

7.3 小批量随机梯度下降

在每一次迭代中，梯度下降使用整个训练数据集来计算梯度，因此它有时也被称为批量梯度下降（batch gradient descent）。而随机梯度下降在每次迭代中只随机采样一个样本来计算梯度。正如我们在前几章中所看到的，我们还可以在每轮迭代中随机均匀采样多个样本来组成一个小批量，然后使用这个小批量来计算梯度。下面就来描述小批量随机梯度下降。

设目标函数 $f(\mathbf{x}) : \mathbb{R}^d \rightarrow \mathbb{R}$ 。在迭代开始前的时间步设为0。该时间步的自变量记为 $\mathbf{x}_0 \in \mathbb{R}^d$ ，通常由随机初始化得到。在接下来的每一个时间步 $t > 0$ 中，小批量随机梯度下降随机均匀采样一个由训练数据样本索引组成的小批量 B_t 。我们可以通过重复采样（sampling with replacement）或者不重复采样（sampling without replacement）得到一个小批量中的各个样本。前者允许同一个小批量中出现重复的样本，后者则不允许如此，且更常见。对于这两者间的任一种方式，都可以使用

$$\mathbf{g}_t \leftarrow \nabla f_{B_t}(\mathbf{x}_{t-1}) = \frac{1}{|B_t|} \sum_{i \in B_t} \nabla f_i(\mathbf{x}_{t-1})$$

来计算时间步 t 的小批量 B_t 上目标函数位于 \mathbf{x}_{t-1} 处的梯度 \mathbf{g}_t 。这里 $|B_t|$ 代表批量大小，即小批量中样本的个数，是一个超参数。同随机梯度一样，重复采样所得的小批量随机梯度 \mathbf{g}_t 也是对梯度 $\nabla f(\mathbf{x}_{t-1})$ 的无偏估计。给定学习率 η_t （取正数），小批量随机梯度下降对自变量的迭代如下：

$$\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} - \eta_t \mathbf{g}_t$$

基于随机采样得到的梯度的方差在迭代过程中无法减小，因此在实际中，（小批量）随机梯度下降的学习率可以在迭代过程中自我衰减，例如 $\eta_t = \eta t^\alpha$ （通常 $\alpha = -1$ 或者 -0.5 ）、 $\eta_t = \eta \alpha^t$ （如 $\alpha = 0.95$ ）或者每迭代若干次后将学习率衰减一次。如此一来，学习率和（小批量）随机梯度乘积的方差会减小。而梯度下降在迭代过程中一直使用目标函数的真实梯度，无须自我衰减学习率。

小批量随机梯度下降中每次迭代的计算开销为 $O(|B|)$ 。当批量大小为1时，该算法即为随机梯度下降；当批量大小等于训练数据样本数时，该算法即为梯度下降。当批量较小时，每次迭代中使用的样本少，这会导致并行处理和内存使用效率变低。这使得在计算同样数目样本的情况下比使用更大批量时所花时间更多。当批量较大时，每个小批量梯度里可能含有更多的冗余信息。为了得到较好的解，批量较大时比批量较小时需要计算的样本数目可能更多，例如增大迭代周期数。

7.3.1 读取数据

本章里我们将使用一个来自NASA的测试不同飞机机翼噪音的数据集来比较各个优化算法 [1]。我们使用该数据集的前1,500个样本和5个特征，并使用标准化对数据进行预处理。

```
%matplotlib inline
import numpy as np
```

```

import time
import torch
from torch import nn, optim
import sys
sys.path.append("..")
import d2lzh_pytorch as d2l

def get_data_ch7(): # 本函数已保存在d2lzh_pytorch包中方便以后使用
    data = np.genfromtxt('../data/airfoil_self_noise.dat', delimiter='\t')
    data = (data - data.mean(axis=0)) / data.std(axis=0)
    return torch.tensor(data[:1500, :-1], dtype=torch.float32), \
        torch.tensor(data[:1500, -1], dtype=torch.float32) # 前1500个样本 (每个样本5个特征)

features, labels = get_data_ch7()
features.shape # torch.Size([1500, 5])

```

7.3.2 从零开始实现

3.2节（线性回归的从零开始实现）中已经实现过小批量随机梯度下降算法。我们在这里将它的输入参数变得更加通用，主要是为了方便本章后面介绍的其他优化算法也可以使用同样的输入。具体来说，我们添加了一个状态输入 `states` 并将超参数放在字典 `hyperparams` 里。此外，我们将在训练函数里对各个小批量样本的损失求平均，因此优化算法里的梯度不需要除以批量大小。

```

def sgd(params, states, hyperparams):
    for p in params:
        p.data -= hyperparams['lr'] * p.grad.data

