

## 7.0. Convolutional Neural Networks

### 7.1. From Fully Connected Layers to Convolutions

- **Invariance**: When detecting objects in an image, it is important to recognize them regardless of where they appear in the image. This concept is known as invariance. For example, planes usually fly, and pigs usually stay on the ground, but the system should still recognize a pig even if it appears at the top of the image. Convolutional neural networks (CNNs) leverage this property to recognize patterns efficiently.
- **Constraining the MLP**: Consider an MLP (Multilayer Perceptron) that processes 2D images as inputs. In this case, the hidden representations also need to be structured as 2D matrices, and each hidden unit would receive input from every pixel of the image. This creates an immense number of parameters, making the model impractical. For example, processing a one-megapixel image with 1000 hidden units requires billions of parameters.

The equation is :

$$\begin{aligned} [\mathbf{H}]_{i,j} &= [\mathbf{U}]_{i,j} + \sum_k \sum_l [\mathbf{W}]_{i,j,k,l} [\mathbf{X}]_{k,l} \\ &= [\mathbf{U}]_{i,j} + \sum_a \sum_b [\mathbf{V}]_{i,j,a,b} [\mathbf{X}]_{i+a,j+b}. \end{aligned}$$

- **Translation Invariance**: Translation invariance means that if an input image shifts in space, the hidden representation shifts correspondingly without changing the structure of the learned features. This can be achieved by ensuring that the weight tensor is not dependent on the position in the image. Thus, the weights are constant across the entire image, simplifying the equation to:

$$[\mathbf{H}]_{i,j} = u + \sum_a \sum_b [\mathbf{V}]_{a,b} [\mathbf{X}]_{i+a,j+b}.$$

- **Locality**: Locality refers to the assumption that the relevant information for processing a particular pixel in the image can be found in the surrounding region, rather than distant parts of the image. By restricting the influence of distant pixels, CNNs reduce the number of parameters and increase the efficiency of the model. This assumption leads to the use of convolutional filters that focus on small neighborhoods in the image.

$$[\mathbf{H}]_{i,j} = u + \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} [\mathbf{V}]_{a,b} [\mathbf{X}]_{i+a,j+b}.$$

- **Convolutions**: Convolution is a fundamental operation in CNNs. It involves applying a filter (kernel) to small regions of the input image to learn spatial features. The 2D convolution is mathematically defined as:

$$(f * g)(i, j) = \sum_a \sum_b f(a, b) g(i - a, j - b).$$

This operation calculates the output by summing the weighted pixel values within the filter's local receptive field, resulting in an efficient reduction of parameters compared to fully connected layers.

- **Channels**: Images typically consist of three color channels: red, green, and blue (RGB). In CNNs, these channels are treated as additional dimensions, making the input a 3D tensor. The hidden representations are similarly structured as 3D tensors with multiple feature maps (channels) stacked together. The equation for convolution with multiple channels is:

$$[\mathbf{H}]_{i,j,d} = \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} \sum_c [\mathbf{V}]_{a,b,c,d} [\mathbf{X}]_{i+a,j+b,c},$$

Here,  $d$  indexes the output channels in the hidden representation  $\mathbf{H}$ , allowing the model to learn multiple feature representations at each spatial location.

### 7.2. Convolutions for Images

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

- The Cross-Correlation Operation

```
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.]])
```

- Convolutional Layers

```
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
```

- Object Edge Detection in Images

```
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.]])
```

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

```
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.]])
```

```
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.]])
```

- Learning a Kernel

```
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 # Learning rate
```

```
for i in range(10):
    Y_hat = conv2d(X)
    l = (Y_hat - Y) ** 2
    conv2d.zero_grad()
    l.sum().backward()
    # Update the kernel
    conv2d.weight.data[:] -= lr * conv2d.weight.grad
    if (i + 1) % 2 == 0:
        print(f'epoch {i + 1}, loss {l.sum():.3f}')
```

```
epoch 2, loss 13.800
epoch 4, loss 3.481
epoch 6, loss 1.061
epoch 8, loss 0.374
epoch 10, loss 0.143
C:\ProgramData\Anaconda3\lib\site-packages\torch\nn\modules\lazy.py:178: UserWarning: Lazy modules are a new feature under heavy development so changes to the AP
warnings.warn('Lazy modules are a new feature under heavy development ')
```

