, 3), stride=(2, 3), padding=(0, 1))
pool2d(X)
```

⤷ 
```
tensor([[[[ 5.,  7.],
          [13., 15.]]]])
```

```
X = torch.cat((X, X + 1), 1)
X
```

⤷ 
```
tensor([[[[ 0.,  1.,  2.,  3.],
          [ 4.,  5.,  6.,  7.],
          [ 8.,  9., 10., 11.],
          [12., 13., 14., 15.]],

         [[ 1.,  2.,  3.,  4.],
          [ 5.,  6.,  7.,  8.],
          [ 9., 10., 11., 12.],
          [13., 14., 15., 16.]]]])
```

```
pool2d = nn.MaxPool2d(3, padding=1, stride=2)
pool2d(X)
```

⤷ 
```
tensor([[[[ 5.,  7.],
          [13., 15.]],

         [[ 6.,  8.],
          [14., 16.]]]])
```

## ⌄ 7.6 Convolutional Neural Networks (LeNet)

```
import torch
from torch import nn
from d2l import torch as d2l
```

```
def init_cnn(module):
    if type(module) == nn.Linear or type(module) == nn.Conv2d:
        nn.init.xavier_uniform_(module.weight)
```

```
class LeNet(d2l.Classifier):
    def __init__(self, lr=0.1, num_classes=10):
        super().__init__()
        self.save_hyperparameters()
        self.net = nn.Sequential(
            nn.LazyConv2d(6, kernel_size=5, padding=2),
            nn.Sigmoid(),
            nn.AvgPool2d(kernel_size=2, stride=2),
            nn.LazyConv2d(16, kernel_size=5),
            nn.Sigmoid(),
            nn.AvgPool2d(kernel_size=2, stride=2),
            nn.Flatten(),
            nn.LazyLinear(120),
            nn.Sigmoid(),
            nn.LazyLinear(84),
            nn.Sigmoid(),
            nn.LazyLinear(num_classes),
        )
```

```
@d2l.add_to_class(d2l.Classifier)
def layer_summary(self, X_shape):
    X = torch.randn(*X_shape)
    for layer in self.net:
        X = layer(X)
        print(layer.__class__.__name__, 'output shape:\t', X.shape)

model = LeNet()
model.layer_summary((1, 1, 28, 28))
```

⤷ 
```
Conv2d output shape:     torch.Size([1, 6, 28, 28])
Sigmoid output shape:    torch.Size([1, 6, 28, 28])
AvgPool2d output shape:  torch.Size([1, 6, 14, 14])
Conv2d output shape:     torch.Size([1, 16, 10, 10])
Sigmoid output shape:    torch.Size([1, 16, 10, 10])
```
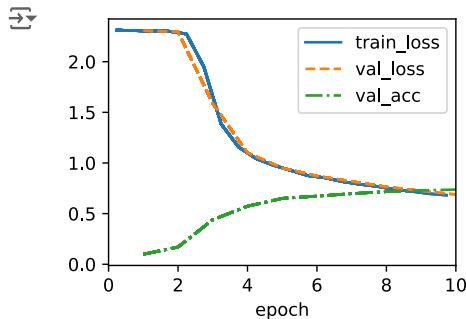
```
AvgPool2d output shape:    torch.Size([1, 16, 5, 5])
Flatten output shape:      torch.Size([1, 400])
Linear output shape:       torch.Size([1, 120])
Sigmoid output shape:      torch.Size([1, 120])
Linear output shape:       torch.Size([1, 84])
Sigmoid output shape:      torch.Size([1, 84])
Linear output shape:       torch.Size([1, 10])
```

```
trainer = d2l.Trainer(max_epochs=10, num_gpus=1)
data = d2l.FashionMNIST(batch_size=128)
model = LeNet(lr=0.1)
model.apply_init([next(iter(data.get_dataloader(True)))[0]], init_cnn)
trainer.fit(model, data)
```



## 8.2 Networks Using Blocks (VGG)

