

CNN

- Useful Models for Images, Videos and Texts

CNN

- 卷积神经网络（convolutional neural network）是含有卷积层（convolutional layer）的神经网络。本节中介绍的卷积神经网络均使用最常见的二维卷积层。它有高和宽两个空间维度，常用来处理图像数据。我们将介绍简单形式的二维卷积层的工作原理。

二维互相关运算

- 虽然卷积层得名于卷积（convolution）运算，但我们通常在卷积层中使用更加直观的互相关（cross-correlation）运算。在二维卷积层中，一个二维输入数组和一个二维核（kernel）数组通过互相关运算输出一个二维数组。
- 我们用一个具体例子来解释二维互相关运算的含义。如图5.1所示，输入是一个高和宽均为3的二维数组。我们将该数组的形状记为 3×3 或 $(3, 3)$ 。核数组的高和宽分别为2。该数组在卷积计算中又称卷积核或过滤器（filter）。卷积核窗口（又称卷积窗口）的形状取决于卷积核的高和宽，即 2×2 。图5.1中的阴影部分为第一个输出元素及其计算所使用的输入和核数组元素： $0 \times 0 + 1 \times 1 + 3 \times 2 + 4 \times 3 = 19$ 。

0	1	2
3	4	5
6	7	8

*

0	1
2	3

=

19	25
37	43

二维互相关运算

- 在二维互相关运算中，卷积窗口从输入数组的最左上方开始，按从左往右、从上往下的顺序，依次在输入数组上滑动。当卷积窗口滑动到某一位置时，窗口中的输入子数组与核数组按元素相乘并求和，得到输出数组中相应位置的元素。图5.1中的输出数组高和宽分别为2，其中的4个元素由二维互相关运算得出：
- $0 \times 0 + 1 \times 1 + 3 \times 2 + 4 \times 3 = 19$,
- $1 \times 0 + 2 \times 1 + 4 \times 2 + 5 \times 3 = 25$,
- $3 \times 0 + 4 \times 1 + 6 \times 2 + 7 \times 3 = 37$,
- $4 \times 0 + 5 \times 1 + 7 \times 2 + 8 \times 3 = 43$.

0	1	2
3	4	5
6	7	8

 *

0	1
2	3

 =

19	25
37	43

二维互相关运算

- 下面我们将上述过程实现在corr2d函数里。它接受输入数组X与核数组K，并输出数组Y。

```
import torch
```

```
from torch import nn
```

```
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
```

二维互相关运算

```
X = torch.tensor([[0, 1, 2], [3, 4, 5], [6, 7, 8]])
```

```
K = torch.tensor([[0, 1], [2, 3]])
```

```
corr2d(X, K)
```

- 输出：

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

二维卷积层

二维卷积层将输入和卷积核做互相关运算，并加上一个标量偏差来得到输出。下面基于corr2d函数来实现一个自定义的二维卷积层。在构造函数__init__里我们声明weight和bias这两个模型参数。前向计算函数forward则是直接调用corr2d函数再加上偏差。

```
class Conv2D(nn.Module):  
    def __init__(self, kernel_size):  
        super(Conv2D, self).__init__()  
        self.weight = nn.Parameter(torch.randn(kernel_size))  
        self.bias = nn.Parameter(torch.randn(1))  
  
    def forward(self, x):  
        return corr2d(x, self.weight) + self.bias
```

特征图和感受野

0	1	2
3	4	5
6	7	8

 \star

0	1
2	3

 $=$

19	25
37	43

二维卷积层输出的二维数组可以看作是输入在空间维度（宽和高）上某一级的表征，也叫特征图（feature map）。影响元素x的前向计算的所有可能输入区域（可能大于输入的实际尺寸）叫做x的感受野（receptive field）。

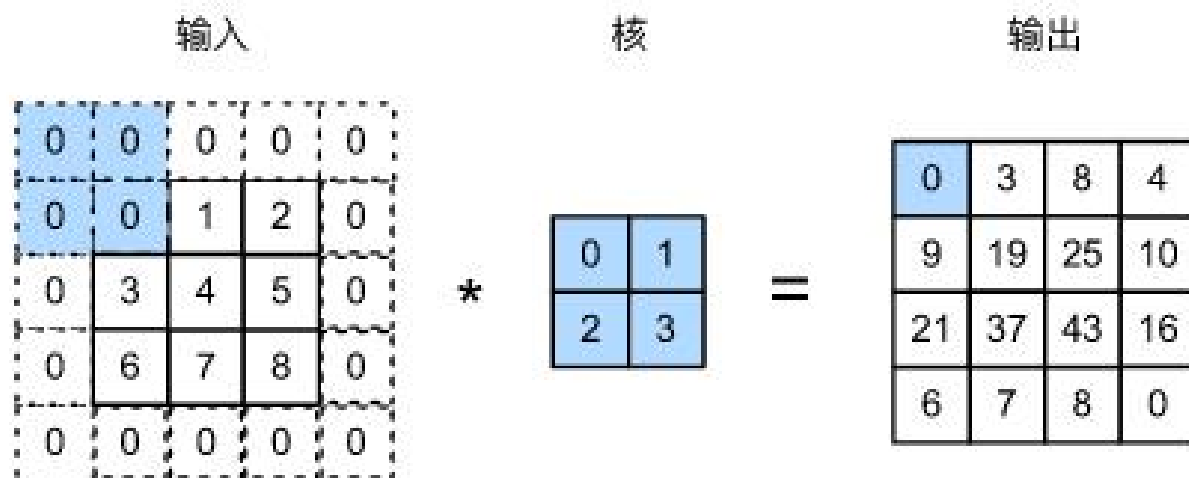
以图为例，输入中阴影部分的四个元素是输出中阴影部分元素的感受野。我们将图中形状为 2×2 的输出记为Y，并考虑一个更深的卷积神经网络：将Y与另一个形状为 2×2 的核数组做互相关运算，输出单个元素z。那么，z在Y上的感受野包括Y的全部四个元素，在输入上的感受野包括其中全部9个元素。

填充和步幅

- 我们使用高和宽为3的输入与高和宽为2的卷积核得到高和宽为2的输出。一般来说，假设输入形状是 $n_h \times n_w$ ，卷积核窗口形状是 $k_h \times k_w$ ，那么输出形状将会是
- $(n_h - k_h + 1) \times (n_w - k_w + 1)$.
- 所以卷积层的输出形状由输入形状和卷积核窗口形状决定。卷积层的两个超参数，即填充和步幅，可以对给定形状的输入和卷积核改变输出形状。

填充

- 填充（padding）是指在输入高和宽的两侧填充元素（通常是0元素）。我们在原输入高和宽的两侧分别添加了值为0的元素，使得输入高和宽从3变成了5，并导致输出高和宽由2增加到4。图中的阴影部分为第一个输出元素及其计算所使用的输入和核数组元素： $0 \times 0 + 0 \times 1 + 0 \times 2 + 0 \times 3 = 0$ 。



填充

- 如果在高的两侧一共填充 p_h 行，在宽的两侧一共填充 p_w 列，那么输出形状将会是

$$(n_h - k_h + p_h + 1) \times (n_w - k_w + p_w + 1),$$

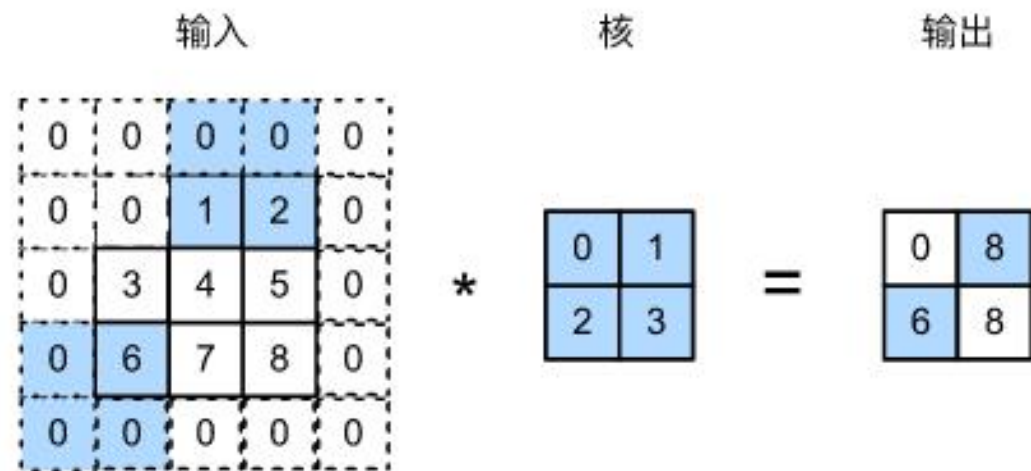
- 也就是说，输出的高和宽会分别增加 p_h 和 p_w 。
- 在很多情况下，我们会设置 $p_h = k_h - 1$ 和 $p_w = k_w - 1$ 来使输入和输出具有相同的高和宽。这样会方便在构造网络时推测每个层的输出形状。假设这里 k_h 是奇数，我们会在高的两侧分别填充 $p_h/2$ 行。
- 卷积神经网络经常使用奇数高宽的卷积核，如1、3、5和7，所以两端上的填充个数相等。对任意的二维数组 X ，设它的第 i 行第 j 列的元素为 $X[i,j]$ 。当输入和输出具有相同的高和宽时，输出 $Y[i,j]$ 是由输入以 $X[i,j]$ 为中心的窗口同卷积核进行互相关计算得到的。

