

## 参考文献

- (1) Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., & Anguelov, D. & Rabinovich, A.(2015). Going deeper with convolutions. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 1-9).
- (2) Ioffe, S., & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. arXiv preprint arXiv:1502.03167.
- (3) Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., & Wojna, Z. (2016). Rethinking the inception architecture for computer vision. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 2818-2826).
- (4) Szegedy, C., Ioffe, S., Vanhoucke, V., & Alemi, A. A. (2017, February). Inception-v4, inception-resnet and the impact of residual connections on learning. In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 4, p. 12).

---

注：除代码外本节与原书此节基本相同，[原书传送门](#)

## 5.10 批量归一化

本节我们介绍批量归一化（batch normalization）层，它能让较深的神经网络的训练变得更加容易（1）。在3.16节（实战Kaggle比赛：预测房价）里，我们对输入数据做了标准化处理：处理后的任意一个特征在数据集中所有样本上的均值为0、标准差为1。标准化处理输入数据使各个特征的分布相近：这往往更容易训练出有效的模型。

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。批量归一化和下一节将要介绍的残差网络为训练和设计深度模型提供了两类重要思路。

### 5.10.1 批量归一化层

对全连接层和卷积层做批量归一化的方法稍有不同。下面我们将分别介绍这两种情况下的批量归一化。

### 5.10.1.1 对全连接层做批量归一化

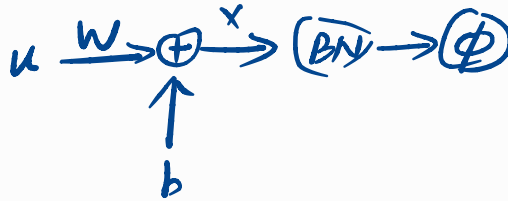
我们先考虑如何对全连接层做批量归一化。通常，我们将批量归一化层置于全连接层中的仿射变换和激活函数之间。设全连接层的输入为  $\mathbf{u}$ ，权重参数和偏差参数分别为  $\mathbf{W}$  和  $\mathbf{b}$ ，激活函数为  $\phi$ 。设批量归一化的运算符为 BN。那么，使用批量归一化的全连接层的输出为

Batch-size \* 输出神经元个数

output =  $\phi(\text{BN}(\mathbf{x}))$ ,  
对中间值  $\mathbf{x}$  做标准化，让其在 batch\_size 上的均值为 0，标准差为 1。

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  $\mathcal{B} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ 。它们正是批量归一化层的输入。对于小批量  $\mathcal{B}$  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

$$\mathbf{y}^{(i)} = \text{BN}(\mathbf{x}^{(i)}),$$

并由以下几步求得。首先，对小批量  $\mathcal{B}$  求均值和方差：

m: batch\_size 大小

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)}, \sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)} - \mu_{\mathcal{B}})^2,$$

d维

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：

$$\hat{\mathbf{x}}^{(i)} \leftarrow \frac{\mathbf{x}^{(i)} - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}},$$

这里  $\epsilon > 0$  是一个很小的常数，保证分母大于 0。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  $\gamma$  和偏移 (shift) 参数  $\beta$ 。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：

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

至此，我们得到了  $\mathbf{x}^{(i)}$  的批量归一化的输出  $\mathbf{y}^{(i)}$ 。值得注意的是，可学习的拉伸和偏移参数保留了不对  $\hat{\mathbf{x}}^{(i)}$  做批量归一化的可能：此时只需学出  $\gamma = \sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}$  和  $\beta = \mu_{\mathcal{B}}$ 。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。  
即如果两个参数等于上面这样，x就回到了归一化前的值，自己带到上面的式子就知道。

### 5.10.1.2 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  $p$  和  $q$ 。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。