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

```
tensor([[ 1.0225, -0.9464]])
```

- **Cross-Correlation and Convolution** : Cross-correlation is a common operation used in CNNs, but convolution can also be performed by flipping the kernel both horizontally and vertically before applying the cross-correlation operation.
  - The kernels used in CNNs are learned from the data, and the output remains the same whether the layer performs cross-correlation or strict convolution.
  - If a CNN layer learns a kernel through cross-correlation, flipping it results in the equivalent convolutional kernel. Thus, the operation is computationally similar, and in practice, deep learning literature often refers to cross-correlation as convolution for simplicity.
- **Feature Map and Receptive Field** : In CNNs, the output of a convolutional layer is often called a feature map, which represents the learned spatial features of an image. The receptive field refers to the region of the input that contributes to a particular output element in a CNN. As layers deepen in the network, the receptive field grows, allowing the model to capture features from a larger portion of the image.

### 7.3. Padding and Stride

#### • 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
```

```
⚡ C:\ProgramData\Anaconda3\lib\site-packages\torch\nn\modules\lazy.py:178: UserWarning: Lazy modules are a new feature under heavy development so changes to the API
  warnings.warn('Lazy modules are a new feature under heavy development '
  torch.Size([8, 8])
```

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

```
comp_conv2d(conv2d, X).shape
```

```
⚡ torch.Size([8, 8])
```

#### • Stride

```
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
```

#### • Multiple Input Channels

```
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.]])
```

#### • Multiple Output Channels

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

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

```
↗ torch.Size([3, 2, 2, 2])
```

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

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

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

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

#### • 1×1 Convolutional Layer

```
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))
```

```
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

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

#### • Maximum Pooling and Average Pooling

```
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
```

```
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.]])
```

```
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.]])
```

#### • Padding and Stride

```
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.]]]])
```

```
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.]]]])
```

#### • Multiple Channels

```
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.]]]])
```

```
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.]],
         [[ 1.,  2.,  3.,  4.],
           [ 5.,  6.,  7.,  8.],
           [ 9., 10., 11., 12.],
           [13., 14., 15., 16.]],
         [[ 2.,  3.,  4.,  5.],
           [ 6.,  7.,  8.,  9.],
           [10., 11., 12., 13.],
           [14., 15., 16., 17.]]]])
```

### 7.6. Convolutional Neural Networks(LeNet)

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

```
def init_cnn(module):
    """Initialize weights for CNNs."""
    if type(module) == nn.Linear or type(module) == nn.Conv2d:
        nn.init.xavier_uniform_(module.weight)

class LeNet(d2l.Classifier):
    """The LeNet-5 model."""
    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: %s' % 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])

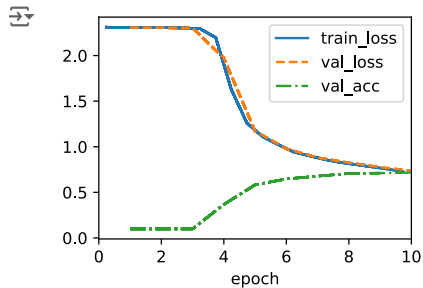
```