步幅

- 我们介绍二维互相关运算。卷积窗口从输入数组的最左上方开始，按从左往右、从上往下的顺序，依次在输入数组上滑动。我们将每次滑动的行数和列数称为步幅（stride）。
- 目前我们看到的例子里，在高和宽两个方向上步幅均为1。我们也可以使用更大步幅。下图展示了在高上步幅为3、在宽上步幅为2的二维互相关运算。图中的阴影部分为输出元素及其计算所使用的输入和核数组元素：

$$0 \times 0 + 0 \times 1 + 1 \times 2 + 2 \times 3 = 8、$$

$$0 \times 0 + 6 \times 1 + 0 \times 2 + 0 \times 3 = 6。$$



步幅

- 一般来说，当高上步幅为 s_h ，宽上步幅为 s_w 时，输出形状为 $(n_h - k_h + p_h + s_h) / s_h \times (n_w - k_w + p_w + s_w) / s_w$
- 如果设置 $p_h = k_h - 1$ 和 $p_w = k_w - 1$ ，并且输入的高和宽能分别被高和宽上的步幅整除，那么输出形状将是 $(n_h / s_h) \times (n_w / s_w)$ 。
- 填充可以增加输出的高和宽。这常用来使输出与输入具有相同的高和宽。
- 步幅可以减小输出的高和宽，例如输出的高和宽仅为输入的高和宽的 $1/n$ （ n 为大于1的整数）。

多输入通道和多输出通道

- 前面两节里我们用到的输入和输出都是二维数组，但真实数据的维度经常更高。例如，彩色图像在高和宽2个维度外还有RGB（红、绿、蓝）3个颜色通道。假设彩色图像的高和宽分别是 h 和 w （像素），那么它可以表示为一个 $3 \times h \times w$ 的多维数组。我们将大小为3的这一维称为通道（channel）维。我们将介绍含多个输入通道或多个输出通道的卷积核。

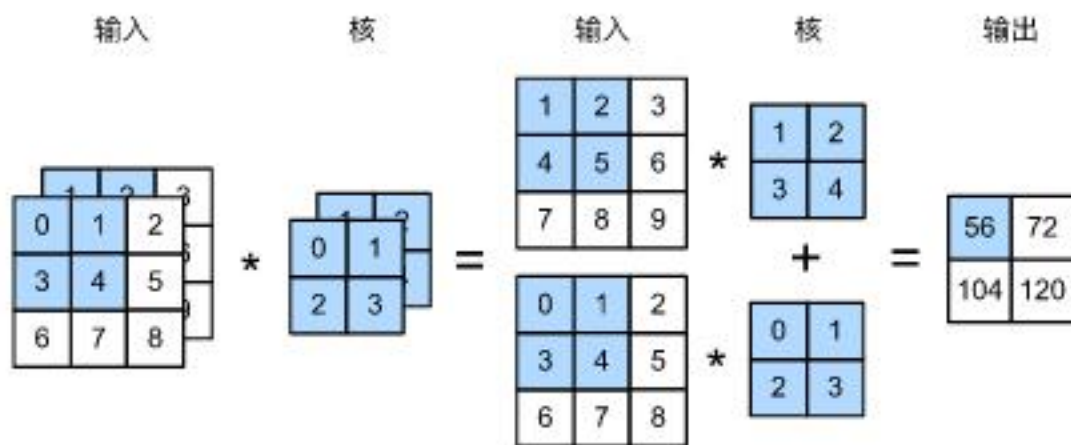
多输入通道

- 当输入数据含多个通道时，我们需要构造一个输入通道数与输入数据的通道数相同的卷积核，从而能够与含多通道的输入数据做互相关运算。假设输入数据的通道数为 c_i ，那么卷积核的输入通道数同样为 c_i 。设卷积核窗口形状为 $k_h \times k_w$ 。当 $c_i=1$ 时，我们知道卷积核只包含一个形状为 $k_h \times k_w$ 的二维数组。当 $c_i > 1$ 时，我们将会为每个输入通道各分配一个形状为 $k_h \times k_w$ 的核数组。把这 c_i 个数组在输入通道维上连结，即得到一个形状为 $c_i \times k_h \times k_w$ 的卷积核。由于输入和卷积核各有 c_i 个通道，我们可以在各个通道上对输入的二维数组和卷积核的二维核数组做互相关运算，再将这 c_i 个互相关运算的二维输出按通道相加，得到一个二维数组。这就是含多个通道的输入数据与多输入通道的卷积核做二维互相关运算的输出。

多输入通道

- 下图展示了含2个输入通道的二维互相关计算的例子。在每个通道上，二维输入数组与二维核数组做互相关运算，再按通道相加即得到输出。图中阴影部分为第一个输出元素及其计算所使用的输入和核数组元素：

$$(1 \times 1 + 2 \times 2 + 4 \times 3 + 5 \times 4) + (0 \times 0 + 1 \times 1 + 3 \times 2 + 4 \times 3) = 56$$



多输出通道

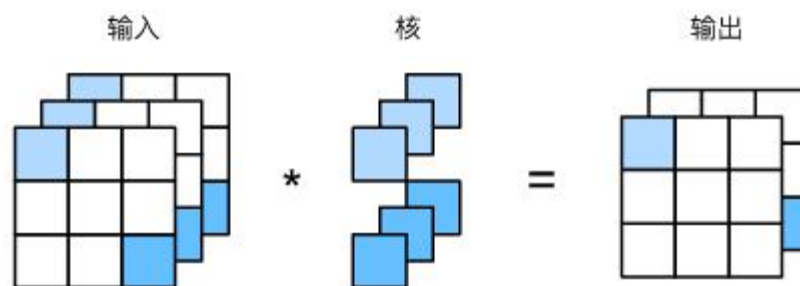
- 当输入通道有多个时，因为我们对各个通道的结果做了累加，所以不论输入通道数是多少，输出通道数总是为1。设卷积核输入通道数和输出通道数分别为 c_i 和 c_o ，高和宽分别为 k_h 和 k_w 。如果希望得到含多个通道的输出，我们可以为每个输出通道分别创建形状为 $c_i \times k_h \times k_w$ 的核数组。将它们在输出通道维上连结，卷积核的形状即 $c_o \times c_i \times k_h \times k_w$ 。在做互相关运算时，每个输出通道上的结果由卷积核在该输出通道上的核数组与整个输入数组计算而来。

1 x 1卷积层

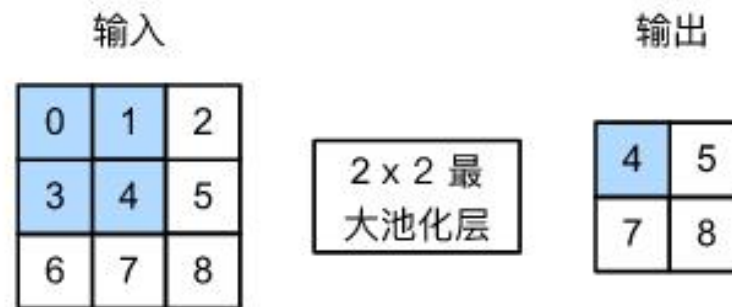
- 最后我们讨论卷积窗口形状为1 x 1 ($k_h=k_w=1$) 的多通道卷积层。我们通常称之为1 x 1卷积层，并将其中的卷积运算称为1 x 1卷积。因为使用了最小窗口，1 x 1卷积失去了卷积层可以识别高和宽维度上相邻元素构成的模式的功能。实际上，1 x 1卷积的主要计算发生在通道维上。

1 x 1卷积层

- 下图展示了使用输入通道数为3、输出通道数为2的1 x 1卷积核的互相关计算。值得注意的是，输入和输出具有相同的高和宽。输出中的每个元素来自输入中在高和宽上相同位置的元素在不同通道之间的按权重累加。假设我们将通道维当作特征维，将高和宽维度上的元素当成数据样本，那么1 x 1卷积层的作用与全连接层等价。



池化层



- 卷积核可以精确定位像素变化的位置。但实际图像里，我们感兴趣的物体不会总出现在固定位置：这会导致我们感兴趣的输出可能出现在卷积输出Y中的不同位置，进而对后面的模式识别造成不便。池化（pooling）层的提出是为了缓解卷积层对位置的过度敏感性。
- 同卷积层一样，池化层每次对输入数据的一个固定形状窗口（又称池化窗口）中的元素计算输出。不同于卷积层里计算输入和核的互相关性，池化层直接计算池化窗口内元素的最大值或者平均值。该运算也分别叫做最大池化或平均池化。在二维最大池化中，池化窗口从输入数组的最左上方开始，按从左往右、从上往下的顺序，依次在输入数组上滑动。当池化窗口滑动到某一位置时，窗口中的输入子数组的最大值即输出数组中相应位置的元素。

池化层

- 最大池化和平均池化分别取池化窗口中输入元素的最大值和平均值作为输出。
- 池化层的一个主要作用是缓解卷积层对位置的过度敏感性。
- 可以指定池化层的填充和步幅。
- 池化层的输出通道数跟输入通道数相同。