### 5.10.1.3 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。

## 5.10.2 从零开始实现

下面我们自己实现批量归一化层。

```
1 import time
2 import torch
3 from torch import nn, optim
4 import torch.nn.functional as F
5
6 import sys
7 sys.path.append("..")
8 import d2lzh_pytorch as d2l
9 device = torch.device('cuda' if torch.cuda.is_available() else
   'cpu')
10
11 def batch_norm(is_training, X, gamma, beta, moving_mean,
   moving_var, eps, momentum):
12     # 判断当前模式是训练模式还是预测模式
13     if not is_training:
14         # 如果是在预测模式下，直接使用传入的移动平均所得的均值和方差
15         X_hat = (X - moving_mean) / torch.sqrt(moving_var + eps)
16     else:
17         # 以batch为单位计算均值和方差
18         assert len(X.shape) in (2, 4) # 先判断X的形状是2还是4;
19         if len(X.shape) == 2: # 若是2，则是全连接层仿射变换的输出，则在其第0维上求平均。因为全连接层的维度是m*d，m是batch_size，
20             # 所以对全连接层就是在其batch_size的维度上求平均。
21             # 使用全连接层的情况，计算特征维上的均值和方差
22             mean = X.mean(dim=0) # 求出的平均值第1维的值是d
23             var = ((X - mean) ** 2).mean(dim=0)
24         else: # 若X.shape = 4，则X是卷积计算后得到的输出，即batch_size*通道数*高*宽
25             # 使用二维卷积层的情况，计算通道维上 (axis=1) 的均值和方差。这里
26             # 即对batch_size*h*w个元素同时做批量归一化，使用相同的均值和方差。
27             # 我们需要保持
28             # x的形状以便后面可以做广播运算
29             mean = X.mean(dim=0, keepdim=True).mean(dim=2,
30 keepdim=True).mean(dim=3, keepdim=True) # 后→前，先在3rd维上求均值，再在2nd维上求均值，后在0th维上求均值，分别对应应在w、
31 h、batch_size上求均值。
32             var = ((X - mean) ** 2).mean(dim=0,
33 keepdim=True).mean(dim=2, keepdim=True).mean(dim=3, keepdim=True)
34             # 训练模式下用当前的均值和方差做标准化
35             # 公式: X_hat = (X - mean) / torch.sqrt(var + eps)
36             # 更新移动平均的均值和方差
```

```

30         超参, 自己设置
        moving_mean = momentum * moving_mean + (1.0 - momentum) *
mean
31         moving_var = momentum * moving_var + (1.0 - momentum) * var
32         Y = gamma * X_hat + beta # 拉伸和偏移
33         return Y, moving_mean, moving_var

```

接下来, 我们自定义一个 `BatchNorm` 层。它保存参与求梯度和迭代的拉伸参数 `gamma` 和偏移参数 `beta`, 同时也维护移动平均得到的均值和方差, 以便能够在模型预测时被使用。`BatchNorm` 实例所需指定的 `num_features` 参数对于全连接层来说应为输出个数, 对于卷积层来说则为输出通道数。该实例所需指定的 `num_dims` 参数对于全连接层和卷积层来说分别为2和4。

```

1  class BatchNorm(nn.Module): 全连接层: 代表输出神经元的个数;
2      def __init__(self, num_features, num_dims): 卷积层的情况下代表通道数; 2或4, 对应全连接层和卷积层
3          super(BatchNorm, self).__init__()
4          if num_dims == 2:
5              shape = (1, num_features)
6          else:
7              shape = (1, num_features, 1, 1)
8          # 参与求梯度和迭代的拉伸和偏移参数, 分别初始化成0和1
9          self.gamma = nn.Parameter(torch.ones(shape))
10         self.beta = nn.Parameter(torch.zeros(shape))
11         # 不参与求梯度和迭代的变量, 全在内存上初始化成0
12         self.moving_mean = torch.zeros(shape)
13         self.moving_var = torch.zeros(shape)
14
15     def forward(self, X):
16         # 如果x不在内存上, 将moving_mean和moving_var复制到x所在显存上
17         if self.moving_mean.device != X.device: # 如果X不在GPU上, 就复制到GPU上
18             self.moving_mean = self.moving_mean.to(X.device)
19             self.moving_var = self.moving_var.to(X.device)
20         # 保存更新过的moving_mean和moving_var, Module实例的training属性默
认为true, 调用.eval()后设成false
21         Y, self.moving_mean, self.moving_var =
batch_norm(self.training,
22             X, self.gamma, self.beta, self.moving_mean,
23             self.moving_var, eps=1e-5, momentum=0.9) # 在forward中已经调用了batch_norm了
24         return Y

