

> Installation

[] ↳ 숨겨진 셀 1개

✓ 7. Convolutional Neural Networks

✓ 7.1. From Fully Connected Layers to Convolutions

✓ Discussions & Exercises

Convolutional neural networks (CNNs) are one creative way that machine learning has embraced for exploiting some of the known structure in natural images.

7.1.1. Invariance

CNNs systematize this idea of *spatial invariance*, exploiting it to learn useful representations with fewer parameters.

7.1.2.1. Translation Invariance

This implies that a shift in the input \mathbf{X} should simply lead to a shift in the hidden representation \mathbf{H} . This is only possible if \mathbf{V} and \mathbf{U} do not actually depend on (i, j) .

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

7.1.2.2. Locality

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

7.1.4. Channels

Adding channels allow us to bring back some of the complexity that was lost due to the restrictions imposed on the convolutional kernel by locality and translation invariance.

✓ 7.2. Convolutions for Images

```
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.]])
```

```
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.]])
```

```
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.]])
```

```
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 9.008
   epoch 4, loss 1.515
   epoch 6, loss 0.256
   epoch 8, loss 0.043
   epoch 10, loss 0.008
```

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

```
⇒ tensor([[ 0.9865, -0.9823]])
```

✓ Discussions & Exercises

7.2.2. Convolutional Layers

A convolutional layer cross-correlates the input and kernel and adds a scalar bias to produce an output. The two parameters of a convolutional layer are the kernel and the scalar bias. When training models based on convolutional layers, we typically initialize the kernels randomly, just as we would with a fully connected layer.

In $h \times w$ convolution or an $h \times w$ convolution kernel, the height and width of the convolution kernel are h and w , respectively. We also refer to a convolutional layer with an $h \times w$ convolution kernel simply as an $h \times w$ convolutional layer.

7.2.6. Feature Map and Receptive Field

The convolutional layer output is sometimes called a *feature map*, as it can be regarded as the learned representations (features) in the spatial dimensions (e.g., width and height) to the subsequent layer. In CNNs, for any element x of some layer, its *receptive field* refers to all the

elements (from all the previous layers) that may affect the calculation of x during the forward propagation. Note that the receptive field may be larger than the actual size of the input.

✓ 7.3. Padding and Stride

```
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])
```

✓ Discussions & Exercises

7.3.1. Padding

Padding can increase the height and width of the output. This is often used to give the output the same height and width as the input to avoid undesirable shrinkage of the output. Moreover, it ensures that all pixels are used equally frequently. Typically we pick symmetric padding on both sides of the input height and width.

7.3.2. Stride

We refer to the number of rows and columns traversed per slide as *stride*.

7.3.3. Summary and Discussion

By default, the padding is 0 and the stride is 1.

✓ 7.4. Multiple Input and Multiple Output 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.]])
```

```
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.]])
```

```
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
```

✓ Discussions & Exercises

7.4.4. Discussion

Channels allow us to combine the best of both worlds: MLPs that allow for significant nonlinearities and convolutions that allow for *localized* analysis of features. In particular, channels allow the CNN to reason with multiple features, such as edge and shape detectors at the same time. They also offer a practical trade-off between the drastic parameter reduction arising from translation invariance and locality, and the need for expressive and diverse models in computer vision.

✓ 7.5. 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.]])

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

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

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.]]]])
```

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

✓ Discussions & Exercises

7.5.1. Maximum Pooling and Average Pooling

The idea of *Average pooling* is akin to downsampling an image. Rather than just taking the value of every second (or third) pixel for the lower resolution image, we can average over adjacent pixels to obtain an image with better signal-to-noise ratio since we are combining the information from multiple adjacent pixels. *Max-pooling* was introduced in Riesenhuber and Poggio (1999) in the context of cognitive neuroscience to describe how information aggregation might be aggregated hierarchically for the purpose of object recognition. In almost all cases, max-pooling is preferable to average pooling.

7.5.3. Multiple Channels

When processing multi-channel input data, the pooling layer pools each input channel separately,

rather than summing the inputs up over channels as in a convolutional layer. This means that the number of output channels for the pooling layer is the same as the number of input channels.

✓ 7.6. Convolutional Neural Networks (LeNet)

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

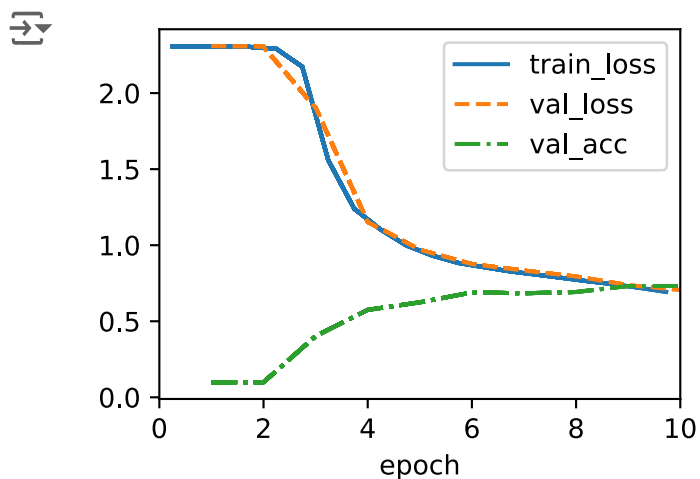
class LeNet(d2l.Classifier):
    def __init__(self, lr=0.1, num_classes=10):
        super().__init__()
        self.save_hyperparameters()
        self.net = nn.Sequential(
            nn.LazyConv2d(6, kernel_size=5, padding=2), nn.Sigmoid(),
            nn.AvgPool2d(kernel_size=2, stride=2),
     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