批量归一化

- 通常来说，数据标准化预处理对于浅层模型就足够有效了。随着模型训练的进行，当每层中参数更新时，靠近输出层的输出较难出现剧烈变化。但对深层神经网络来说，即使输入数据已做标准化，训练中模型参数的更新依然很容易造成靠近输出层输出的剧烈变化。这种计算数值的不稳定性通常令我们难以训练出有效的深度模型。
- 批量归一化的提出正是为了应对深度模型训练的挑战。在模型训练时，批量归一化利用小批量上的均值和标准差，不断调整神经网络中间输出，从而使整个神经网络在各层的中间输出的数值更稳定。

对全连接层做批量归一化

- 通常，我们将批量归一化层置于全连接层中的仿射变换和激活函数之间。设全连接层的输入为 \mathbf{u} ，权重参数和偏差参数分别为 \mathbf{W} 和 \mathbf{b} ，激活函数为 ϕ 。设批量归一化的运算符为 BN。那么，使用批量归一化的全连接层的输出为 $\phi(\text{BN}(\mathbf{x}))$ ，其中批量归一化输入 \mathbf{x} 由仿射变换 $\mathbf{x} = \mathbf{W}\mathbf{u} + \mathbf{b}$ 得到。考虑一个由 m 个样本组成的小批量，仿射变换的输出为一个小批量 $B = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ 。它们正是批量归一化层的输入。对于小批量 B 中任意样本 $\mathbf{x}^{(i)} \in \mathbb{R}^d, 1 \leq i \leq m$ ，批量归一化层的输出同样是 d 维向量 $\mathbf{y}^{(i)} = \text{BN}(\mathbf{x}^{(i)})$ ，并由以下步骤求得。首先，对小批量 B 求均值和方差：

$$\boldsymbol{\mu}_B \leftarrow \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)},$$

$$\boldsymbol{\sigma}_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)} - \boldsymbol{\mu}_B)^2,$$

对全连接层做批量归一化

- 其中的平方计算是按元素求平方。接下来，使用按元素开方和按元素除法对 $\mathbf{x}^{(i)}$ 标准化：

$$\hat{\mathbf{x}}^{(i)} \leftarrow \frac{\mathbf{x}^{(i)} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}},$$

- 在上面标准化的基础上，批量归一化层引入了两个可以学习的模型参数，拉伸（scale）参数 γ 和偏移（shift）参数 β 。这两个参数和 $\mathbf{x}^{(i)}$ 形状相同，皆为 d 维向量。它们与 $\hat{\mathbf{x}}^{(i)}$ 分别做按元素乘法（符号 \odot ）和加法计算：

$$\mathbf{y}^{(i)} \leftarrow \gamma \odot \hat{\mathbf{x}}^{(i)} + \beta.$$

- 值得注意的是，可学习的拉伸和偏移参数保留了不做批量归一化的可能：此时只需学出 $\gamma = \sqrt{\sigma_B^2 + \epsilon}$ 和 $\beta = \mu_B$ 。我们可以对此这样理解：如果批量归一化无益，理论上，学出的模型可以不使用批量归一化。

对卷积层做批量归一化

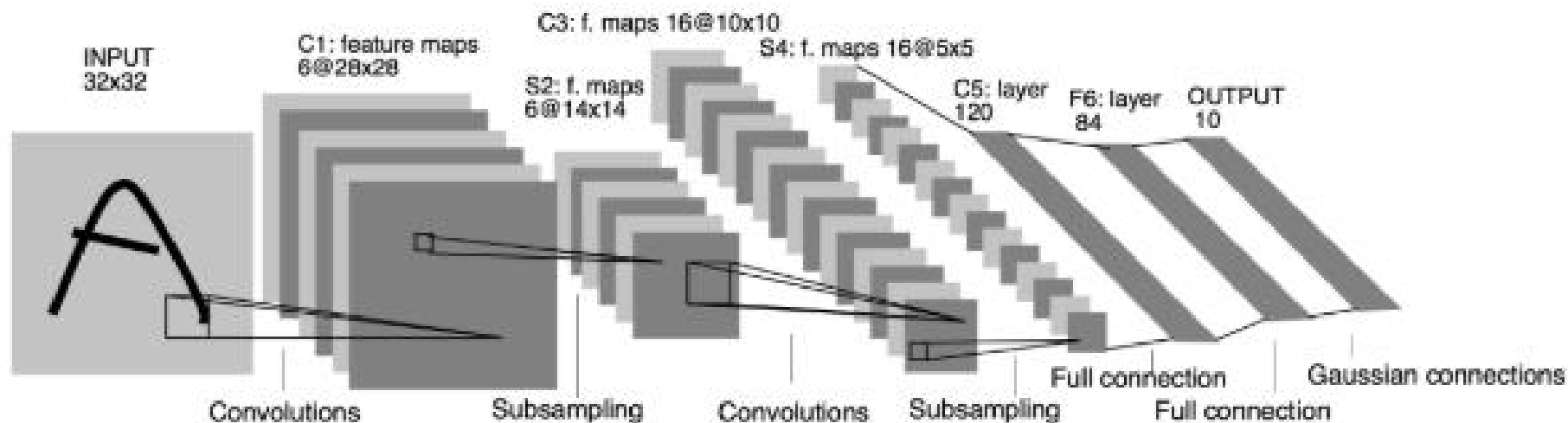
- 对卷积层来说，批量归一化发生在卷积计算之后、应用激活函数之前。如果卷积计算输出多个通道，我们需要对这些通道的输出分别做批量归一化，且**每个通道都拥有独立的拉伸和偏移参数，并均为标量**。设小批量中有 m 个样本。在单个通道上，假设卷积计算输出的高和宽分别为 p, q 。我们需要对该通道中 $m \times p \times q$ 个元素同时做批量归一化。对这些元素做标准化计算时，我们使用相同的均值和方差，即该通道中 $m \times p \times q$ 个元素的均值和方差。

预测时的批量归一化

- 使用批量归一化训练时，我们可以将批量大小设得大一点，从而使批量内样本的均值和方差的计算都较为准确。将训练好的模型用于预测时，我们希望模型对于任意输入都有确定的输出。因此，单个样本的输出不应取决于批量归一化所需要的随机小批量中的均值和方差。一种常用的方法是通过移动平均估算整个训练数据集的样本均值和方差，并在预测时使用它们得到确定的输出。可见，和丢弃层一样，批量归一化层在训练模式和预测模式下的计算结果也是不一样的。

卷积神经网络 (LeNet)

- 接下来介绍一个早期用来识别手写数字图像的卷积神经网络：LeNet。这个名字来源于LeNet论文的第一作者Yann LeCun。LeNet展示了通过梯度下降训练卷积神经网络可以达到手写数字识别在当时最先进的结果。这个奠基性的工作第一次将卷积神经网络推上舞台，为世人所知。LeNet的网络结构如下图所示。



LeNet模型

- LeNet分为卷积层块和全连接层块两个部分。下面我们分别介绍这两个模块。
- 卷积层块里的基本单位是卷积层后接最大池化层：卷积层用来识别图像里的空间模式，如线条和物体局部，之后的最大池化层则用来降低卷积层对位置的敏感性。卷积层块由两个这样的基本单位重复堆叠构成。在卷积层块中，每个卷积层都使用 5×5 的窗口，并在输出上使用sigmoid激活函数。第一个卷积层输出通道数为6，第二个卷积层输出通道数则增加到16。这是因为第二个卷积层比第一个卷积层的输入的高和宽要小，所以增加输出通道使两个卷积层的参数尺寸类似。卷积层块的两个最大池化层的窗口形状均为 2×2 ，且步幅为2。由于池化窗口与步幅形状相同，池化窗口在输入上每次滑动所覆盖的区域互不重叠。

LeNet模型

- 卷积层块的输出形状为(批量大小, 通道, 高, 宽)。当卷积层块的输出传入全连接层块时, 全连接层块会将小批量中每个样本变平 (flatten)。也就是说, 全连接层的输入形状将变成二维, 其中第一维是小批量中的样本, 第二维是每个样本变平后的向量表示, 且向量长度为通道、高和宽的乘积。全连接层块含3个全连接层。它们的输出个数分别是120、84和10, 其中10为输出的类别个数。

```
class LeNet(nn.Module):
    def __init__(self):
        super(LeNet, self).__init__()
        self.conv = nn.Sequential(
            nn.Conv2d(1, 6, 5), # in_channels, out_channels, kernel_size
            nn.Sigmoid(),
            nn.MaxPool2d(2, 2), # kernel_size, stride
            nn.Conv2d(6, 16, 5),
            nn.Sigmoid(),
            nn.MaxPool2d(2, 2)      )
        self.fc = nn.Sequential(
            nn.Linear(16*4*4, 120),
            nn.Sigmoid(),
            nn.Linear(120, 84),
            nn.Sigmoid(),
            nn.Linear(84, 10)      )

    def forward(self, img):
        feature = self.conv(img)
        output = self.fc(feature.view(img.shape[0], -1))
        return output
```

```
net = LeNet()
```