```

### 5.10.2.1 使用批量归一化层的LeNet

下面我们修改5.5节（卷积神经网络（LeNet））介绍的LeNet模型, 从而应用批量归一化层。我们在所有的卷积层或全连接层之后、激活层之前加入批量归一化层。

```

1  net = nn.Sequential(
2      nn.Conv2d(1, 6, 5), # in_channels, out_channels,
kernel_size
3      BatchNorm(6, num_dims=4), # 在这一层里加上BatchNorm层即可。
4      nn.Sigmoid(),
5      nn.MaxPool2d(2, 2), # kernel_size, stride

```

```

6         nn.Conv2d(6, 16, 5),
7         BatchNorm(16, num_dims=4),
8         nn.Sigmoid(),
9         nn.MaxPool2d(2, 2),
10        d2l.FlattenLayer(),
11        nn.Linear(16*4*4, 120),
12        BatchNorm(120, num_dims=2),
13        nn.Sigmoid(),
14        nn.Linear(120, 84),
15        BatchNorm(84, num_dims=2),
16        nn.Sigmoid(),
17        nn.Linear(84, 10)
18    )

```

下面我们训练修改后的模型。

```

1  batch_size = 256
2  train_iter, test_iter =
    d2l.load_data_fashion_mnist(batch_size=batch_size)
3
4  lr, num_epochs = 0.001, 5
5  optimizer = torch.optim.Adam(net.parameters(), lr=lr)
6  d2l.train_ch5(net, train_iter, test_iter, batch_size, optimizer,
    device, num_epochs)

```

输出：

```

1  training on  cuda
2  epoch 1, loss 0.0039, train acc 0.790, test acc 0.835, time 2.9 sec
3  epoch 2, loss 0.0018, train acc 0.866, test acc 0.821, time 3.2 sec
4  epoch 3, loss 0.0014, train acc 0.879, test acc 0.857, time 2.6 sec
5  epoch 4, loss 0.0013, train acc 0.886, test acc 0.820, time 2.7 sec
6  epoch 5, loss 0.0012, train acc 0.891, test acc 0.859, time 2.8 sec

```

最后我们查看第一个批量归一化层学习到的拉伸参数 `gamma` 和偏移参数 `beta`。

```

1  net[1].gamma.view((-1,)), net[1].beta.view((-1,))

```

输出：

```

1  (tensor([ 1.2537,  1.2284,  1.0100,  1.0171,  0.9809,  1.1870],
    device='cuda:0'),
2  tensor([ 0.0962,  0.3299, -0.5506,  0.1522, -0.1556,  0.2240],
    device='cuda:0'))

```

## 5.10.3 简洁实现

与我们刚刚自己定义的 `BatchNorm` 类相比，Pytorch 中 `nn` 模块定义的 `BatchNorm1d` 和 `BatchNorm2d` 类使用起来更加简单，二者分别用于全连接层和卷积层，都需要指定输入的 `num_features` 参数值。下面我们使用 PyTorch 实现使用批量归一化的 LeNet。

```
1 net = nn.Sequential(  
2     nn.Conv2d(1, 6, 5), # in_channels, out_channels,  
   kernel_size  
3     nn.BatchNorm2d(6),  
4     nn.Sigmoid(),  
5     nn.MaxPool2d(2, 2), # kernel_size, stride  
6     nn.Conv2d(6, 16, 5),  
7     nn.BatchNorm2d(16),  
   卷积层的batch_norm  
8     nn.Sigmoid(),  
9     nn.MaxPool2d(2, 2),  
10    d2l.FlattenLayer(),  
11    nn.Linear(16*4*4, 120),  
12    nn.BatchNorm1d(120),  
13    nn.Sigmoid(),  
14    nn.Linear(120, 84),  
15    nn.BatchNorm1d(84),  
   全连接层的batch_norm  
16    nn.Sigmoid(),  
17    nn.Linear(84, 10)  
18 )
```

使用同样的超参数进行训练。

```
1 batch_size = 256  
2 train_iter, test_iter =  
   d2l.load_data_fashion_mnist(batch_size=batch_size)  
3  
4 lr, num_epochs = 0.001, 5  
5 optimizer = torch.optim.Adam(net.parameters(), lr=lr)  
6 d2l.train_ch5(net, train_iter, test_iter, batch_size, optimizer,  
   device, num_epochs)
```

