

批量规范化

训练深层神经网络是十分困难的，特别是在较短的时间内使他们收敛更加棘手。在本节中，我们将介绍 *批量规范化*（batch normalization），这是一种流行且有效的技术，可持续加速深层网络的收敛速度。再结合在之后在 ResNet 中将介绍的残差块，批量规范化使得研究人员能够训练100层以上的网络。

训练深层网络

为什么需要批量规范化层呢？

首先，数据预处理的方式通常会对最终结果产生巨大影响。回想一下我们应用多层感知机来预测房价的例子。使用真实数据时，我们的第一步是标准化输入特征，使其平均值为0，方差为1。直观地说，这种标准化可以很好地与我们的优化器配合使用，因为它可以将参数的量级进行统一。

第二，对于典型的多层感知机或卷积神经网络。当我们训练时，中间层中的变量（例如，多层感知机中的仿射变换输出）可能具有更广的变化范围：不论是沿着从输入到输出的层，跨同一层中的单元，或是随着时间的推移，模型参数的随着训练更新变幻莫测。批量规范化的发明者非正式地假设，这些变量分布中的这种偏移可能会阻碍网络的收敛。直观地说，我们可能会猜想，如果一个层的可变值是另一层的100倍，这可能需要对学习率进行补偿调整。

第三，更深层的网络很复杂，容易过拟合。这意味着正则化变得更加重要。

批量规范化应用于单个可选层（也可以应用到所有层），其原理如下：在每次训练迭代中，我们首先规范化输入，即通过减去其均值并除以其标准差，其中两者均基于当前小批量处理。接下来，我们应用比例系数和比例偏移。正是由于这个基于 *批量* 统计的 *标准化*，才有了 *批量规范化* 的名称。

请注意，如果我们尝试使用大小为1的小批量应用批量规范化，我们将无法学到任何东西。这是因为在减去均值之后，每个隐藏单元将为0。所以，只有使用足够大的小批量，批量规范化这种方法才是有效且稳定的。请注意，在应用批量规范化时，批量大小的选择可能比没有批量规范化时更重要。

从形式上来说，用 $\mathbf{x} \in B$ 表示一个来自小批量 B 的输入，批量规范化BN根据以下表达式转换 \mathbf{x} ：

$$\text{BN}(\mathbf{x}) = \gamma \odot \frac{\mathbf{x} - \hat{\boldsymbol{\mu}}_B}{\hat{\boldsymbol{\sigma}}_B} + \boldsymbol{\beta}.$$

$$\text{BN}(\mathbf{x}) = \gamma \odot \frac{\mathbf{x} - \hat{\boldsymbol{\mu}}_B}{\hat{\boldsymbol{\sigma}}_B} + \boldsymbol{\beta}.$$

在上式中， $\hat{\boldsymbol{\mu}}_B$ 是小批量 B 的样本均值， $\hat{\boldsymbol{\sigma}}_B$ 是小批量 B 的样本标准差。应用标准化后，生成的小批量的平均值为0和单位方差为1。由于单位方差（与其他一些魔法数）是一个主观的选择，因此我们通常包含 *拉伸参数*（scale） γ 和 *偏移参数*（shift） $\boldsymbol{\beta}$ ，它们的形状与 \mathbf{x} 相同。请注意， γ 和 $\boldsymbol{\beta}$ 是需要与其他模型参数一起学习的参数。

由于在训练过程中，中间层的变化幅度不能过于剧烈，而批量规范化将每一层主动居中，并将它们重新调整为给定的平均值和大小（通过 $\hat{\boldsymbol{\mu}}_B$ 和 $\hat{\boldsymbol{\sigma}}_B$ ）。

从形式上来看，我们计算出 $\hat{\boldsymbol{\mu}}_B$ 和 $\hat{\boldsymbol{\sigma}}_B$ ，如下所示：

$$\begin{aligned}\hat{\boldsymbol{\mu}}_B &= \frac{1}{|B|} \sum_{\mathbf{x} \in B} \mathbf{x}, \\ \hat{\boldsymbol{\sigma}}_B^2 &= \frac{1}{|B|} \sum_{\mathbf{x} \in B} (\mathbf{x} - \hat{\boldsymbol{\mu}}_B)^2 + \epsilon.\end{aligned}$$

请注意，我们在方差估计值中添加一个小的常量 $\epsilon > 0$ ，以确保我们永远不会尝试除以零，即使在经验方差估计值可能消失的情况下也是如此。估计值 $\hat{\mu}_B$ 和 $\hat{\sigma}_B$ 通过使用平均值和方差的噪声（noise）估计来抵消缩放问题。你可能会认为这种噪声是一个问题，而事实上它是有益的。

