

✓ COSE474-2024F: DEEP LEARNING HW2

✓ 7. Convolutional Neural Networks

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

 [Show hidden output](#)


✓ 7.2. Convolutions for Images

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

✓ 7.2.1. The Cross-Correlation Operation

```
def corr2d(X, K):
    """Compute 2D cross-correlation."""
    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.]])


✓ 7.2.2. 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
```


✓ 7.2.3. 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.]])
```

✓ 7.2.4. Learning a Kernel

```
conv2d = nn.LazyConv2d(1, kernel_size=(1, 2), bias=False)
```

```
# (example, channel, height, width)
```

```
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()
    # 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 4.519
   epoch 4, loss 1.377
   epoch 6, loss 0.485
   epoch 8, loss 0.185
   epoch 10, loss 0.074
```

```
# kernel tensor we learned
```

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

```
↵ tensor([[ 0.9650, -1.0204]])
```

✓ 7.3. Padding and Stride

✓ 7.3.1. Padding

```
def comp_conv2d(conv2d, X):
    # (1, 1) => batch size, num channels = 1
    X = X.reshape((1, 1) + X.shape)
    Y = conv2d(X)
    # Strip examples and channels
    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])
```

```
# different size padding for height and width
```

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

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

```
↵ torch.Size([8, 8])
```

✓ 7.3.2. 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

✓ 7.4.1. Multiple Input Channels

```
# cross-correlation operation per channel and then adding up the results
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.]])
```

✓ 7.4.2. Multiple Output Channels

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

```
K.shape
```

```
↗ torch.Size([2, 2, 2])
```

```
# 3 times [2, 2, 2]
K = torch.stack((K, K + 1, K + 2), 0)
K.shape
```

```
↗ torch.Size([3, 2, 2, 2])
```

```
# get output for each kernel
corr2d_multi_in_out(X, K)
```

```
↗ tensor([[[[ 56.,  72.],
              [104., 120.]],
           [[ 76., 100.],
              [148., 172.]],
           [[ 96., 128.],
              [192., 224.]]]])
```

✓ 7.4.3. 1x1 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

✓ 7.5.1. 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.]])
```

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

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

✓ 7.5.2. 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(2)
pool2d(X)
```

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

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

✓ 7.5.3. 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.]]]]])
```

```
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)

7.6.1. LeNet

```

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:\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])

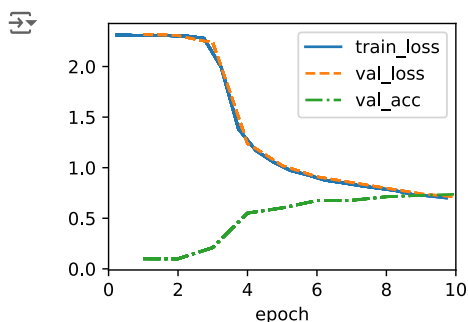
```

7.6.2. 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)

```



✓ 7. Discussions & Exercises

7.1. Disadvantages of Fully Connected Layers (Discussion)

- too many parameters (infeasible to train)
- do not assume any structure regarding how the features interact (losing spatial information)

✓ 7.1 CNN (Discussion)

1. **translation invariance** -> similar response to the same patch in the earliest layers, regardless of it's location in the image.
2. **locality principle** -> the earliest layers focus on local regions
3. deeper layers capture longer-range features of the image ("can see the bigger picture")

Reducing the Number of parameters using convolution Layers

MLP with two-dimensional images:

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

- requires 10^{12} parameters

Translation Invariance (Convolution Formula):

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

- requires 4×10^6 parameters

Locality (Convolutional Layer):

- outside of range given by Δ we set $[\mathbf{V}]_{a,b} = 0$

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

- requires $4\Delta^2$ parameters (Δ is typically smaller than 10)

Support for multiple channels (feature maps):

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

Channels:

- channels allow the CNN to reason with multiple features, such as edge and shape detectors at the same time
- the time complexity of computing $k \times k$ convolution given an image of size $h \times w$ with input and output channels is:

$$\mathcal{O}(c_{\text{in}} \cdot c_{\text{out}} \cdot k^2 \cdot h \cdot w)$$
- for image of size 256×256 and 5×5 kernel and 128 in and out channels, it equals to over 53 billion operations

✓ 7.3.4. Exercises

Given the final code example in this section with kernel size (3,5), padding (0,1), and stride (3,4), calculate the output shape to check if it is consistent with the experimental result.