输出：

```
1 training on  cuda  
2 epoch 1, loss 0.0054, train acc 0.767, test acc 0.795, time 2.0 sec  
3 epoch 2, loss 0.0024, train acc 0.851, test acc 0.748, time 2.0 sec  
4 epoch 3, loss 0.0017, train acc 0.872, test acc 0.814, time 2.2 sec  
5 epoch 4, loss 0.0014, train acc 0.883, test acc 0.818, time 2.1 sec  
6 epoch 5, loss 0.0013, train acc 0.889, test acc 0.734, time 1.8 sec
```

## 小结

- 在模型训练时，批量归一化利用小批量上的均值和标准差，不断调整神经网络的中间输出，从而使整个神经网络在各层的中间输出的数值更稳定。
- 对全连接层和卷积层做批量归一化的方法稍有不同。
- 批量归一化层和丢弃层一样，在训练模式和预测模式的计算结果是不一样的。
- PyTorch提供了BatchNorm类方便使用。

## 参考文献

(1) Ioffe, S., & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. arXiv preprint arXiv:1502.03167.

---

注：除代码外本节与原书此节基本相同，[原书传送门](#)

## 5.11 残差网络 (RESNET)

让我们先思考一个问题：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) (1)。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。

### 5.11.2 残差块

让我们聚焦于神经网络局部。如图5.9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5.9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  $f(x) - x$ 。残差映射在实际中往往更容易优化。以本节开头提到的恒等映射作为我们希望学出的理想映射  $f(x)$ 。我们只需将图5.9中右图虚线框内上方的加权运算（如仿射）的权重和偏差参数设为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。图5.9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。



