## ∨ LIBRARIES

```
pip install torch==2.0.0 torchvision==0.15.1
```

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

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

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

## ∨ 7. CNN

### ∨ 7.1 From Fully Connected Layers to Convolutions

Key takeaways:

1. Our goal is to minimize number of parameters needed without limiting its expressive power or dimensions.

2. Channels could bring the lost of complexity due to locality and translation invariance.

### ∨ 7.2 Convolutions for Images

```
#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.]])
```

```
#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 13.189
epoch 4, loss 4.461
epoch 6, loss 1.669
epoch 8, loss 0.657
epoch 10, loss 0.265
/usr/local/lib/python3.10/dist-packages/torch/nn/modules/lazy.py:180: UserWarning: Lazy modules are a new feature under heavy developmen
  warnings.warn('Lazy modules are a new feature under heavy development '
```

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

```
tensor([[ 0.9369, -1.0426]])
```

Key takeaways:

1. Input tensor and a kernel tensor are combined to produce an output tensor through a cross-correlation operation in convolutional layers.

2. Kernels will be initialized randonmly.

3. Horizontally adjacent elements are the same, the output is 0. Otherwise, the output is nonzero in cross-correlation operation with the inputs.

4. When Convolutional layer performs strict cross-correlation, this layer performs strict convolution where the learned kernel K' will be the same as K after K' is flipped both horizontally and vertically.

5. Convolutional layer output is sometimes called a feature map, as it can be regarded as the learned representations (features) in the spatial dimensions.

## ⌄ 7.3 Padding and Stride

```
#Padding

def comp_conv2d(conv2d, X):
    # (1, 1) indicates that batch size and the number of channels are both 1
    X = X.reshape((1, 1) + X.shape)
    Y = conv2d(X)
    # Strip the first two dimensions: 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])

```
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])

Key takeaways:

1. We can decrease the number of unuse pixels of inputs by following padding techniques which by adding up few extra dummy pixels at the boundary that will be set to 0.

2. Choosing odd kernel sizes has the benefit that we can preserve the dimensionality while padding with the same number of rows on top and bottom, and the same number of columns on left and right.

3. The stride can reduce the resolution of the output, for example reducing the height and width of the output to only 1/n of the height and width of the input for n > 1.


## ⌄ 7.4 Multiple Input and Multiple Output Channels

```
#Multiple Channel inputs

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 Channel outputs

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.]]])
```

```
#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))
    # Matrix multiplication in the fully connected layer
    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
```

Key takeaways:

1. The depth (number of channels) and width (number of filters) in layers significantly affect the feature extraction capability.

2. Channels allow CNNs to process multiple types of data simultaneously, enhancing the model's ability to learn complex patterns.

## ˅ 7.5 Pooling

```
#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.]])
```

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

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

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

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

Key takeaways:

1. Unlike CNN, pooling layer has no parameters (kernel) for sliding. And calculating maximum (Maximum pooling) or average (Average pooling) value of the elements in the pooling window.

2. Average pooling similar to old CNN, to downsampling an image but without kernel.

3. Has the similar concept with CNN, Padding and Stride.

4. Deep learning by default automatically match the stride value according to desired pooling window size. However, it can still be overridden.

## 7.6 Convolutional Neural Networks (LeNet)

```
#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])
/usr/local/lib/python3.10/dist-packages/torch/nn/modules/lazy.py:180: UserWarning: Lazy modules are a new feature under heavy developmen
  warnings.warn('Lazy modules are a new feature under heavy development '
```
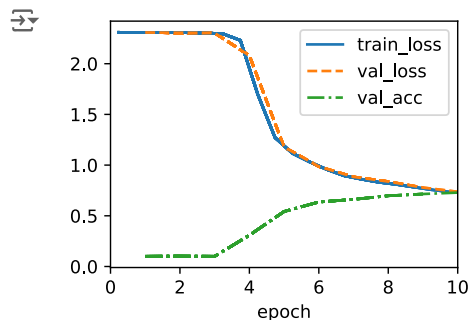
```
#Training

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)
```



Key takeaways:

1. LeNet, the first published CNN.
2. Consists 2 parts: Convolutional encoder (2 convolutional layers) and Dense block (3 Fully connected layers).
3. In order to pass output from the convolutional block to the dense block, we must flatten each example in the minibatch.
4. In LeNet, Gaussian activation layer is replaced by a softmax layer

## 8. Modern Convolutional Neural Networks

## 8.2. Networks Using Blocks (VGG)

```
#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)
```

```
#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))
```
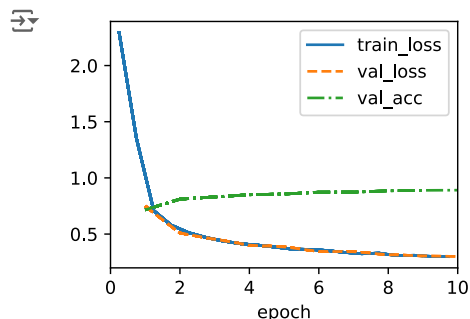
```
Sequential output shape:         torch.Size([1, 64, 112, 112])
Sequential output shape:         torch.Size([1, 128, 56, 56])
Sequential output shape:         torch.Size([1, 256, 28, 28])
Sequential output shape:         torch.Size([1, 512, 14, 14])
Sequential output shape:         torch.Size([1, 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])
```

```
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)
```



Key takeaways:

1. Traditional CNN has disadvantage which has a rapid decrement in spatial resolution.
2. Use multiple convolutions in between downsampling via max-pooling in the form of a block.
3. VGG block consists of a sequence of convolutions with kernels with padding of 1 (keeping height and width) followed by a max-pooling layer with stride of 2 (halving height and width after each block).
4. It calls2 argument, Number of convolutional layers and Number of output channels.

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

```python
#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)
```

```python
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
#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))
```

```python
@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)
```

```python
@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)
```

```python
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])
```
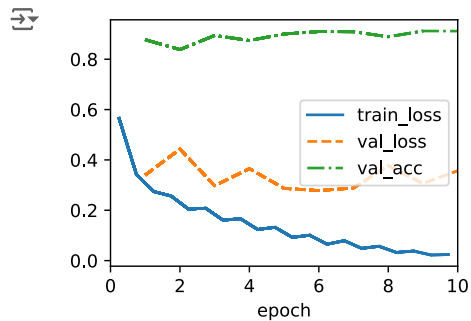
```
#Training

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)
```

Key takeaways:

1. Larger size of training data generally leads to better f*(F)
2. Assumption that if there is more powerful architecture of F', equals to better outcome.
3. Theory shows that larger function (Non-nested function) class not necessarily always move closer to truth function f*.
4. Nested function classes can avoid issue exists in non-nested function classes.
5. With residual blocks, inputs can forward propagate faster through the residual connections across layers.
6. ResNet has VGG's full 3x3 convolutional layer design. (Two 3x3 convolutional layers with the same number output channels + ReLU).
7. Output of two convolutional layers must be the same shape to be added together.