事实证明，这是深度学习中一个反复出现的主题。由于尚未在理论上明确的原因，优化中的各种噪声源通常会导致更快的训练和较少的过拟合：这种变化似乎是正则化的一种形式。在一些初步研究中有工作将批量规范化的性质与贝叶斯先验相关联。这些理论揭示了为什么批量规范化最适应50 ~ 100范围中的中等批量大小的难题。

另外，批量规范化层在“训练模式”（通过小批量统计数据规范化）和“预测模式”（通过数据集统计规范化）中的功能不同。在训练过程中，我们无法得知使用整个数据集来估计平均值和方差，所以只能根据每个小批次的平均值和方差不断训练模型。而在预测模式下，可以根据整个数据集精确计算批量规范化所需的平均值和方差。

现在，我们了解一下批量规范化在实践中是如何工作的。

批量规范化层

回想一下，批量规范化和其他层之间的一个关键区别是，由于批量规范化在完整的小批量上运行，因此我们不能像以前在引入其他层时那样忽略批量大小。我们在下面讨论这两种情况：全连接层和卷积层，他们的批量规范化实现略有不同。

全连接层

通常，我们将批量规范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} ，权重参数和偏置参数分别为 \mathbf{W} 和 \mathbf{b} ，激活函数为 ϕ ，批量规范化的运算符为BN。那么，使用批量规范化的全连接层的输出的计算详情如下：

$$\mathbf{h} = \phi(\text{BN}(\mathbf{W}\mathbf{x} + \mathbf{b})).$$

回想一下，均值和方差是在应用变换的“相同”小批量上计算的。

卷积层

同样，对于卷积层，我们可以在卷积层之后和非线性激活函数之前应用批量规范化。当卷积有多个输出通道时，我们需要对这些通道的“每个”输出执行批量规范化，每个通道都有自己的拉伸（scale）和偏移（shift）参数，这两个参数都是标量。假设我们的小批量包含 m 个样本，并且对于每个通道，卷积的输出具有高度 p 和宽度 q 。那么对于卷积层，我们在每个输出通道的 $m \cdot p \cdot q$ 个元素上同时执行每个批量规范化。因此，在计算平均值和方差时，我们会收集所有空间位置的值，然后在给定通道内应用相同的均值和方差，以便在每个空间位置对值进行规范化。

预测过程中的批量规范化

正如我们前面提到的，批量规范化在训练模式和预测模式下的行为通常不同。首先，将训练好的模型用于预测时，我们不再需要样本均值中的噪声以及在微批次上估计每个小批次产生的样本方差了。其次，例如，我们可能需要使用我们的模型对逐个样本进行预测。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，和暂退法一样，批量规范化层在训练模式和预测模式下的计算结果也是不一样的。

从零实现

下面，我们从头开始实现一个具有张量的批量规范化层。

```
In [1]: 1 import torch
2 from torch import nn
3 from d2l import torch as d2l
4
5
6 def batch_norm(X, gamma, beta, moving_mean, moving_var, eps, momentum):
7     # 通过is_grad_enabled来判断当前模式是训练模式还是预测模式
8     if not torch.is_grad_enabled():
9         # 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
10        X_hat = (X - moving_mean) / torch.sqrt(moving_var + eps)
11    else:
12        assert len(X.shape) in (2, 4)
13        if len(X.shape) == 2:
14            # 使用全连接层的情况，计算特征维上的均值和方差
15            mean = X.mean(dim=0)
16            var = ((X - mean) ** 2).mean(dim=0)
17        else:
18            # 使用二维卷积层的情况，计算通道维上 (axis=1) 的均值和方差。
19            # 这里我们需要保持X的形状以便后面可以做广播运算
20            mean = X.mean(dim=(0, 2, 3), keepdim=True)
21            var = ((X - mean) ** 2).mean(dim=(0, 2, 3), keepdim=True)
22        # 训练模式下，用当前的均值和方差做标准化
23        X_hat = (X - mean) / torch.sqrt(var + eps)
24        # 更新移动平均的均值和方差
25        moving_mean = momentum * moving_mean + (1.0 - momentum) * mean
26        moving_var = momentum * moving_var + (1.0 - momentum) * var
27    Y = gamma * X_hat + beta # 缩放和移位
28    return Y, moving_mean.data, moving_var.data
```