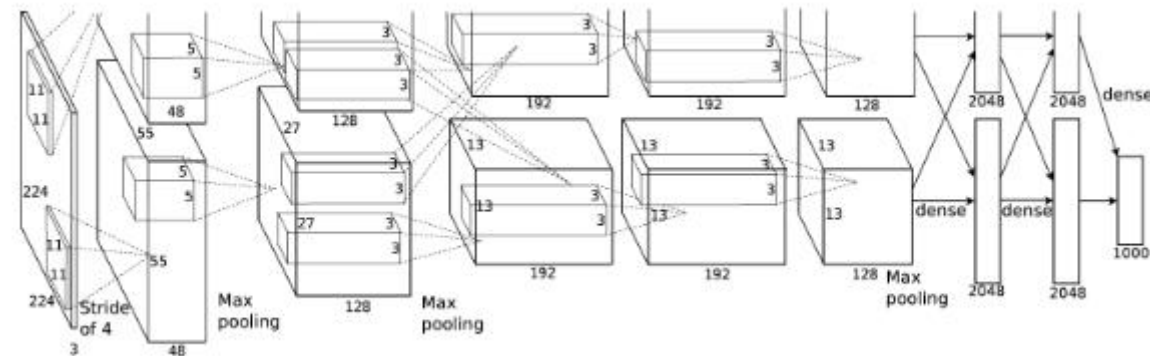
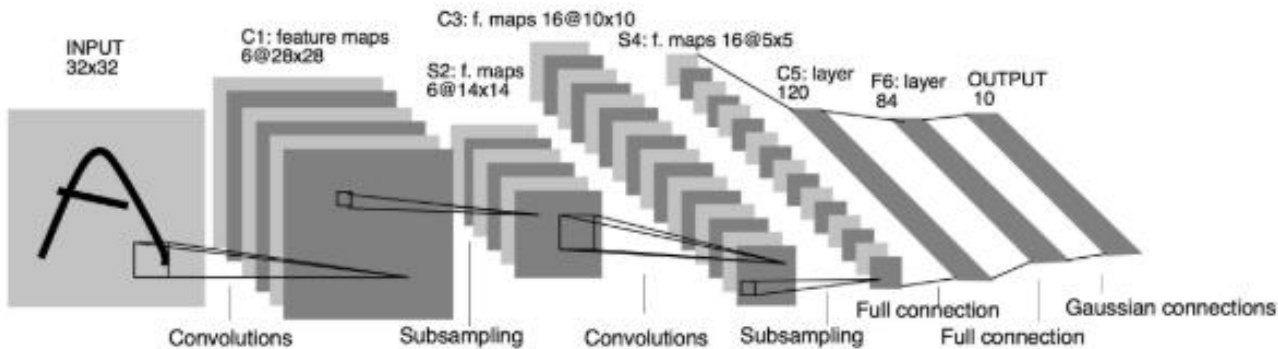
```
print(net)
```

输出：

```
LeNet(  
  (conv): Sequential(  
    (0): Conv2d(1, 6, kernel_size=(5, 5), stride=(1, 1))  
    (1): Sigmoid()  
    (2): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)  
    (3): Conv2d(6, 16, kernel_size=(5, 5), stride=(1, 1))  
    (4): Sigmoid()  
    (5): MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1, ceil_mode=False)  
  )  
  (fc): Sequential(  
    (0): Linear(in_features=256, out_features=120, bias=True)  
    (1): Sigmoid()  
    (2): Linear(in_features=120, out_features=84, bias=True)  
    (3): Sigmoid()  
    (4): Linear(in_features=84, out_features=10, bias=True)  
  )  
)
```

深度卷积神经网络（AlexNet）

- 在LeNet提出后的将近20年里，神经网络一度被其他机器学习方法超越，如支持向量机。虽然LeNet可以在早期的小数据集上取得好的成绩，但是在更大的真实数据集上的表现并不尽如人意。一方面，神经网络计算复杂。虽然20世纪90年代也有过一些针对神经网络的加速硬件，但并没有像之后GPU那样大量普及。因此，训练一个多通道、多层和有大量参数的卷积神经网络在当年很难完成。另一方面，当年研究者还没有大量深入研究参数初始化和非凸优化算法等诸多领域，导致复杂的神经网络的训练通常较困难。
- 2012年，AlexNet横空出世。这个模型的名字来源于论文第一作者的姓名Alex Krizhevsky。AlexNet使用了8层卷积神经网络，并以很大的优势赢得了ImageNet 2012图像识别挑战赛。它首次证明了学习到的特征可以超越手工设计的特征，从而一举打破计算机视觉研究的现状。



- 第二， AlexNet将sigmoid激活函数改成了更加简单的ReLU激活函数。一方面， ReLU激活函数的计算更简单， 例如它并没有sigmoid激活函数中的求幂运算。另一方面， ReLU激活函数在不同的参数初始化方法下使模型更容易训练。这是由于当sigmoid激活函数输出极接近0或1时， 这些区域的梯度几乎为0， 从而造成反向传播无法继续更新部分模型参数；而ReLU激活函数在正区间的梯度恒为1。因此， 若模型参数初始化不当， sigmoid函数可能在正区间得到几乎为0的梯度， 从而令模型无法得到有效训练。
- 第三， AlexNet通过丢弃法来控制全连接层的模型复杂度。而LeNet并没有使用丢弃法。
- 第四， AlexNet引入了大量的图像增广， 如翻转、裁剪和颜色变化， 从而进一步扩大数据集来缓解过拟合。

```

class AlexNet(nn.Module):
    def __init__(self):
        super(AlexNet, self).__init__()
        self.conv = nn.Sequential(
            nn.Conv2d(1, 96, 11, 4),
            # in_channels, out_channels, kernel_size, stride, padding
            nn.ReLU(),
            nn.MaxPool2d(3, 2), # kernel_size, stride
            # 减小卷积窗口，使用填充为2来使得输入与输出的高
            # 和宽一致，且增大输出通道数
            nn.Conv2d(96, 256, 5, 1, 2),
            nn.ReLU(),
            nn.MaxPool2d(3, 2),
            # 连续3个卷积层，且使用更小的卷积窗口。除了最后
            # 的卷积层外，进一步增大了输出通道数。
            # 前两个卷积层后不使用池化层来减小输入的高和宽
            nn.Conv2d(256, 384, 3, 1, 1),
            nn.ReLU(),
            nn.Conv2d(384, 384, 3, 1, 1),
            nn.ReLU(),
            nn.Conv2d(384, 256, 3, 1, 1),
            nn.ReLU(),
            nn.MaxPool2d(3, 2)
        )

```

```

# 这里全连接层的输出个数比LeNet中的大数倍。
# 使用丢弃层来缓解过拟合
self.fc = nn.Sequential(
    nn.Linear(256*5*5, 4096),
    nn.ReLU(),
    nn.Dropout(0.5),
    nn.Linear(4096, 4096),
    nn.ReLU(),
    nn.Dropout(0.5),
    # 输出层。由于这里使用Fashion-MNIST，
    # 所以用类别数为10，而非论文中的1000
    nn.Linear(4096, 10),
)

```

```

def forward(self, img):
    feature = self.conv(img)
    output = self.fc(feature.view(img.shape[0], -1))
    return output

```

```
net = AlexNet()
print(net)
输出 :
AlexNet(
  (conv): Sequential(
    (0): Conv2d(1, 96, kernel_size=(11, 11), stride=(4, 4))
    (1): ReLU()
    (2): MaxPool2d(kernel_size=3, stride=2, padding=0, dilation=1)
    (3): Conv2d(96, 256, kernel_size=(5, 5), stride=(1, 1), padding=(2, 2))
    (4): ReLU()
    (5): MaxPool2d(kernel_size=3, stride=2, padding=0, dilation=1)
    (6): Conv2d(256, 384, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (7): ReLU()
    (8): Conv2d(384, 384, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (9): ReLU()
    (10): Conv2d(384, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (11): ReLU()
    (12): MaxPool2d(kernel_size=3, stride=2, padding=0, dilation=1, ceil_mode=False)
  )
  (fc): Sequential(
    (0): Linear(in_features=6400, out_features=4096, bias=True)
    (1): ReLU()
    (2): Dropout(p=0.5)
    (3): Linear(in_features=4096, out_features=4096, bias=True)
    (4): ReLU()
    (5): Dropout(p=0.5)
    (6): Linear(in_features=4096, out_features=10, bias=True)
  )
)
```

使用重复元素的网络（VGG）

- AlexNet在LeNet的基础上增加了3个卷积层。但AlexNet作者对它们的卷积窗口、输出通道数和构造顺序均做了大量的调整。虽然AlexNet指明了深度卷积神经网络可以取得出色的结果，但并没有提供简单的规则以指导后来的研究者如何设计新的网络。
- 下面介绍VGG，它的名字来源于论文作者所在的实验室Visual Geometry Group。VGG提出了可以通过重复使用简单的基础块来构建深度模型的思路。

VGG块

- VGG块的组成规律是：连续使用数个相同的填充为1、窗口形状为 3×3 的卷积层后接上一个步幅为2、窗口形状为 2×2 的最大池化层。卷积层保持输入的高和宽不变，而池化层则对其减半。我们使用vgg_block函数来实现这个基础的VGG块，它指定了卷积层数量和输入输出通道数。

```
def vgg_block(num_convs, in_channels, out_channels):  
    blk = []  
    for i in range(num_convs):  
        if i == 0:  
            blk.append(nn.Conv2d(in_channels, out_channels, kernel_size=3, padding=1))  
        else:  
            blk.append(nn.Conv2d(out_channels, out_channels, kernel_size=3, padding=1))  
        blk.append(nn.ReLU())  
    blk.append(nn.MaxPool2d(kernel_size=2, stride=2)) # 这里会使宽高减半  
    return nn.Sequential(*blk)
```