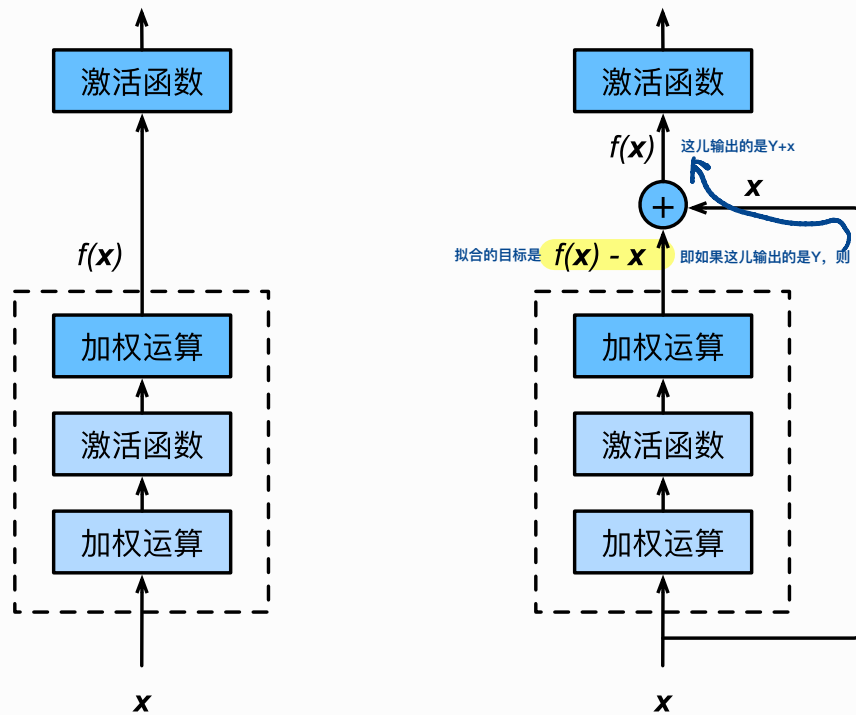


图5.9 普通的网络结构（左）与加入残差连接的网络结构（右）

ResNet沿用了VGG全 $3 \times 3$ 卷积层的设计。残差块里首先有2个有相同输出通道数的 $3 \times 3$ 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 $1 \times 1$ 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 $1 \times 1$ 卷积层来修改通道数以及卷积层的步幅。

```

1  import time
2  import torch
3  from torch import nn, optim
4  import torch.nn.functional as F
5
6  import sys
7  sys.path.append("..")
8  import d2lzh_pytorch as d2l
9  device = torch.device('cuda' if torch.cuda.is_available() else
    'cpu')
10
11 class Residual(nn.Module): # 本类已保存在d2lzh_pytorch包中方便以后使用
12     def __init__(self, in_channels, out_channels,
13 use_1x1conv=False, stride=1):
14         super(Residual, self).__init__()
15         self.conv1 = nn.Conv2d(in_channels, out_channels,
16 kernel_size=3, padding=1, stride=stride)
17         self.conv2 = nn.Conv2d(out_channels, out_channels,
18 kernel_size=3, padding=1)
19         if use_1x1conv:
20             self.conv3 = nn.Conv2d(in_channels, out_channels,
21 kernel_size=1, stride=stride)

```



```

18         else:
19             self.conv3 = None
20             self.bn1 = nn.BatchNorm2d(out_channels)
21             self.bn2 = nn.BatchNorm2d(out_channels)
22
23     重点看forward方法: def forward(self, X):
24         Y = F.relu(self.bn1(self.conv1(X)))
25         Y = self.bn2(self.conv2(Y)) # 这一步得出的Y是一个结果，如果没有残差，Y就可以直接作为结果。但现在需要把X的形状变成Y，因为
26         if self.conv3:                                     残差块最终返回的是Y+x
27             X = self.conv3(X) # 改变X的通道数，与Y相同
28         return F.relu(Y + X)

```

下面我们来查看输入和输出形状一致的情况。

```

1 blk = Residual(3, 3) # 输入通道和输出通道相同
2 X = torch.rand((4, 3, 6, 6))
3 blk(X).shape # torch.Size([4, 3, 6, 6])

```

我们也可以在增加输出通道数的同时减半输出的高和宽。

```

1 blk = Residual(3, 6, use_1x1conv=True, stride=2) # 当输入和输出通道数不同时，就需要1*1全连接层改变通道数
2 blk(X).shape # torch.Size([4, 6, 3, 3])

```

## 5.11.2 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 $7 \times 7$ 卷积层后接步幅为2的 $3 \times 3$ 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。

```

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

```

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。之后的每个模块在第一个残差块里将上一个模块的通道数翻倍，并将高和宽减半。

下面我们来实现这个模块。注意，这里对第一个模块做了特别处理。

```

1  def resnet_block(in_channels, out_channels, num_residuals,
    first_block=False):
2      if first_block:
3          assert in_channels == out_channels # 第一个模块的通道数同输入通道数一致
4      blk = []
5      for i in range(num_residuals):
6          if i == 0 and not first_block: # 若为各block中的第0个残差块
7              blk.append(Residual(in_channels, out_channels,
    use_1x1conv=True, stride=2))
8          else: # 在非第0个残差块中, 保证残差块的输入和输出通道都是out_channels
9              blk.append(Residual(out_channels, out_channels))
10     return nn.Sequential(*blk)

```

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。

添加4个残差block

```

1  net.add_module("resnet_block1", resnet_block(64, 64, 2,
    first_block=True))
2  net.add_module("resnet_block2", resnet_block(64, 128, 2))
3  net.add_module("resnet_block3", resnet_block(128, 256, 2))
4  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。

```

1  net.add_module("global_avg_pool", d2l.GlobalAvgPool2d()) #
    GlobalAvgPool2d的输出: (Batch, 512, 1, 1)
2  net.add_module("fc", nn.Sequential(d2l.FlattenLayer(),
    nn.Linear(512, 10)))

```

这里每个模块里有4个卷积层（不计算 $1 \times 1$ 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。

在训练ResNet之前，我们来观察一下输入形状在ResNet不同模块之间的变化。

输入一个例子：

```

1  X = torch.rand((1, 1, 224, 224))
2  for name, layer in net.named_children():
3      X = layer(X)
4      print(name, ' output shape:\t', X.shape)

```

输出：

每一层的大小输出:

```
1 0 output shape:      torch.Size([1, 64, 112, 112])
2 1 output shape:      torch.Size([1, 64, 112, 112])
3 2 output shape:      torch.Size([1, 64, 112, 112])
4 3 output shape:      torch.Size([1, 64, 56, 56])
5 resnet_block1 output shape:      torch.Size([1, 64, 56, 56])
6 resnet_block2 output shape:      torch.Size([1, 128, 28, 28])
7 resnet_block3 output shape:      torch.Size([1, 256, 14, 14])
8 resnet_block4 output shape:      torch.Size([1, 512, 7, 7])
9 global_avg_pool output shape:    torch.Size([1, 512, 1, 1])
10 fc output shape:     torch.Size([1, 10])
```

### 5.11.3 获取数据和训练模型

下面我们在Fashion-MNIST数据集上训练ResNet。

```
1 batch_size = 256
2 # 如出现“out of memory”的报错信息，可减小batch_size或resize
3 train_iter, test_iter = d2l.load_data_fashion_mnist(batch_size,
4               resize=96)
5
6 lr, num_epochs = 0.001, 5
7 optimizer = torch.optim.Adam(net.parameters(), lr=lr)
8 d2l.train_ch5(net, train_iter, test_iter, batch_size, optimizer,
9               device, num_epochs)
```

输出:

```
1 training on  cuda
2 epoch 1, loss 0.0015, train acc 0.853, test acc 0.885, time 31.0 sec
3 epoch 2, loss 0.0010, train acc 0.910, test acc 0.899, time 31.8 sec
4 epoch 3, loss 0.0008, train acc 0.926, test acc 0.911, time 31.6 sec
5 epoch 4, loss 0.0007, train acc 0.936, test acc 0.916, time 31.8 sec
6 epoch 5, loss 0.0006, train acc 0.944, test acc 0.926, time 31.5 sec
```

## 小结

- 残差块通过跨层的数据通道从而能够训练出有效的深度神经网络。
- ResNet深刻影响了后来的深度神经网络的设计。

## 参考文献

(1) He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 770-778).

(2) He, K., Zhang, X., Ren, S., & Sun, J. (2016, October). Identity mappings in deep residual networks. In European Conference on Computer Vision (pp. 630-645). Springer, Cham.

注：除代码外本节与原书此节基本相同，[原书传送门](#)

## 5.12 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。本节我们介绍其中的一个：稠密连接网络 (DenseNet) (1)。它与ResNet的主要区别如图5.10所示。

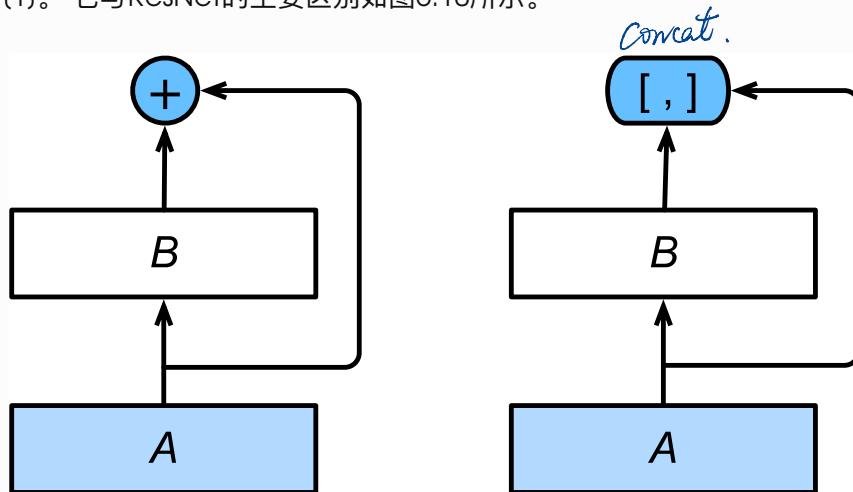


图5.10 ResNet (左) 与DenseNet (右) 在跨层连接上的主要区别：使用相加和使用连结

图5.10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。因为A和B不断连接，那么通道数会越来越大。也就是上右图看到的部分

### 5.12.1 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。