```
import torch
from torch import nn
from d2l import torch as d2l
```

```
def vgg_block(num_convs, out_channels):
    layers = []
    for _ in range(num_convs):
        layers.append(nn.LazyConv2d(out_channels, kernel_size=3, padding=1))
        layers.append(nn.ReLU())
    layers.append(nn.MaxPool2d(kernel_size=2, stride=2))
    return nn.Sequential(*layers)
```

```
class VGG(d2l.Classifier):
    def __init__(self, arch, lr=0.1, num_classes=10):
        super().__init__()
        self.save_hyperparameters()
        conv_blks = []
        for (num_convs, out_channels) in arch:
            conv_blks.append(vgg_block(num_convs, out_channels))
        self.net = nn.Sequential(
            *conv_blks, nn.Flatten(),
            nn.LazyLinear(4096), nn.ReLU(), nn.Dropout(0.5),
            nn.LazyLinear(4096), nn.ReLU(), nn.Dropout(0.5),
            nn.LazyLinear(num_classes))
        self.net.apply(d2l.init_cnn)
```

```
VGG(arch=((1, 64), (1, 128), (2, 256), (2, 512), (2, 512))).layer_summary(
    (1, 1, 224, 224))
```
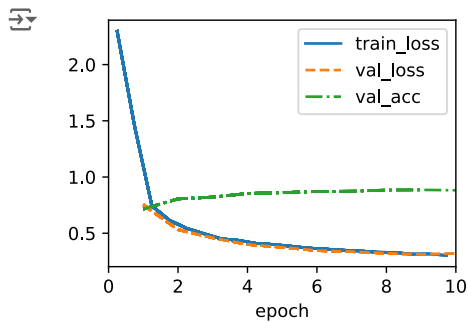
```
Sequential output shape:        torch.Size([1, 64, 112, 112])
Sequential output shape:        torch.Size([1, 128, 56, 56])
Sequential output shape:        torch.Size([1, 256, 28, 28])
Sequential output shape:        torch.Size([1, 512, 14, 14])
Sequential output shape:        torch.Size([1, 512, 7, 7])
Flatten output shape:      torch.Size([1, 25088])
Linear output shape:       torch.Size([1, 4096])
ReLU output shape:         torch.Size([1, 4096])
Dropout output shape:      torch.Size([1, 4096])
Linear output shape:       torch.Size([1, 4096])
ReLU output shape:         torch.Size([1, 4096])
Dropout output shape:      torch.Size([1, 4096])
Linear output shape:       torch.Size([1, 10])
```

```
model = VGG(arch=((1, 16), (1, 32), (2, 64), (2, 128), (2, 128)), lr=0.01)
trainer = d2l.Trainer(max_epochs=10, num_gpus=1)
data = d2l.FashionMNIST(batch_size=128, resize=(224, 224))
model.apply_init([next(iter(data.get_dataloader(True)))[0]], d2l.init_cnn)
trainer.fit(model, data)
```

## 8.6 Residual Networks (ResNet) and ResNeXt

```python
import torch
from torch import nn
from torch.nn import functional as F
from d2l import torch as d2l


class Residual(nn.Module):
    def __init__(self, num_channels, use_1x1conv=False, strides=1):
        super().__init__()
        self.conv1 = nn.LazyConv2d(num_channels, kernel_size=3, padding=1,
                                   stride=strides)
        self.conv2 = nn.LazyConv2d(num_channels, kernel_size=3, padding=1)
        if use_1x1conv:
            self.conv3 = nn.LazyConv2d(num_channels, kernel_size=1,
                                       stride=strides)
        else:
            self.conv3 = None
        self.bn1 = nn.LazyBatchNorm2d()
        self.bn2 = nn.LazyBatchNorm2d()

    def forward(self, X):
        Y = F.relu(self.bn1(self.conv1(X)))
        Y = self.bn2(self.conv2(Y))
        if self.conv3:
            X = self.conv3(X)
        Y += X
        return F.relu(Y)