VGG网络

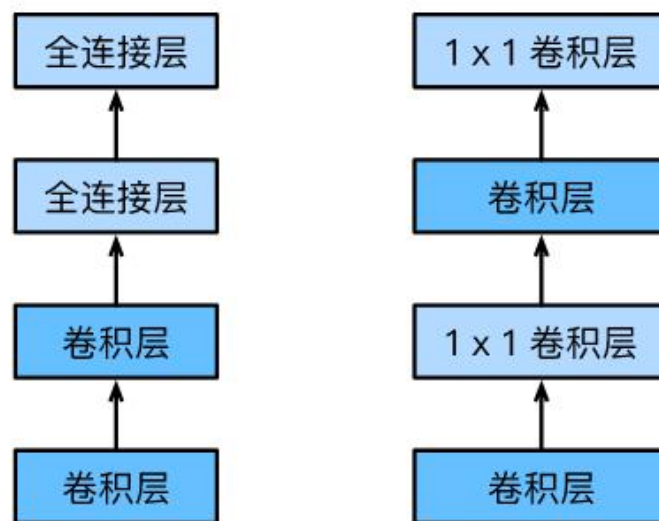
- 与AlexNet和LeNet一样，VGG网络由卷积层模块后接全连接层模块构成。卷积层模块串联数个vgg_block，其超参数由变量conv_arch定义。该变量指定了每个VGG块里卷积层个数和输入输出通道数。全连接模块则跟AlexNet中的一样。
- 现在我们构造一个VGG网络。它有5个卷积块，前2块使用单卷积层，而后3块使用双卷积层。第一块的输入输出通道分别是1（因为下面要使用的Fashion-MNIST数据的通道数为1）和64，之后每次对输出通道数翻倍，直到变为512。因为这个网络使用了8个卷积层和3个全连接层，所以经常被称为VGG-11。

网络中的网络 (NiN)

- LeNet、AlexNet和VGG在设计上的共同之处是：先以由卷积层构成的模块充分抽取空间特征，再以由全连接层构成的模块来输出分类结果。其中，AlexNet和VGG对LeNet的改进主要在于如何对这两个模块加宽（增加通道数）和加深。本节我们介绍网络中的网络（NiN）。它提出了另外一个思路，即串联多个由卷积层和“全连接”层构成的小网络来构建一个深层网络。

NiN块

- 卷积层的输入和输出通常是四维数组（样本，通道，高，宽），而全连接层的输入和输出则通常是二维数组（样本，特征）。如果想在全连接层后再接上卷积层，则需要将全连接层的输出变换为四维。回忆 1×1 卷积层。它可以看成全连接层，其中空间维度（高和宽）上的每个元素相当于样本，通道相当于特征。因此，NiN使用 1×1 卷积层来替代全连接层，从而使空间信息能够自然传递到后面的层中去。下图对比了NiN同AlexNet和VGG等网络在结构上的主要区别。



NiN块

- NiN块是NiN中的基础块。它由一个卷积层加两个充当全连接层的 1×1 卷积层串联而成。其中第一个卷积层的超参数可以自行设置，而第二和第三个卷积层的超参数一般是固定的。

```
def nin_block(in_channels, out_channels, kernel_size, stride, padding):
    blk = nn.Sequential(nn.Conv2d(in_channels, out_channels, kernel_size, stride, padding),
                        nn.ReLU(),
                        nn.Conv2d(out_channels, out_channels, kernel_size=1),
                        nn.ReLU(),
                        nn.Conv2d(out_channels, out_channels, kernel_size=1),
                        nn.ReLU())
    return blk
```

NiN模型

- NiN是在AlexNet问世不久后提出的。它们的卷积层设定有类似之处。NiN使用卷积窗口形状分别为 11×11 、 5×5 和 3×3 的卷积层，相应的输出通道数也与AlexNet中的一致。每个NiN块后接一个步幅为2、窗口形状为 3×3 的最大池化层。
- 除使用NiN块以外，NiN还有一个设计与AlexNet显著不同：NiN去掉了AlexNet最后的3个全连接层，取而代之地，NiN使用了输出通道数等于标签类别数的NiN块，然后使用全局平均池化层对每个通道中所有元素求平均并直接用于分类。这里的全局平均池化层即窗口形状等于输入空间维形状的平均池化层。NiN的这个设计的好处是可以显著减小模型参数尺寸，从而缓解过拟合。然而，该设计有时会造成获得有效模型的训练时间的增加。

```
class GlobalAvgPool2d(nn.Module):
```

```
    # 全局平均池化层可通过将池化窗口形状设置成输入的高和宽实现
```

```
    def __init__(self):
```

```
        super(GlobalAvgPool2d, self).__init__()
```

```
    def forward(self, x):
```

```
        return F.avg_pool2d(x, kernel_size=x.size()[2:])
```

```
net = nn.Sequential(
```

```
    nin_block(1, 96, kernel_size=11, stride=4, padding=0),
```

```
    nn.MaxPool2d(kernel_size=3, stride=2),
```

```
    nin_block(96, 256, kernel_size=5, stride=1, padding=2),
```

```
    nn.MaxPool2d(kernel_size=3, stride=2),
```

```
    nin_block(256, 384, kernel_size=3, stride=1, padding=1),
```

```
    nn.MaxPool2d(kernel_size=3, stride=2),
```

```
    nn.Dropout(0.5),
```

```
    nin_block(384, 10, kernel_size=3, stride=1, padding=1), # 标签类别数是10
```

```
    GlobalAvgPool2d(),
```

```
    d2l.FlattenLayer()) # 将四维的输出转成二维的输出，其形状为(批量大小, 10)
```

NiN模型

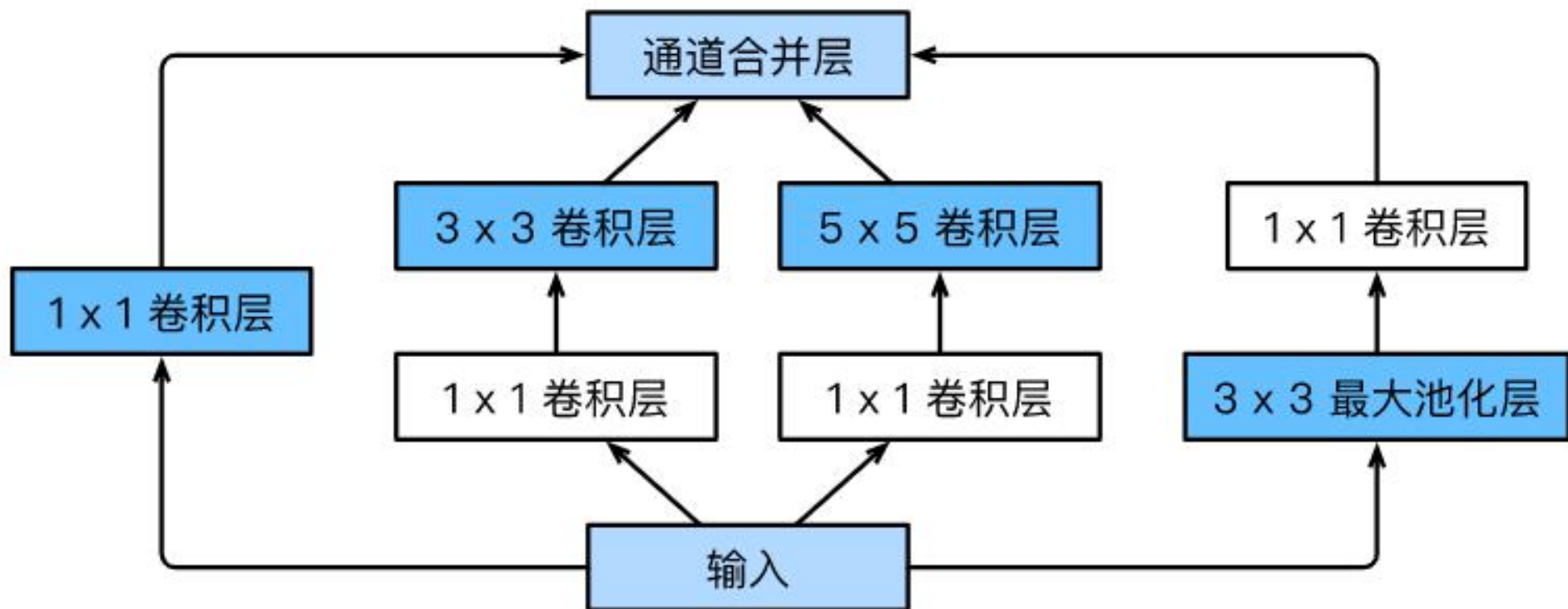
- NiN重复使用由卷积层和代替全连接层的 1×1 卷积层构成的NiN块来构建深层网络。
- NiN去除了容易造成过拟合的全连接输出层，而是将其替换成输出通道数等于标签类别数的NiN块和全局平均池化层。
- NiN的以上设计思想影响了后面一系列卷积神经网络的设计。

含并行连结的网络（GoogLeNet）

- 在2014年的ImageNet图像识别挑战赛中，一个名叫GoogLeNet的网络结构大放异彩。它虽然在名字上向LeNet致敬，但在网络结构上已经很难看到LeNet的影子。GoogLeNet吸收了NiN中网络串联网络的思想，并在此基础上做了很大改进。在随后的几年里，研究人员对GoogLeNet进行了数次改进，这里介绍这个模型系列的第一个版本。