The formulas to calculate the output size of height h_{out} and width w_{out} are:

$$h_{\text{out}} = \left\lfloor \frac{h_{\text{in}} - k_h + p_h + s_h}{s_h} \right\rfloor$$

$$w_{\text{out}} = \left\lfloor \frac{w_{\text{in}} - k_w + p_w + s_w}{s_w} \right\rfloor$$

If we substitute with the given numbers we get

$$h_{\text{out}} = \left\lfloor \frac{8 - 3 + 0 + 3}{3} \right\rfloor = 2$$

$$w_{\text{out}} = \left\lfloor \frac{8 - 5 + 1 + 4}{4} \right\rfloor = 2$$

The result is consistent with the experimental result

✓ 7.4.1 Cross-correlation computation with multiple input and output channels (Discussion)

- perform a cross-correlation (convolution) operation per channel and then **add up the results**
- the shape of the convolution kernel is $c_i \times k_h \times k_w$
- if we add multiple output channels, we need to add kernel for each output. The shape of the result kernel is $c_o \times c_i \times k_h \times k_w$

✓ 7.4.3 1x1 Convolution Layer (Discussion)

- a computation of the 1x1 convolution occurs on the **channel dimension**
- the convolution is followed by nonlinearity
- requires $c_o \times c_i$ weights

✓ 7.6.1. LeNet (Discussion & Exercise)

- among the first published CNNs to capture wide attention for its performance on computer vision tasks
- introduced by Yann LeCun a researcher at AT&T Bell Labs, for the purpose of recognizing handwritten digits in images (LeCun et al., 1998)
 - eventually adapted to recognize digits for processing deposits in ATM machines
- consists of:
 - 2 convolution layers (5x5 convolution with padding 2, 2x2 AvgPool with stride 2)
 - 3 fully connected layers
- originally it used sigmoid activation function and average pooling (ReLU and max-pooling work better)

Modernizing LeNet: (Exercise)

- to modernize LeNet the sigmoid activation function and AvgPool can be replaced by **ReLU and MaxPool**

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

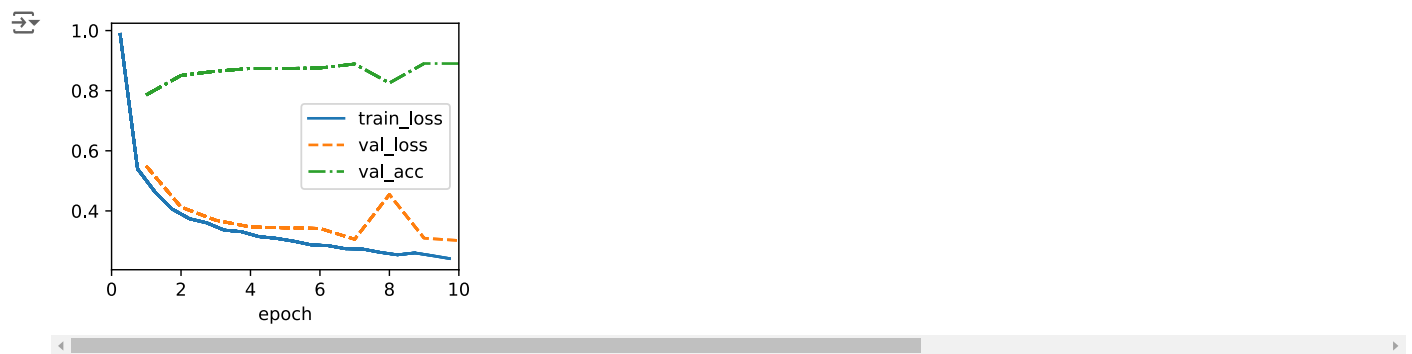
class LeNet(d2l.Classifier):
    def __init__(self, lr=0.1, num_classes=10):
        super().__init__()
        self.save_hyperparameters()
        self.net = nn.Sequential(
            nn.LazyConv2d(6, kernel_size=5, padding=2), nn.ReLU(),
            nn.MaxPool2d(kernel_size=2, stride=2),
            nn.LazyConv2d(16, kernel_size=5), nn.ReLU(),
            nn.MaxPool2d(kernel_size=2, stride=2),
            nn.Flatten(),
            nn.LazyLinear(120), nn.ReLU(),
            nn.LazyLinear(84), nn.ReLU(),
            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])
ReLU output shape: torch.Size([1, 6, 28, 28])
MaxPool2d output shape: torch.Size([1, 6, 14, 14])
Conv2d output shape: torch.Size([1, 16, 10, 10])
ReLU output shape: torch.Size([1, 16, 10, 10])
MaxPool2d output shape: torch.Size([1, 16, 5, 5])
Flatten output shape: torch.Size([1, 400])
Linear output shape: torch.Size([1, 120])
ReLU output shape: torch.Size([1, 120])
Linear output shape: torch.Size([1, 84])
ReLU 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)
```



✓ COSE474-2024F: DEEP LEARNING HW2

✓ 8. Modern Convolutional Neural Networks

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

 [Show hidden output](#)

✓ 8.2. Networks Using Blocks (VGG)

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


✓ 8.2.1. VGG Blocks