### • Training

```

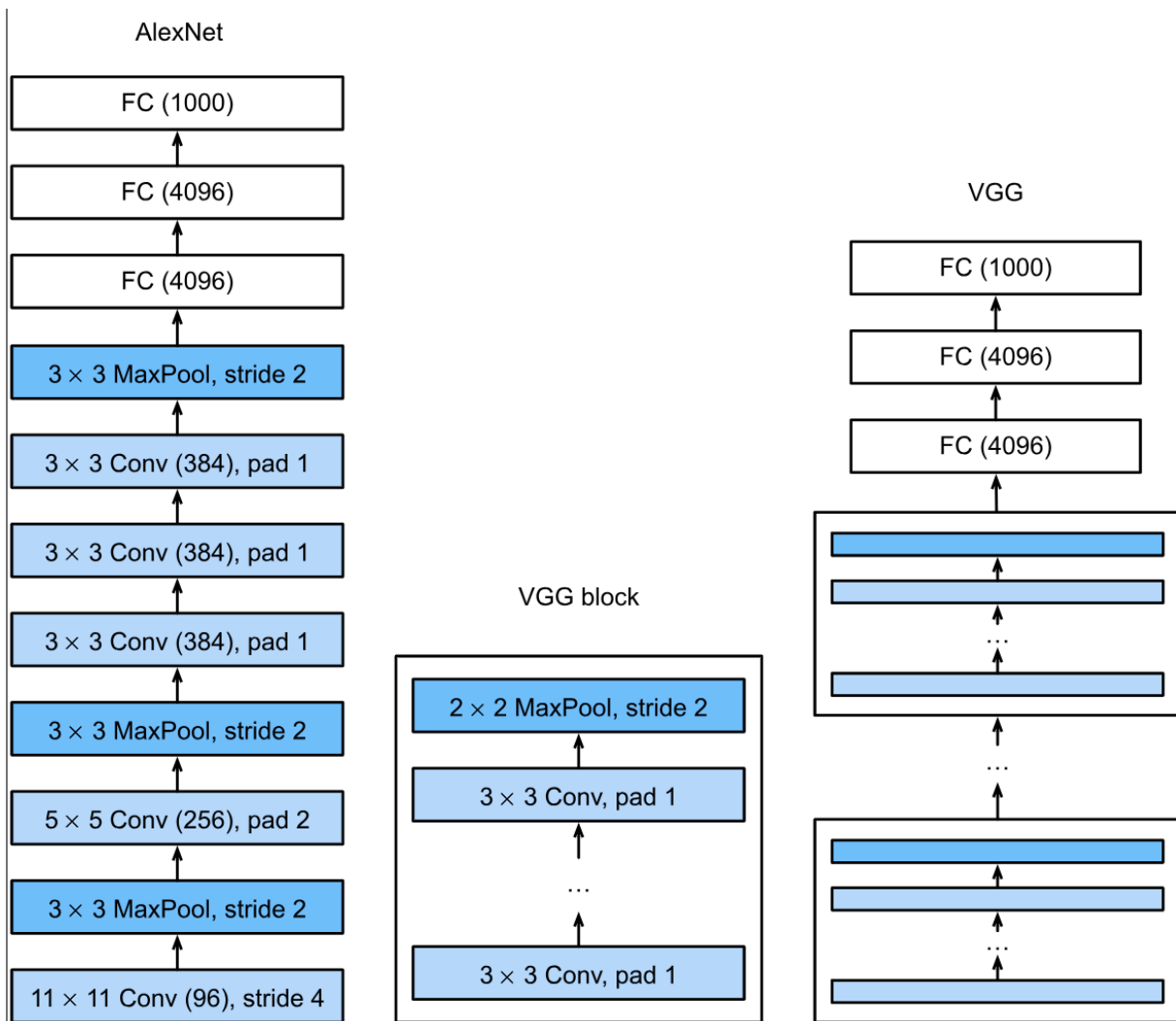
trainer = d2l.Trainer(max_epochs=10, num_gpus=1)
data = d2l.FashionMNIST(batch_size=128)
model = LeNet(lr=0.1)
model.apply_init([next(iter(data.get_dataloader(True)))[0]), init_cnn]
trainer.fit(model, data)