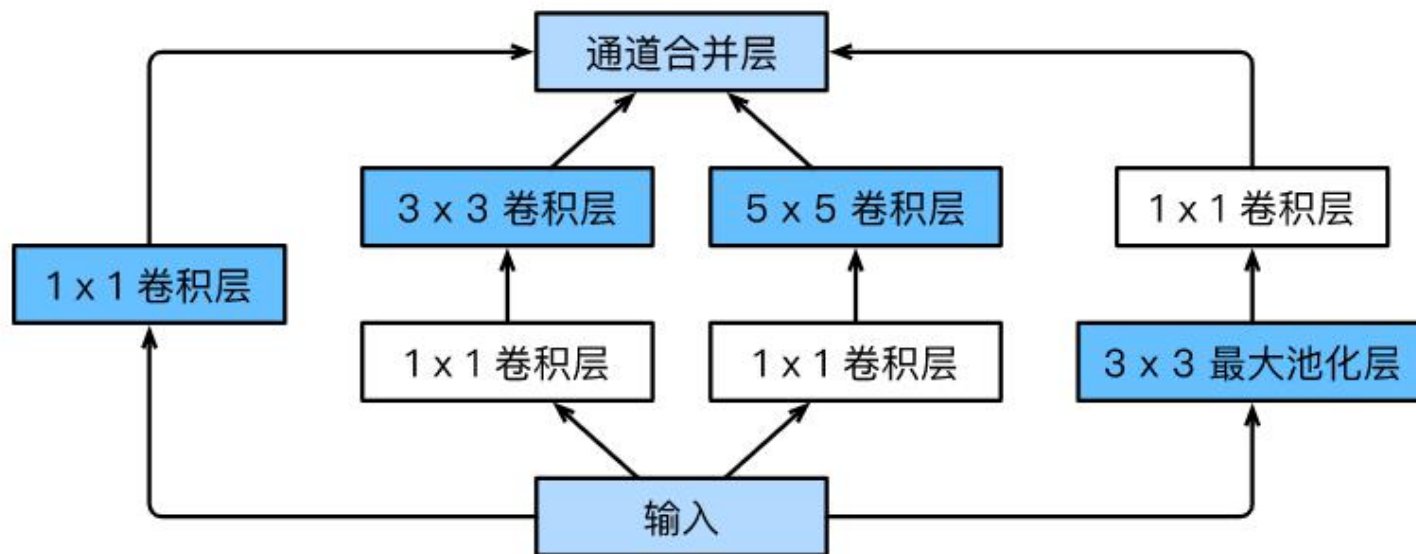
Inception 块

- GoogLeNet中的基础卷积块叫作Inception块，得名于同名电影《盗梦空间》（Inception）。与介绍的NiN块相比，这个基础块在结构上更加复杂，如图所示。



Inception 块

- 由图可以看出，Inception块里有4条并行的线路。前3条线路使用窗口大小分别是 1×1 、 3×3 和 5×5 的卷积层来抽取不同空间尺寸下的信息，其中中间2个线路会对输入先做 1×1 卷积来减少输入通道数，以降低模型复杂度。第四条线路则使用 3×3 最大池化层，后接 1×1 卷积层来改变通道数。4条线路都使用了合适的填充来使输入与输出的高和宽一致。最后我们将每条线路的输出在通道维上连结，并输入接下来的层中去。



```
class Inception(nn.Module):
```

```
    def __init__(self, in_c, c1, c2, c3, c4): # c1 - c4为每条线路里的层的输出通道数
```

```
        super(Inception, self).__init__()
```

```
        self.p1_1 = nn.Conv2d(in_c, c1, kernel_size=1) # 线路1, 单1 x 1卷积层
```

```
        self.p2_1 = nn.Conv2d(in_c, c2[0], kernel_size=1) # 线路2, 1 x 1卷积层后接3 x 3卷积层
```

```
        self.p2_2 = nn.Conv2d(c2[0], c2[1], kernel_size=3, padding=1)
```

```
        self.p3_1 = nn.Conv2d(in_c, c3[0], kernel_size=1) # 线路3, 1 x 1卷积层后接5 x 5卷积层
```

```
        self.p3_2 = nn.Conv2d(c3[0], c3[1], kernel_size=5, padding=2)
```

```
        self.p4_1 = nn.MaxPool2d(kernel_size=3, stride=1, padding=1) # 线路4, 3 x 3最大池化层后接1 x 1卷积层
```

```
        self.p4_2 = nn.Conv2d(in_c, c4, kernel_size=1)
```

```
    def forward(self, x):
```

```
        p1 = F.relu(self.p1_1(x))
```

```
        p2 = F.relu(self.p2_2(F.relu(self.p2_1(x))))
```

```
        p3 = F.relu(self.p3_2(F.relu(self.p3_1(x))))
```

```
        p4 = F.relu(self.p4_2(self.p4_1(x)))
```

```
        return torch.cat((p1, p2, p3, p4), dim=1) # 在通道维上连结输出
```

GoogLeNet模型

- GoogLeNet跟VGG一样，在主体卷积部分中使用5个模块（block），每个模块之间使用步幅为2的 3×3 最大池化层来减小输出高宽。第一模块使用一个64通道的 7×7 卷积层。

```
b1 = nn.Sequential(nn.Conv2d(1, 64, kernel_size=7, stride=2, padding=3),  
                  nn.ReLU(),  
                  nn.MaxPool2d(kernel_size=3, stride=2, padding=1))
```

GoogLeNet模型

- 第二模块使用2个卷积层：首先是64通道的 1×1 卷积层，然后将通道增大3倍的 3×3 卷积层。它对应Inception块中的第二条线路。

```
b2 = nn.Sequential(nn.Conv2d(64, 64, kernel_size=1),  
                  nn.Conv2d(64, 192, kernel_size=3, padding=1),  
                  nn.MaxPool2d(kernel_size=3, stride=2, padding=1))
```

GoogLeNet模型

- 第三模块串联2个完整的Inception块。第一个Inception块的输出通道数为 $64+128+32+32=256$ ，其中4条线路的输出通道数比例为 $64:128:32:32=2:4:1:1$ 。其中第二、第三条线路先分别将输入通道数减小至 $96/192=1/2$ 和 $16/192=1/12$ 后，再接上第二层卷积层。第二个Inception块输出通道数增至 $128+192+96+64=480$ ，每条线路的输出通道数之比为 $128:192:96:64 = 4:6:3:2$ 。其中第二、第三条线路先分别将输入通道数减小至 $128/256=1/2$ 和 $32/256=1/8$ 。

```
b3 = nn.Sequential(Inception(192, 64, (96, 128), (16, 32), 32),  
                  Inception(256, 128, (128, 192), (32, 96), 64),  
                  nn.MaxPool2d(kernel_size=3, stride=2, padding=1))
```

GoogLeNet模型

- 第四模块更加复杂。它串联了5个Inception块，其输出通道数分别是 $192+208+48+64=512$ 、 $160+224+64+64=512$ 、 $128+256+64+64=512$ 、 $112+288+64+64=528$ 和 $256+320+128+128=832$ 。这些线路的通道数分配和第三模块中的类似，首先含 3×3 卷积层的第二条线路输出最多通道，其次是仅含 1×1 卷积层的第一条线路，之后是含 5×5 卷积层的第三条线路和含 3×3 最大池化层的第四条线路。其中第二、第三条线路都会先按比例减小通道数。这些比例在各个Inception块中都略有不同。

GoogLeNet模型

```
b4 = nn.Sequential(Inception(480, 192, (96, 208), (16, 48), 64),  
                  Inception(512, 160, (112, 224), (24, 64), 64),  
                  Inception(512, 128, (128, 256), (24, 64), 64),  
                  Inception(512, 112, (144, 288), (32, 64), 64),  
                  Inception(528, 256, (160, 320), (32, 128), 128),  
                  nn.MaxPool2d(kernel_size=3, stride=2, padding=1))
```

GoogLeNet模型

- 第五模块有输出通道数为 $256+320+128+128=832$ 和 $384+384+128+128=1024$ 的两个Inception块。其中每条线路的通道数的分配思路和第三、第四模块中的一致，只是在具体数值上有所不同。需要注意的是，第五模块的后面紧跟输出层，该模块同NiN一样使用全局平均池化层来将每个通道的高和宽变成1。最后我们将输出变成二维数组后接上一个输出个数为标签类别数的全连接层。

```
b5 = nn.Sequential(Inception(832, 256, (160, 320), (32, 128), 128),  
                  Inception(832, 384, (192, 384), (48, 128), 128),  
                  d2l.GlobalAvgPool2d())
```

```
net = nn.Sequential(b1, b2, b3, b4, b5,  
                   d2l.FlattenLayer(), nn.Linear(1024, 10))
```