```
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)
```

✓ 8.2.2. VGG Network

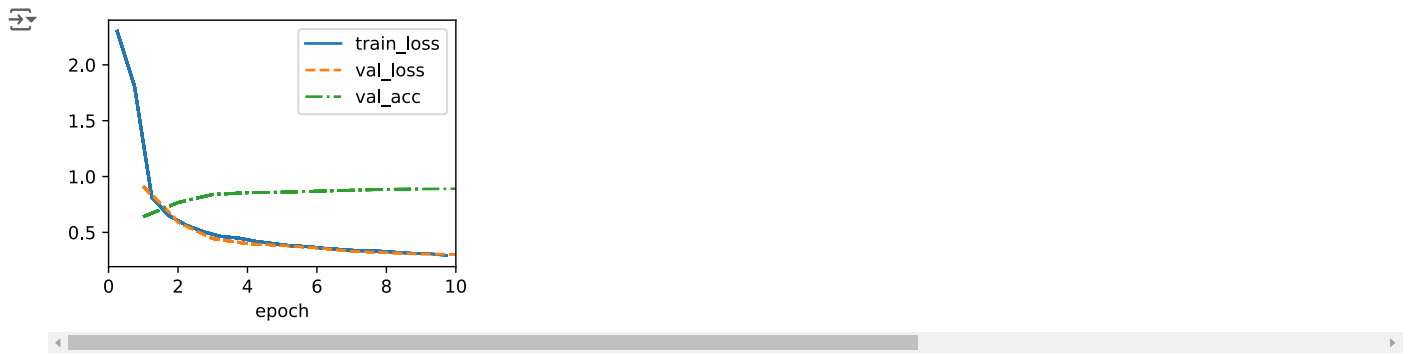
```
class VGG(d2l.Classifier):
    def __init__(self, arch, lr=0.1, num_classes=10):
        super().__init__()
        self.save_hyperparameters()
        conv_blks = []
        for (num_convs, out_channels) in arch:
            conv_blks.append(vgg_block(num_convs, out_channels))
        self.net =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

✓ 8.2.3. 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

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

8.6.1. Function Classes (Discussion)

8.6.2. Residual Blocks

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

8.6.3. 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.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.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])

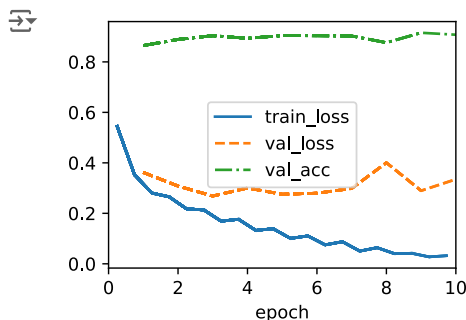
```

8.6.4. 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_data_loader(True)))[0]], d2l.init_cnn)
trainer.fit(model, data)

```



8. Discussions & Exercises

8.2.2. VGG Network (Discussion)

- first truly **modern CNN** developed by the Visual Geometry Group (VGG) at Oxford University
- introduces **VGG blocks** - building blocks of the network
- the VGG block consists of:
 1. convolutional layer with a **sequence of 3x3 convolutions** and padding (maintains the resolution)
 2. each convolution is followed by a nonlinearity (ReLU)
 3. a pooling layer (2x2 max-pooling with stride 2)
- 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.

8.6.3. ResNet Model (Discussion)

- ResNet developed by Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun (2016) from Microsoft Research.
- introduces **Residual blocks**

- instead of directly learning $y = f(x)$ the residual block needs to learn $y = f(x) + x$ or the residual mapping $g(x) = f(x) - x$. The function $f(x)$ represents the residual function that models **the difference between the input and the output** rather than the entire transformation.
- The residual block includes **residual connection** (or shortcut connection) that carries the input x to an addition at the end of the block before the last activation function.
- The ResNet network consists of residual blocks (3x3 convolution -> batch norm -> ReLU -> 3x3 convolution -> batch norm -> residual connection -> ReLU)
- if we need to downsample the images we can use **convolution with stride**