我们现在可以**创建一个正确的 BatchNorm 层**。这个层将保持适当的参数：拉伸 `gamma` 和偏移 `beta`，这两个参数将在训练过程中更新。此外，我们的层将保存均值和方差的移动平均值，以便在模型预测期间随后使用。

撇开算法细节，注意我们实现层的基础设计模式。通常情况下，我们用一个单独的函数定义其数学原理，比如说 `batch_norm`。然后，我们将此功能集成到一个自定义层中，其代码主要处理数据移动到训练设备（如GPU）、分配和初始化任何必需的变量、跟踪移动平均线（此处为均值和方差）等问题。为了方便起见，我们并不担心在这里自动推断输入形状，因此我们需要指定整个特征的数量。不用担心，深度学习框架中的批量规范化API将为我们解决上述问题，我们稍后将展示这一点。

```
In [2]: 1 class BatchNorm(nn.Module):
2         # num_features: 完全连接层的输出数量或卷积层的输出通道数。
3         # num_dims: 2表示完全连接层，4表示卷积层
4         def __init__(self, num_features, num_dims):
5             super().__init__()
6             if num_dims == 2:
7                 shape = (1, num_features)
8             else:
9                 shape = (1, num_features, 1, 1)
10            # 参与求梯度和迭代的拉伸和偏移参数，分别初始化成1和0
11            self.gamma = nn.Parameter(torch.ones(shape))
12            self.beta = nn.Parameter(torch.zeros(shape))
13            # 非模型参数的变量初始化为0和1
14            self.moving_mean = torch.zeros(shape)
15            self.moving_var = torch.ones(shape)
16
17        def forward(self, X):
18            # 如果X不在内存上，将moving_mean和moving_var
19            # 复制到X所在显存上
20            if self.moving_mean.device != X.device:
21                self.moving_mean = self.moving_mean.to(X.device)
22                self.moving_var = self.moving_var.to(X.device)
23            # 保存更新过的moving_mean和moving_var
24            Y, self.moving_mean, self.moving_var = batch_norm(
25                X, self.gamma, self.beta, self.moving_mean,
26                self.moving_var, eps=1e-5, momentum=0.9)
27            return Y
```

使用批量规范化层的 LeNet

为了更好地理解如何应用 BatchNorm，下面我们将其应用于**LeNet模型**。回想一下，批量规范化是在卷积层或全连接层之后、相应的激活函数之前应用的。

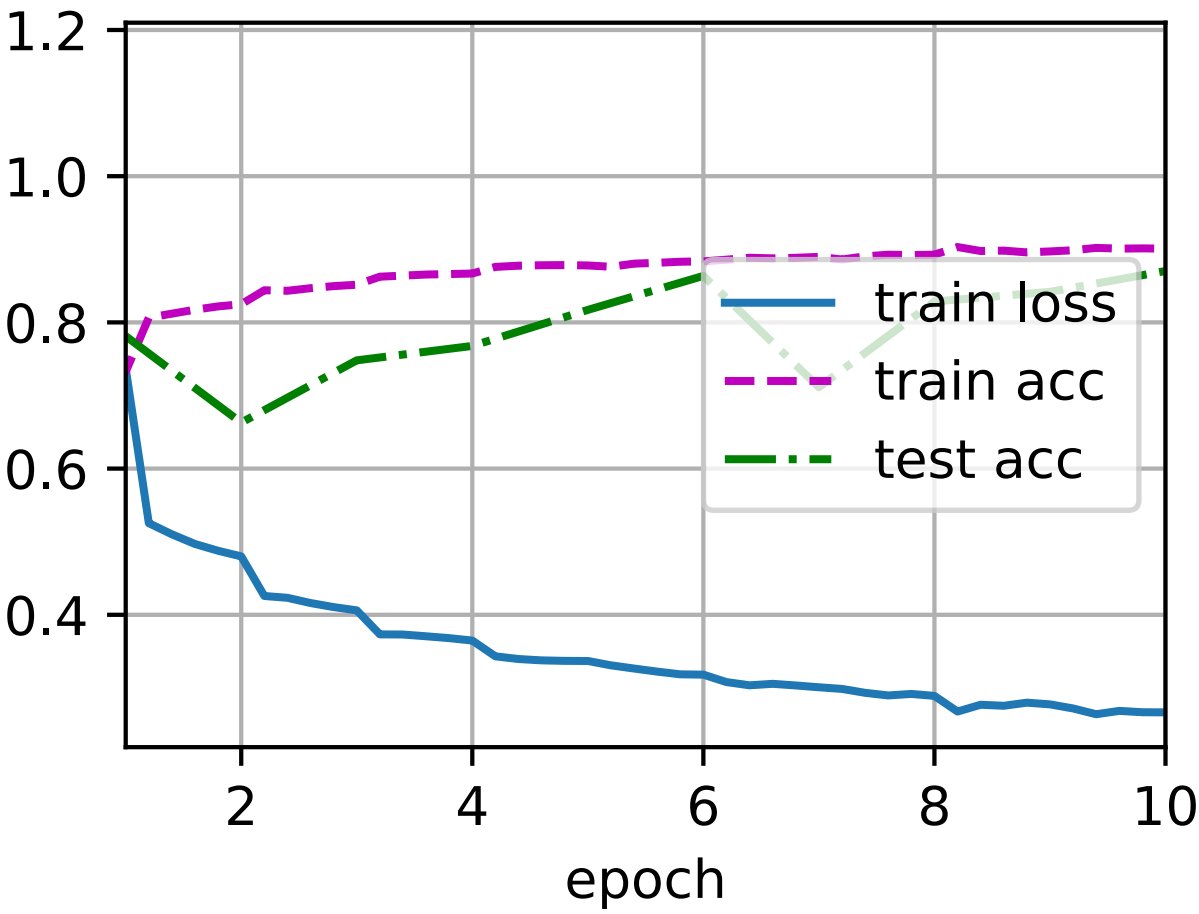
```
In [3]: 1 net = nn.Sequential(
2         nn.Conv2d(1, 6, kernel_size=5), BatchNorm(6, num_dims=4), nn.Sigmoid(),
3         nn.AvgPool2d(kernel_size=2, stride=2),
4         nn.Conv2d(6, 16, kernel_size=5), BatchNorm(16, num_dims=4), nn.Sigmoid(),
5         nn.AvgPool2d(kernel_size=2, stride=2), nn.Flatten(),
6         nn.Linear(16*4*4, 120), BatchNorm(120, num_dims=2), nn.Sigmoid(),
7         nn.Linear(120, 84), BatchNorm(84, num_dims=2), nn.Sigmoid(),
8         nn.Linear(84, 10))
```