blk = Residual(3)
X = torch.randn(4, 3, 6, 6)
blk(X).shape
```

```
torch.Size([4, 3, 6, 6])
```

```python
blk = Residual(6, use_1x1conv=True, strides=2)
blk(X).shape
```

```
torch.Size([4, 6, 3, 3])
```

```python
class ResNet(d2l.Classifier):
    def b1(self):
        return nn.Sequential(
            nn.LazyConv2d(64, kernel_size=7, stride=2, padding=3),
            nn.LazyBatchNorm2d(), nn.ReLU(),
            nn.MaxPool2d(kernel_size=3, stride=2, padding=1))


@d2l.add_to_class(ResNet)
def block(self, num_residuals, num_channels, first_block=False):
    blk = []
    for i in range(num_residuals):
        if i == 0 and not first_block:
            blk.append(Residual(num_channels, use_1x1conv=True, strides=2))
        else:
            blk.append(Residual(num_channels))
    return nn.Sequential(*blk)


@d2l.add_to_class(ResNet)
def __init__(self, arch, lr=0.1, num_classes=10):
    super(ResNet, self).__init__()
    self.save_hyperparameters()
    self.net = nn.Sequential(self.b1())
```

```
    for i, b in enumerate(arch):
        self.net.add_module(f'b{i+2}', self.block(*b, first_block=(i==0)))
    self.net.add_module('last', nn.Sequential(
        nn.AdaptiveAvgPool2d((1, 1)), nn.Flatten(),
        nn.LazyLinear(num_classes)))
    self.net.apply(d2l.init_cnn)


class ResNet18(ResNet):
    def __init__(self, lr=0.1, num_classes=10):
        super().__init__(((2, 64), (2, 128), (2, 256), (2, 512)),
                         lr, num_classes)


ResNet18().layer_summary((1, 1, 96, 96))
```
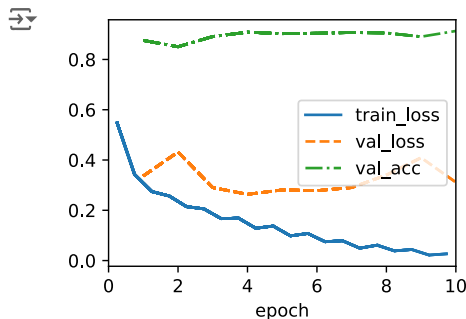
```
Sequential output shape:         torch.Size([1, 64, 24, 24])
Sequential output shape:         torch.Size([1, 64, 24, 24])
Sequential output shape:         torch.Size([1, 128, 12, 12])
Sequential output shape:         torch.Size([1, 256, 6, 6])
Sequential output shape:         torch.Size([1, 512, 3, 3])
Sequential output shape:         torch.Size([1, 10])
```