```



## ✓ 8.0. Modern Convolutional Neural Networks

### ✓ 8.2. Networks Using Blocks(VGG)



```

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

```

- **VGG Block**

```
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)
```

- **VGG Network**

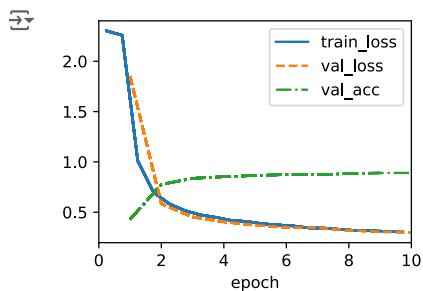
```
class VGG(d2l.Classifier):
    def __init__(self, arch, lr=0.1, num_classes=10):
        super().__init__()
        self.save_hyperparameters()
        conv_blks = []
        for (num_convs, out_channels) in arch:
            conv_blks.append(vgg_block(num_convs, out_channels))
        self.net = 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])
```

- **Training**

```
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)
```



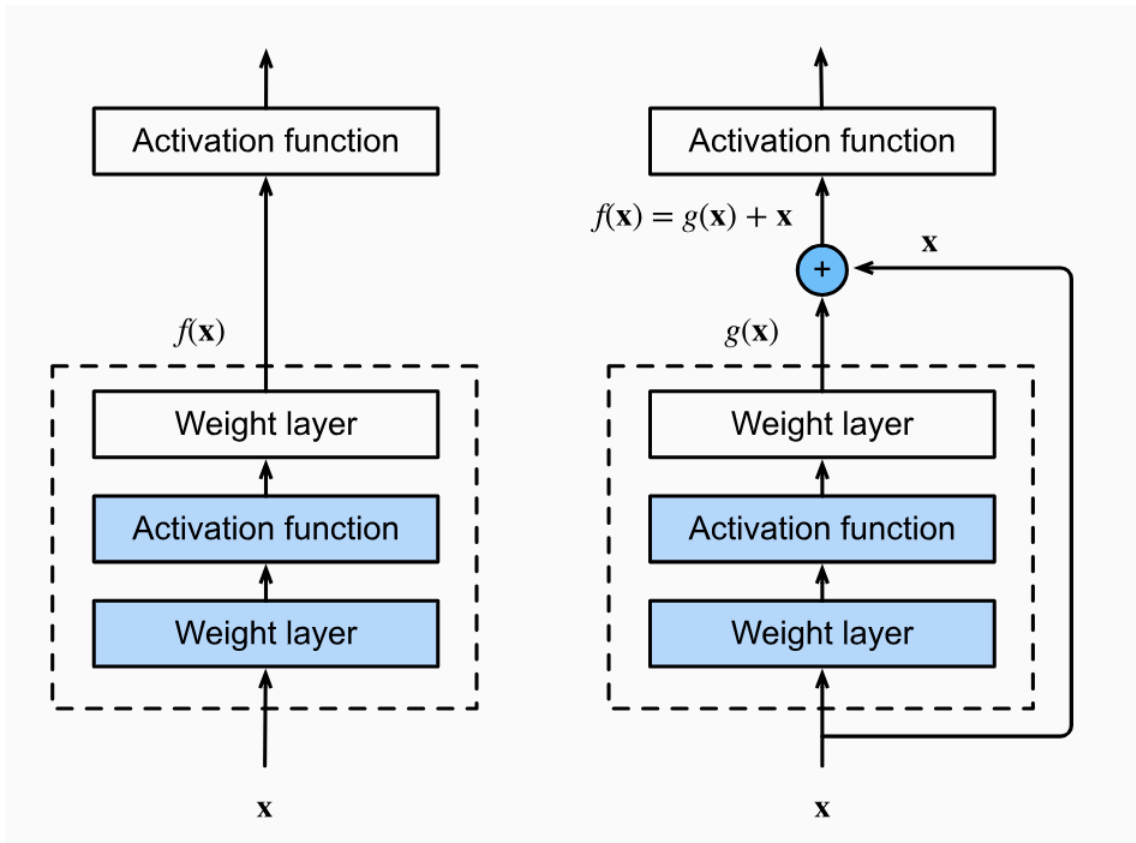
## ✓ 8.6. Residual Networks(ResNet) and ResNeXt