✓ Discussions & Exercises

7.6.1. LeNet

At a high level, LeNet (LeNet-5) consists of two parts:

- (i) a convolutional encoder consisting of two convolutional layers
- (ii) a dense block consisting of three fully connected layers

✓ 8. Modern Convolutional Neural Networks

✓ 8.2. Networks Using Blocks (VGG)

```
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)
```

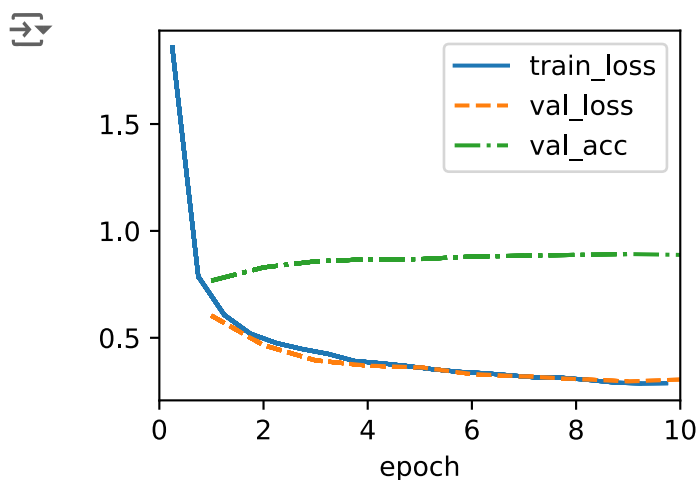
```
class VGG(d2l.Classifier):
    def __init__(self, arch, lr=0.1, num_classes=10):
        super().__init__()
        self.save_hyperparameters()
        conv_blks = []
        for (num_convs, out_channels) in arch:
            conv_blks.append(vgg_block(num_convs, out_channels))
        self.net =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

```
# @title
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)
```



✓ Discussions & Exercises

The design of neural network architectures has grown progressively more abstract, from thinking in terms of individual neurons to whole layers, and now to blocks, repeating patterns of layers.

The idea of using blocks first emerged from the VGG, which is easy to implement by using loops and subroutines. This has now progressed to researchers using entire trained models to repurpose them for different, albeit related, tasks, typically called *foundation models*.

8.2.1. VGG Blocks

The basic building block of CNNs is a sequence of the following:

- (i) a convolutional layer with padding to maintain the resolution
- (ii) a nonlinearity such as a ReLU
- (iii) a pooling layer such as max-pooling to reduce the resolution

The key idea of Simonyan and Zisserman (2014) was to use multiple convolutions in between downsampling via max-pooling in the form of a block. They showed that deep and narrow networks significantly outperform their shallow counterparts. This set deep learning on a quest for ever deeper networks with over 100 layers for typical applications. Stacking 3×3 convolutions has become a gold standard in later deep networks.

8.2.2. VGG Network

Like AlexNet and LeNet, the VGG Network can be partitioned into two parts:

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

The key difference is that the convolutional layers are grouped in nonlinear transformations that leave the dimensionality unchanged, followed by a resolution-reduction step.

✓ 8.6. Residual Networks (ResNet) and ResNeXt

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

```
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])
```

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

```
blk(X).shape
```

```
⇒ torch.Size([4, 6, 3, 3])
```

```
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)
```

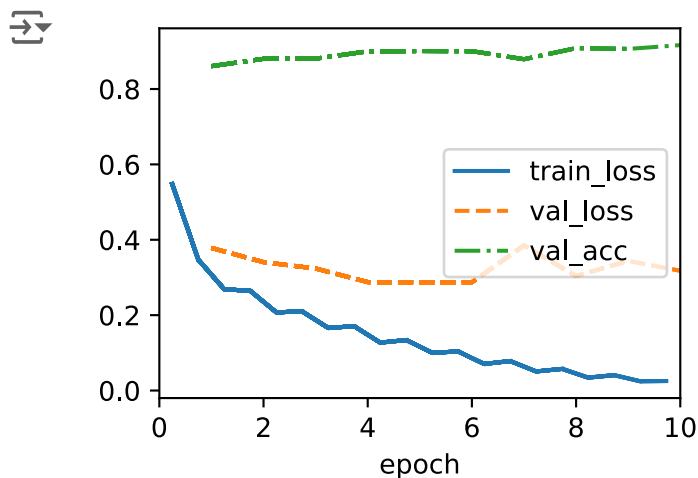
```
ResNet18().layer_summary((1, 1, 96, 96))
```

```
⇒ 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)

```



✓ Discussions & Exercises

8.6.1. Function Classes

Only if larger function classes contain the smaller ones are we guaranteed that increasing them strictly increases the expressive power of the network. For deep neural networks, if we can train the newly-added layer into an identity function $f(x) = x$, the new model will be as effective as the original model. As the new model may get a better solution to fit the training dataset, the added layer might make it easier to reduce training errors.

Residual network (ResNet) is the idea that every additional layer should more easily contain the identity function as one of its elements. These considerations are rather profound but they led to a surprisingly simple solution, a *residual block*.

8.6.2. Residual Blocks

With residual blocks, inputs can forward propagate faster through the residual connections across layers.