和以前一样，我们将在**Fashion-MNIST数据集上训练网络**。这个代码与我们第一次训练LeNet时几乎完全相同，主要区别在于学习率大得多。

In [4]:

```
1 lr, num_epochs, batch_size = 1.0, 10, 256
2 train_iter, test_iter = d2l.load_data_fashion_mnist(batch_size)
3 d2l.train_ch6(net, train_iter, test_iter, num_epochs, lr, d2l.try_gpu())
```

loss 0.267, train acc 0.901, test acc 0.870
31667.5 examples/sec on cuda:0



让我们来看看从第一个批量规范化层中学到的**拉伸参数** `gamma` 和**偏移参数** `beta` 。

In [5]:

```
1 net[1].gamma.reshape((-1,)), net[1].beta.reshape((-1,))
```

Out[5]: (tensor([1.5337, 3.6520, 2.9328, 4.3339, 2.0743, 2.1410], device='cuda:0',
grad_fn=<ReshapeAliasBackward0>),
tensor([2.2010, 3.8538, -1.9424, 2.2748, 0.6700, -1.3545], device='cuda:0',
grad_fn=<ReshapeAliasBackward0>))

简明实现

除了使用我们刚刚定义的 `BatchNorm`，我们也可以直接使用深度学习框架中定义的 `BatchNorm`。该代码看起来几乎与我们上面的代码相同。

In [6]:

```
1 net = nn.Sequential(
2     nn.Conv2d(1, 6, kernel_size=5), nn.BatchNorm2d(6), nn.Sigmoid(),
3     nn.AvgPool2d(kernel_size=2, stride=2),
4     nn.Conv2d(6, 16, kernel_size=5), nn.BatchNorm2d(16), nn.Sigmoid(),
5     nn.AvgPool2d(kernel_size=2, stride=2), nn.Flatten(),
6     nn.Linear(256, 120), nn.BatchNorm1d(120), nn.Sigmoid(),
7     nn.Linear(120, 84), nn.BatchNorm1d(84), nn.Sigmoid(),
8     nn.Linear(84, 10))
```

下面，我们**使用相同超参数来训练模型**。请注意，通常高级API变体运行速度快得多，因为它的代码已编译为C++或CUDA，而我们的自定义代码由Python实现。

In [7]:

1 d2l.train_ch6(net, train_iter, test_iter, num_epochs, lr, d2l.try_gpu())

loss 0.267, train acc 0.902, test acc 0.878
53024.9 examples/sec on cuda:0