```

1 import time
2 import torch
3 from torch import nn, optim
4 import torch.nn.functional as F
5
6 import sys
7 sys.path.append("..")
8 import d2lzh_pytorch as d2l
9 device = torch.device('cuda' if torch.cuda.is_available() else
10 'cpu')
11
12 def conv_block(in_channels, out_channels): # 这个函数将BN、ReLU、Conv2d打包起来
13     blk = nn.Sequential(nn.BatchNorm2d(in_channels),
14                         nn.ReLU(),
15                         nn.Conv2d(in_channels, out_channels,
16                                 kernel_size=3, padding=1))
17     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。

```

1 class DenseBlock(nn.Module):
2     def __init__(self, num_convs, in_channels, out_channels):
3         super(DenseBlock, self).__init__()
4         net = []
5         for i in range(num_convs):
6             in_c = in_channels + i * out_channels
7             net.append(conv_block(in_c, out_channels))
8         self.net = nn.ModuleList(net)
9         self.out_channels = in_channels + num_convs * out_channels
10
11     # 计算输出通道数
12
13     def forward(self, X):
14         for blk in self.net:
15             Y = blk(X) # 每个X都经过一个卷积层
16             X = torch.cat((X, Y), dim=1) # 在通道维上将输入和输出连结
17         return X

```

在该DenseBlock中用到了几个conv\_block 整个DenseBlock输入的通道数 每次卷积之后的输出通道数，即上图B的输出通道数

这一层DenseBlock的输出将作为下一层DenseBlock的输入。自己计算一下，会发现下一层的输出通道是在本层输出的基础之上再加上B卷积层的输出out\_channels，变成了in\_c+2\*out\_c，如此循环。

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）。

```

1 blk = DenseBlock(2, 3, 10) # 依据上一块代码的右图的解释，这儿输出通道数为2*10+3=23
2 X = torch.randn(4, 3, 8, 8) # batch_size 通道数 高 宽
3 Y = blk(X)
4 Y.shape # torch.Size([4, 23, 8, 8]) # 没有过度层的输出通道数是23

```

## 5.12.2 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 $1 \times 1$ 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。

```
1 def transition_block(in_channels, out_channels):
2     blk = nn.Sequential(
3         nn.BatchNorm2d(in_channels),
4         nn.ReLU(),
5         nn.Conv2d(in_channels, out_channels, kernel_size=1),
6         nn.AvgPool2d(kernel_size=2, stride=2)) # 减半高、宽
7     return blk
```

对上一个例子中稠密块的输出使用通道数为10的过渡层。此时输出的通道数减为10，高和宽均减半。

```
1 blk = transition_block(接上上块未接过渡层的输出通道23, 10)
2 blk(Y).shape # torch.Size([4, 10, 4, 4])
```

## 5.12.3 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。

```
1 net = nn.Sequential(
2     nn.Conv2d(1, 64, kernel_size=7, 高、宽减半stride=2, padding=3),
3     nn.BatchNorm2d(64),
4     nn.ReLU(),
5     nn.MaxPool2d(kernel_size=3, 高、宽再减半stride=2, padding=1))
```

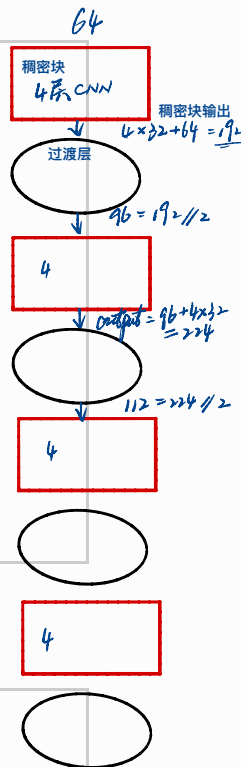
类似于ResNet接下来使用的4个残差块，DenseNet使用的是4个稠密块。同ResNet一样，我们可以设置每个稠密块使用多少个卷积层。这里我们设成4，从而与上一节的ResNet-18保持一致。稠密块里的卷积层通道数（即增长率）设为32，所以每个稠密块将增加128个通道。