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

- **Function Classes** : Function classes refer to the set of functions a network can represent. Ideally, we want to approximate the true function  $f^*$  within this set. When the function classes are nested, adding layers increases the network's expressiveness. However, for non-nested classes, adding complexity might not improve approximation and can even lead to worse results.

$$f_{\mathcal{F}}^* \stackrel{\text{def}}{=} \underset{f}{\operatorname{argmin}} L(\mathbf{X}, \mathbf{y}, f) \text{ subject to } f \in \mathcal{F}.$$

- **Residual Blocks**



```
class Residual(nn.Module):
    """The Residual block of ResNet models."""
    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])
```

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

↗ torch.Size([4, 6, 3, 3])
```

#### • ResNet Model

```
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.blk())
    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.Linear(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])

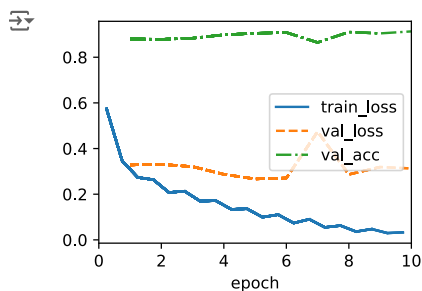
```

### • Training

```

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 and Exercise

### ✓ 7.1.5. Discussion and Exercises

- **Discussion** : A major challenge is reducing the number of parameters without compromising the model's expressiveness. CNNs address this by introducing nested function classes and convolutional kernels, which significantly reduce computational complexity and make previously infeasible problems tractable.

Additionally, the inclusion of multiple channels (such as RGB or even more for specialized imagery) allows CNNs to capture a broader range of information, increasing their power to process complex datasets, like hyperspectral satellite images. This section sets the stage for further exploration of CNNs' capacity to handle high-dimensional data efficiently.

- **Exercise 6** : Prove that the convolution is symmetric, i.e.,  $f * g = g * f$

#### 1. Convolution definition:

$$(f * g)(i) = \sum_a f(a)g(i - a)$$

#### 2. Reindexing the convolution for $(g * f)$ :

$$(g * f)(i) = \sum_a g(a)f(i - a)$$

Using the substitution  $(b = i - a)$ , this becomes:

$$(g * f)(i) = \sum_b g(i - b)f(b)$$

So,

$$(f * g)(i) = \sum_b f(b)g(i - b)$$

This demonstrates that convolution is symmetric:  $f * g = g * f$

### 7.3.3 Discussion

- **Discussion** : Zero-padding is widely used for its computational simplicity and efficiency. It allows CNNs to infer positional information about where the image "whitespace" is located. Alternatives to zero-padding exist, but their use cases are less clear unless artifacts from zero-padding occur.

### 7.4.4. Discussion

- **Discussion** : Channels allow CNNs to analyze multiple features (like edges and shapes) simultaneously, providing a balance between reduced parameters (due to translation invariance and locality) and the need for diverse, expressive models.

However, this flexibility comes with a computational cost. For an image of size  $h \times w$ , the computational complexity of a convolution is  $O(h \cdot w \cdot k^2)$  for a kernel size  $k \times k$ . With input and output channels, this increases to  $O(h \cdot w \cdot k^2 \cdot c_i \cdot c_o)$ . For a 256x256 image, a 5x5 kernel, and 128 channels, this results in over 53 billion operations. Techniques like block-diagonal operations (as seen in architectures like **ResNeXt**) can help reduce these costs.

### 7.5.5. Exercises

#### 2. Prove that max-pooling cannot be implemented through a convolution alone :

Max-pooling cannot be implemented through a convolution alone because max-pooling is a non-linear operation. Max-pooling selects the maximum value from a given region or window, while convolution is a linear operation that computes a weighted sum of the input values. The non-linearity of max-pooling is essential to its function, as it makes decisions based on the maximum value in the region. Since convolution only computes linear combinations, it cannot replicate this behavior. Therefore, max-pooling and convolution perform fundamentally different roles, and convolution alone cannot implement the max-pooling operation.