GoogLeNet模型

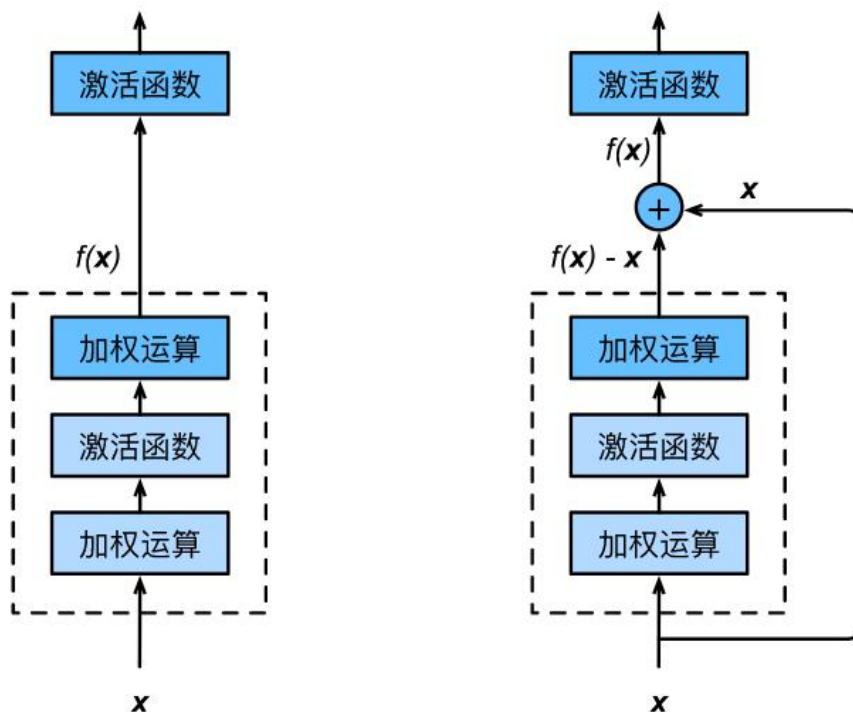
- Inception块相当于一个有4条线路的子网络。它通过不同窗口形状的卷积层和最大池化层来并行抽取信息，并使用 1×1 卷积层减少通道数从而降低模型复杂度。
- GoogLeNet将多个设计精细的Inception块和其他层串联起来。其中Inception块的通道数分配之比是在ImageNet数据集上通过大量的实验得来的。
- GoogLeNet和它的后继者们一度是ImageNet上最高效的模型之一：在类似的测试精度下，它们的计算复杂度往往更低。

残差网络 (ResNet)

- 先思考一个问题：对神经网络模型添加新的层，充分训练后的模型是否只可能更有效地降低训练误差？理论上，原模型解的空间只是新模型解的空间的子空间。也就是说，如果我们能将新添加的层训练成恒等映射 $f(x) = x$ ，新模型和原模型将同样有效。由于新模型可能得出更优的解来拟合训练数据集，因此添加层似乎更容易降低训练误差。然而在实践中，添加过多的层后训练误差往往不降反升。即使利用批量归一化带来的数值稳定性使训练深层模型更加容易，该问题仍然存在。针对这一问题，何恺明等人提出了残差网络 (ResNet)。它在2015年的ImageNet图像识别挑战赛夺魁，并深刻影响了后来的深度神经网络的设计。

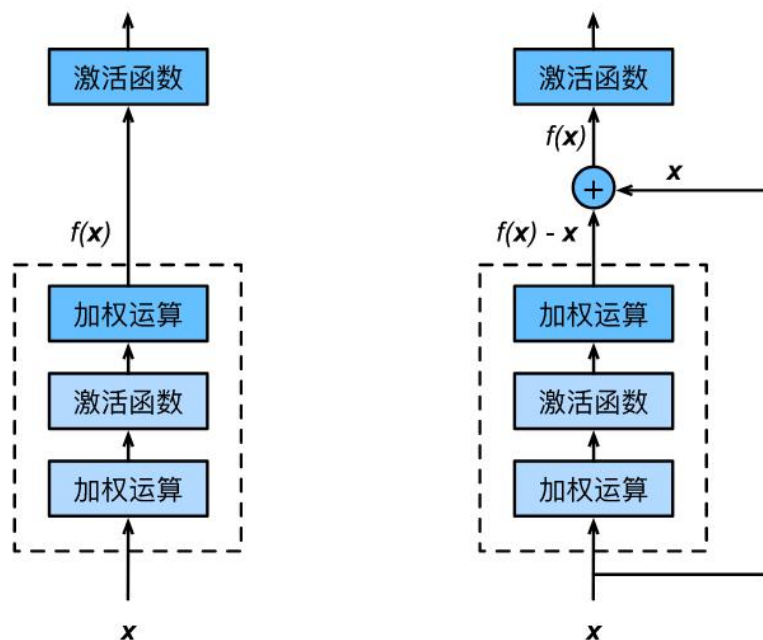
残差块

- 让我们聚焦于神经网络局部。如图所示，设输入为 x 。假设我们希望学出的理想映射为 $f(x)$ ，从而作为图上方激活函数的输入。左图虚线框中的部分需要直接拟合出该映射 $f(x)$ ，而右图虚线框中的部分则需要拟合出有关恒等映射的残差映射 $f(x)-x$ 。残差映射在实际中更容易优化。



残差块

- 以恒等映射作为我们希望学出的理想映射 $f(x)$ 。我们只需将右图虚线框内上方的加权运算（如仿射）的权重和偏差参数学成1和0，那么 $f(x)$ 即为恒等映射。实际中，当理想映射 $f(x)$ 极接近于恒等映射时，残差映射也易于捕捉恒等映射的细微波动。右图也是ResNet的基础块，即残差块（residual block）。在残差块中，输入可通过跨层的数据线路更快地向前传播。



残差块

- ResNet沿用了VGG全 3×3 卷积层的设计。残差块里首先有2个有相同输出通道数的 3×3 卷积层。每个卷积层后接一个批量归一化层和ReLU激活函数。然后将输入跳过这两个卷积运算后直接加在最后的ReLU激活函数前。这样的设计要求两个卷积层的输出与输入形状一样，从而可以相加。如果想改变通道数，就需要引入一个额外的 1×1 卷积层来将输入变换成需要的形状后再做相加运算。
- 残差块的实现如下。它可以设定输出通道数、是否使用额外的 1×1 卷积层来修改通道数以及卷积层的步幅。

```
class Residual(nn.Module): # 本类已保存在d2lzh_pytorch包中方便以后使用
    def __init__(self, in_channels, out_channels, use_1x1conv=False, stride=1):
        super(Residual, self).__init__()
        self.conv1 = nn.Conv2d(in_channels, out_channels, kernel_size=3, padding=1, stride=stride)
        self.conv2 = nn.Conv2d(out_channels, out_channels, kernel_size=3, padding=1)
        if use_1x1conv:
            self.conv3 = nn.Conv2d(in_channels, out_channels, kernel_size=1, stride=stride)
        else:
            self.conv3 = None
        self.bn1 = nn.BatchNorm2d(out_channels)
        self.bn2 = nn.BatchNorm2d(out_channels)

    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)
        return F.relu(Y + X)
```

ResNet模型

- ResNet的前两层跟之前介绍的GoogLeNet中的一样：在输出通道数为64、步幅为2的 7×7 卷积层后接步幅为2的 3×3 的最大池化层。不同之处在于ResNet每个卷积层后增加的批量归一化层。

```
net = nn.Sequential(  
    nn.Conv2d(1, 64, kernel_size=7, stride=2, padding=3),  
    nn.BatchNorm2d(64),  
    nn.ReLU(),  
    nn.MaxPool2d(kernel_size=3, stride=2, padding=1))
```

ResNet模型

- GoogLeNet在后面接了4个由Inception块组成的模块。ResNet则使用4个由残差块组成的模块，每个模块使用若干个同样输出通道数的残差块。第一个模块的通道数同输入通道数一致。由于之前已经使用了步幅为2的最大池化层，所以无须减小高和宽。之后模块在第一个残差块里将上个模块的通道数翻倍，并将高宽减半。

```
def resnet_block(in_channels, out_channels, num_residuals, first_block=False):
    if first_block:
        assert in_channels == out_channels # 第一个模块的通道数同输入通道数一致
    blk = []
    for i in range(num_residuals):
        if i == 0 and not first_block:
            blk.append(Residual(in_channels, out_channels, use_1x1conv=True, stride=2))
        else:
            blk.append(Residual(out_channels, out_channels))
    return nn.Sequential(*blk)
```


ResNet模型

- 接着我们为ResNet加入所有残差块。这里每个模块使用两个残差块。

```
net.add_module("resnet_block1", resnet_block(64, 64, 2, first_block=True))
```

```
net.add_module("resnet_block2", resnet_block(64, 128, 2))
```

```
net.add_module("resnet_block3", resnet_block(128, 256, 2))
```

```
net.add_module("resnet_block4", resnet_block(256, 512, 2))
```

- 最后，与GoogLeNet一样，加入全局平均池化层后接上全连接层输出。

```
net.add_module("global_avg_pool", d2l.GlobalAvgPool2d()) # GlobalAvgPool2d的输出: (Batch, 512, 1, 1)
```