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。

```

1 num_channels, growth_rate = 64, 32 # num_channels为当前的通道数
2 num_convs_in_dense_blocks = [4, 4, 4, 4] # 共4个稠密块，每个稠密块中含4个卷积层
3
4 for i, num_convs in enumerate(num_convs_in_dense_blocks):
5     DB = DenseBlock(num_convs, num_channels, growth_rate)
6     net.add_module("DenseBlosk_%d" % i, DB)
7     # 上一个稠密块的输出通道数 直接调用DB实例的一个属性 获取输出通道数
8     num_channels = DB.out_channels
9     # 在稠密块之间加入通道数减半的过渡层
10    if i != len(num_convs_in_dense_blocks) - 1:
11        net.add_module("transition_block_%d" % i,
12                        transition_block(num_channels, num_channels // 2))
13        num_channels = num_channels // 2 # 更新num_channels

```



同ResNet一样，最后接上全局池化层和全连接层来输出。

```

1 net.add_module("BN", nn.BatchNorm2d(num_channels))
2 net.add_module("relu", nn.ReLU())
3 net.add_module("global_avg_pool", d2l.GlobalAvgPool2d()) #
  GlobalAvgPool2d的输出: (Batch, num_channels, 1, 1)
4 net.add_module("fc", nn.Sequential(d2l.FlattenLayer(),
  nn.Linear(num_channels, 10)))

```

我们尝试打印每个子模块的输出维度确保网络无误：

```

输入: X = torch.rand((1, 1, 96, 96))
2 for name, layer in net.named_children():
3     X = layer(X)
4     print(name, ' output shape:\t', X.shape)

```

输出：

```

1 0 output shape: torch.Size([1, 64, 48, 48]) # conv2d让高、宽减半 (48=96//2)
2 1 output shape: torch.Size([1, 64, 48, 48]) # Batch_Norm和ReLU不改变高、宽
3 2 output shape: torch.Size([1, 64, 48, 48])
4 3 output shape: torch.Size([1, 64, 24, 24]) # MaxPool后高、宽再减半
5 DenseBlosk_0 output shape: torch.Size([1, 192, 24, 24]) # 经过第一个DenseBlock, 通道数成192
6 transition_block_0 output shape: torch.Size([1, 96, 12, 12])
7 DenseBlosk_1 output shape: torch.Size([1, 224, 12, 12])
8 transition_block_1 output shape: torch.Size([1, 112, 6, 6])
9 DenseBlosk_2 output shape: torch.Size([1, 240, 6, 6])
10 transition_block_2 output shape: torch.Size([1, 120, 3, 3])
11 DenseBlosk_3 output shape: torch.Size([1, 248, 3, 3])
12 BN output shape: torch.Size([1, 248, 3, 3])
13 relu output shape: torch.Size([1, 248, 3, 3])
14 global_avg_pool output shape: torch.Size([1, 248, 1, 1])
15 fc output shape: torch.Size([1, 10])

```



## 5.12.4 获取数据并训练模型

由于这里使用了比较深的网络，本节里我们将输入高和宽从224降到96来简化计算。

```
1 batch_size = 256
2 # 如出现“out of memory”的报错信息，可减小batch_size或resize
3 train_iter, test_iter = d2l.load_data_fashion_mnist(batch_size,
4               resize=96)
5 lr, num_epochs = 0.001, 5
6 optimizer = torch.optim.Adam(net.parameters(), lr=lr)
7 d2l.train_ch5(net, train_iter, test_iter, batch_size, optimizer,
8               device, num_epochs)
```

输出：

```
1 training on  cuda
2 epoch 1, loss 0.0020, train acc 0.834, test acc 0.749, time 27.7 sec
3 epoch 2, loss 0.0011, train acc 0.900, test acc 0.824, time 25.5 sec
4 epoch 3, loss 0.0009, train acc 0.913, test acc 0.839, time 23.8 sec
5 epoch 4, loss 0.0008, train acc 0.921, test acc 0.889, time 24.9 sec
6 epoch 5, loss 0.0008, train acc 0.929, test acc 0.884, time 24.3 sec
```

## 小结

- 在跨层连接上，不同于ResNet中将输入与输出相加，DenseNet在通道维上连结输入与输出。
- DenseNet的主要构建模块是稠密块和过渡层。

## 参考文献

(1) Huang, G., Liu, Z., Weinberger, K. Q., & van der Maaten, L. (2017). Densely connected convolutional networks. In Proceedings of the IEEE conference on computer vision and pattern recognition (Vol. 1, No. 2).

---

注：除代码外本节与原书此节基本相同，[原书传送门](#)