```
model = ResNet18(lr=0.01)
trainer = d2l.Trainer(max_epochs=10, num_gpus=1)
data = d2l.FashionMNIST(batch_size=128, resize=(96, 96))
model.apply_init([next(iter(data.get_dataloader(True)))[0]], d2l.init_cnn)
trainer.fit(model, data)
```



## Discussion

## 7.1 From Fully Connected Layers to Convolutions

In the original FC-MLP approach, handling high-dimensional data such as images incurs significant computational costs and results in a large number of parameters, making it highly inefficient.

In the CNN approach, computation is more efficient and requires significantly fewer parameters compared to the Fully-Connected architecture. Additionally, it can be easily parallelized across GPU cores.

Unlike the MLP approach, which is sensitive to the position of objects, the CNN approach has a property called Translation Invariance, allowing it to learn how an object looks regardless of where it is located in the image.

Key concepts used in CNN

- Translation Invariance: It ignores the position-related information of objects.
  As a result, values can be represented using significantly fewer coefficients and parameters.
- Locality: It only gathers information from adjacent pixels.
  This also helps reduce the number of parameters used.

## 7.2 Convolutions for Images

As the kernel slides input data, it creates an output layer that represents the features of the adjacent area defined by the kernel size at each position.

The weights of the kernel are initially set randomly, just like in a Fully Connected layer. When considering larger kernels and consecutive layers, if kernels are designed for specific purposes, it would be difficult to determine whether each kernel is functioning correctly.

Therefore, by calculating the gradient based on the input-output, the model can learn to find the appropriate kernels.

The output of the convolution layer is also referred to as a feature map, as it represents the learned features in the spatial dimensions.

The receptive field represents all the elements from previous layers that can influence the calculation of a specific position 'x' during forward propagation.

As the layer gets deeper, the portion of the input data that a specific part of the layer covers becomes larger, resulting in an increase in the receptive field.

Therefore, as the layer gets deeper, it can be considered to include more global information.

## 7.3 Padding and Stride

Padding is a method used to address the issue of losing information from the outer parts of the input data in conventional approaches.

This issue can become even more severe in a repeated convolution.

Typically, the padding method involves adding extra pixels with a value of zero around the input data.

By setting the kernel size to an odd number, there is the advantage of being able to add an equal number of padding pixels on the left and right, as well as the top and bottom.

When moving the kernel across the input data, it typically shifts by one element.

However, for computational efficiency or downsampling, the kernel may skip intermediate positions and move by larger steps.

At this point, the number of rows or columns the kernel moves in a slide is called the stride.

## 7.4 Multiple Input and Multiple Output Channels

When the input data has multiple channels, the kernel must have the same number of channels as the input data.

This is because the kernel needs to perform cross-correlation with the input data across all channels.

For example, if the input data has a size of 32 x 32 x 3, the kernel's channel size must also be 3.

Additionally, if the kernel size is 5 x 5 x 3, the output data would have a size of 28 x 28 x 1.

Thus, the output data must have 1 channel.

To increase the number of output data channels, the number of kernels must be increased.

By concatenating all the output data, where each kernel's cross-correlation produces 1 channel, the number of output data channels is increased.

The 1x1 Convolution layer does not affect the height and width values but is used to combine information across the channels.

## 7.5 Pooling

As the neural network goes deeper, the receptive field becomes progressively larger.

Reducing the spatial resolution accelerates the process and allows the convolution kernel to cover a larger and more effective area.

The pooling layer reduces the spatial sensitivity of the convolution layer (translation invariance) and allows for downsampling of the spatial dimensions.

Since the pooling layer does not have kernels, it does not contain any parameters.

Instead, it is divided into maximum pooling, which calculates the maximum value, and average pooling, which calculates the average value for each region.

When receiving input data with multiple channels, the pooling layer performs pooling for each channel individually, rather than summing the channels within the convolutional layer.

## 7.6 Convolutional Neural Networks (LeNet)

LeNet is the first CNN that gained widespread attention for its performance in the field of computer vision.

LeNet is divided into two parts
- a convolutional encoder consisting of two convolutional layers
- a dense block consisting of three fully connected layers

## 8.2 Networks Using Blocks (VGG)

VGG demonstrated that deeper and wider networks lead to better performance.
For example, two 3x3 convolutions have the same effect as a single 5x5 convolution.
Stacking such 3x3 convolutions has become the gold standard in deep networks.

The VGG network can also be divided into two parts.

- consisting mostly of convolutional and pooling layers
- consisting of fully connected layers that are identical to those in AlexNet

The VGG network has five convolutional blocks: the first two blocks each contain a single convolutional layer, and the last three blocks contain two convolutional layers each.
The first block has 64 channels, and each subsequent block doubles the number of output channels, eventually reaching 512 channels.

## 8.6 Residual Networks (ResNet) and ResNeXt

ResNet makes it easier to learn by using f(x) - x instead of directly learning the target function f(x), where x is the input.
By using such residual blocks, the input values can propagate more quickly through the network.