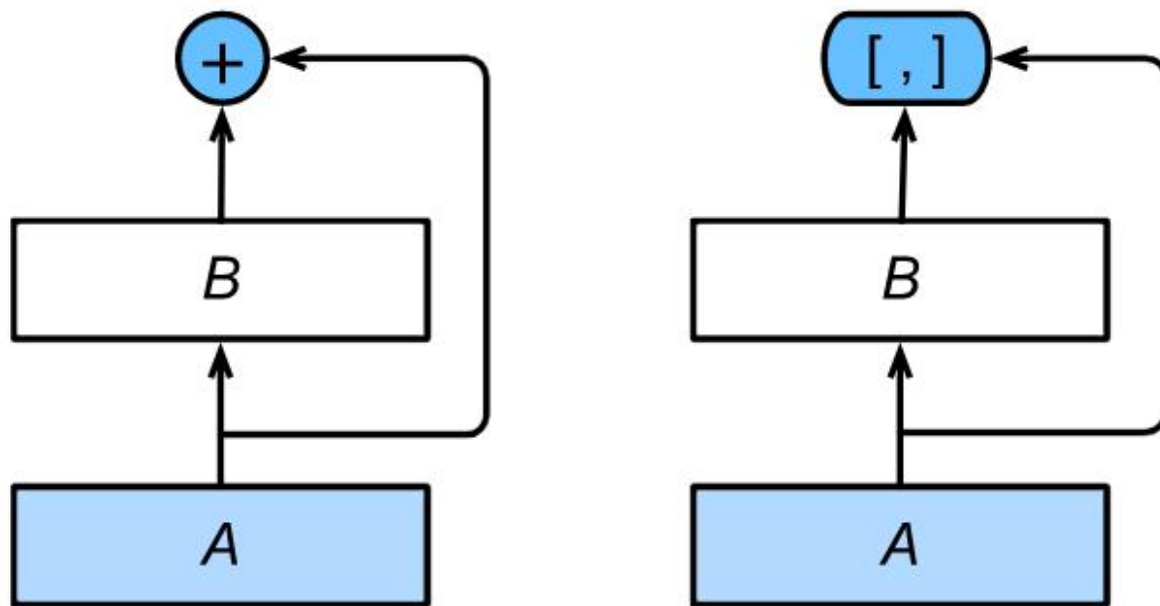
```
net.add_module("fc", nn.Sequential(d2l.FlattenLayer(), nn.Linear(512, 10)))
```

ResNet模型

- 这里每个模块里有4个卷积层（不计算 1×1 卷积层），加上最开始的卷积层和最后的全连接层，共计18层。这个模型通常也被称为ResNet-18。通过配置不同的通道数和模块里的残差块数可以得到不同的ResNet模型，例如更深的含152层的ResNet-152。虽然ResNet的主体架构跟GoogLeNet的类似，但ResNet结构更简单，修改更方便。这些因素都导致了ResNet迅速被广泛使用。

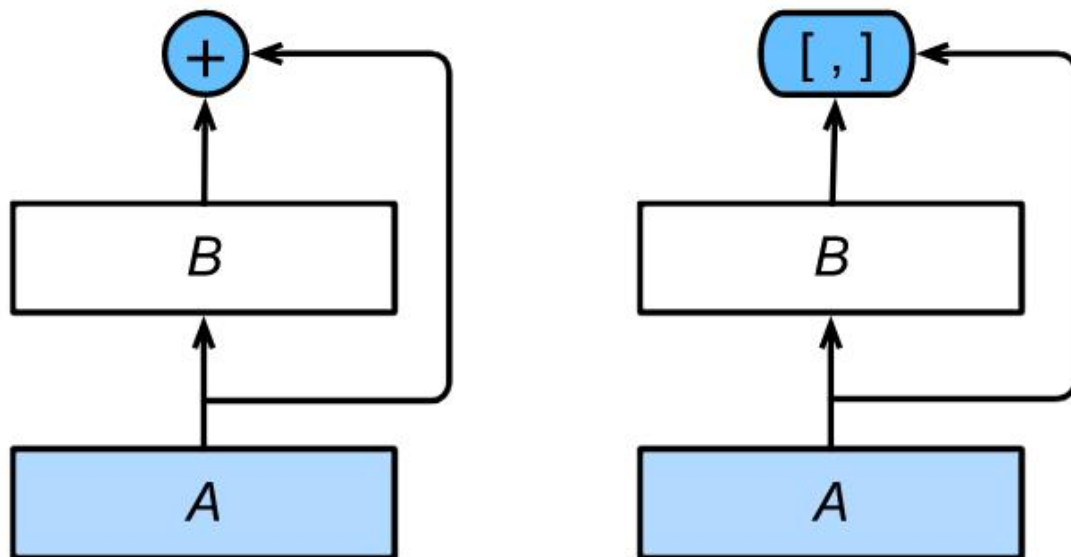
稠密连接网络 (DenseNet)

- ResNet中的跨层连接设计引申出了数个后续工作。我们介绍其中的一个：稠密连接网络 (DenseNet)。它与ResNet的主要区别如图所示。



稠密连接网络 (DenseNet)

- 图中将部分前后相邻的运算抽象为模块A和模块B。与ResNet的主要区别在于，DenseNet里模块B的输出不是像ResNet那样和模块A的输出相加，而是在通道维上连结。这样模块A的输出可以直接传入模块B后面的层。在这个设计里，模块A直接跟模块B后面的所有层连接在了一起。这也是它被称为“稠密连接”的原因。
- DenseNet的主要构建模块是稠密块 (dense block) 和过渡层 (transition layer) 。前者定义了输入和输出是如何连结的，后者则用来控制通道数，使之不过大。



稠密块

- DenseNet使用了ResNet改良版的“批量归一化、激活和卷积”结构，我们首先在conv_block函数里实现这个结构。

```
def conv_block(in_channels, out_channels):  
    blk = nn.Sequential(nn.BatchNorm2d(in_channels),  
                        nn.ReLU(),  
                        nn.Conv2d(in_channels, out_channels, kernel_size=3, padding=1))  
    return blk
```

- 稠密块由多个conv_block组成，每块使用相同的输出通道数。但在前向计算时，我们将每块的输入和输出在通道维上连结。

```
class DenseBlock(nn.Module):
    def __init__(self, num_convs, in_channels, out_channels):
        super(DenseBlock, self).__init__()
        net = []
        for i in range(num_convs):
            in_c = in_channels + i * out_channels
            net.append(conv_block(in_c, out_channels))
        self.net = nn.ModuleList(net)
        self.out_channels = in_channels + num_convs * out_channels # 计算输出通道数

    def forward(self, X):
        for blk in self.net:
            Y = blk(X)
            X = torch.cat((X, Y), dim=1) # 在通道维上将输入和输出连结
        return X
```

稠密块

- 在下面的例子中，我们定义一个有2个输出通道数为10的卷积块。使用通道数为3的输入时，我们会得到通道数为 $3+2\times 10=23$ 的输出。卷积块的通道数控制了输出通道数相对于输入通道数的增长，因此也被称为增长率（growth rate）。

```
blk = DenseBlock(2, 3, 10)
```

```
X = torch.rand(4, 3, 8, 8)
```

```
Y = blk(X)
```

```
Y.shape # torch.Size([4, 23, 8, 8])
```

过渡层

- 由于每个稠密块都会带来通道数的增加，使用过多则会带来过于复杂的模型。过渡层用来控制模型复杂度。它通过 1×1 卷积层来减小通道数，并使用步幅为2的平均池化层减半高和宽，从而进一步降低模型复杂度。

```
def transition_block(in_channels, out_channels):  
    blk = nn.Sequential(  
        nn.BatchNorm2d(in_channels),  
        nn.ReLU(),  
        nn.Conv2d(in_channels, out_channels, kernel_size=1),  
        nn.AvgPool2d(kernel_size=2, stride=2))  
    return blk
```


DenseNet模型

- 我们来构造DenseNet模型。DenseNet首先使用同ResNet一样的单卷积层和最大池化层。

```
net = nn.Sequential(  
    nn.Conv2d(1, 64, kernel_size=7, stride=2, padding=3),  
    nn.BatchNorm2d(64),  
    nn.ReLU(),  
    nn.MaxPool2d(kernel_size=3, stride=2, padding=1))
```

DenseNet模型

- 类似于ResNet接下来使用的4个残差块，DenseNet使用的是4个稠密块。同ResNet一样，我们可以设置每个稠密块使用多少个卷积层。这里我们设成4，从而与ResNet-18保持一致。稠密块里的卷积层通道数（即增长率）设为32，所以每个稠密块将增加128个通道。
- ResNet里通过步幅为2的残差块在每个模块之间减小高和宽。这里我们则使用过渡层来减半高和宽，并减半通道数。

DenseNet模型

num_channels, growth_rate = 64, 32 # num_channels为当前的通道数

num_convs_in_dense_blocks = [4, 4, 4, 4]

for i, num_convs in enumerate(num_convs_in_dense_blocks):

 DB = DenseBlock(num_convs, num_channels, growth_rate)

 net.add_module("DenseBlock_%d" % i, DB)

 # 上一个稠密块的输出通道数

 num_channels = DB.out_channels

 # 在稠密块之间加入通道数减半的过渡层

 if i != len(num_convs_in_dense_blocks) - 1:

 net.add_module("transition_block_%d" % i, transition_block(num_channels, num_channels // 2))

 num_channels = num_channels // 2

DenseNet模型

- 同ResNet一样，最后接上全局池化层和全连接层来输出。

```
net.add_module("BN", nn.BatchNorm2d(num_channels))
```

```
net.add_module("relu", nn.ReLU())
```

```
net.add_module("global_avg_pool", d2l.GlobalAvgPool2d())
```

```
# GlobalAvgPool2d的输出: (Batch, num_channels, 1, 1)
```

```
net.add_module("fc", nn.Sequential(d2l.FlattenLayer(), nn.Linear(num_channels, 10)))
```