```

下面实现一个通用的训练函数，以方便本章后面介绍的其他优化算法使用。它初始化一个线性回归模型，然后可以使用小批量随机梯度下降以及后续小节介绍的其他算法来训练模型。

```

# 本函数已保存在d2lzh_pytorch包中方便以后使用
def train_ch7(optimizer_fn, states, hyperparams, features, labels,
               batch_size=10, num_epochs=2):
    # 初始化模型
    net, loss = d2l.linreg, d2l.squared_loss

    w = torch.nn.Parameter(torch.tensor(np.random.normal(0, 0.01, size=(features.shape[1], 1))))

```

```

requires_grad=True)

b = torch.nn.Parameter(torch.zeros(1, dtype=torch.float32), requires_grad=True)

def eval_loss():
    return loss(net(features, w, b), labels).mean().item()

ls = [eval_loss()]
data_iter = torch.utils.data.DataLoader(
    torch.utils.data.TensorDataset(features, labels), batch_size, shuffle=True)

for _ in range(num_epochs):
    start = time.time()
    for batch_i, (X, y) in enumerate(data_iter):
        l = loss(net(X, w, b), y).mean() # 使用平均损失

        # 梯度清零
        if w.grad is not None:
            w.grad.data.zero_()
            b.grad.data.zero_()

        l.backward()
        optimizer_fn([w, b], states, hyperparams) # 迭代模型参数
        if (batch_i + 1) * batch_size % 100 == 0:
            ls.append(eval_loss()) # 每100个样本记录下当前训练误差

# 打印结果和作图
print('loss: %f, %f sec per epoch' % (ls[-1], time.time() - start))
d2l.set_figsize()
d2l.plt.plot(np.linspace(0, num_epochs, len(ls)), ls)
d2l.plt.xlabel('epoch')
d2l.plt.ylabel('loss')

```

当批量大小为样本总数1,500时，优化使用的是梯度下降。梯度下降的1个迭代周期对模型参数只迭代1次。可以看到6次迭代后目标函数值（训练损失）的下降趋向了平稳。

```

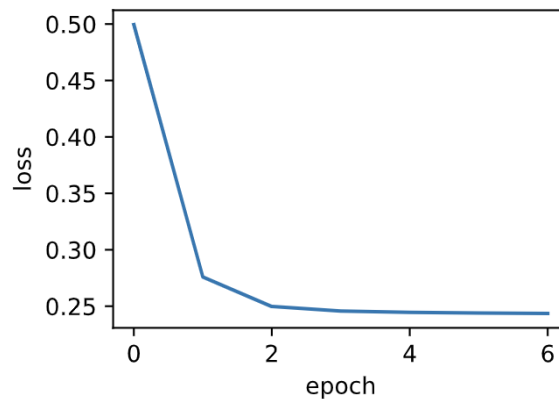
def train_sgd(lr, batch_size, num_epochs=2):
    train_ch7(sgd, None, {'lr': lr}, features, labels, batch_size, num_epochs)

train_sgd(1, 1500, 6)

```

输出:

```
loss: 0.243605, 0.014335 sec per epoch
```

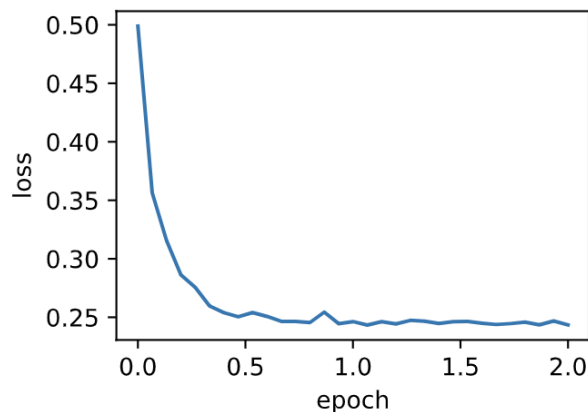


当批量大小为1时，优化使用的是随机梯度下降。为了简化实现，有关（小批量）随机梯度下降的实验中，我们未对学习率进行自我衰减，而是直接采用较小的常数学习率。随机梯度下降中，每处理一个样本会更新一次自变量（模型参数），一个迭代周期里会对自变量进行1,500次更新。可以看到，目标函数值的下降在1个迭代周期后就变得较为平缓。

```
train_sgd(0.005, 1)
```

输出:

```
loss: 0.243433, 0.270011 sec per epoch
```



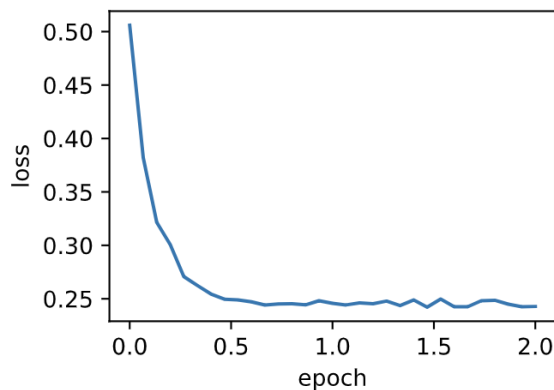
虽然随机梯度下降和梯度下降在一个迭代周期里都处理了1,500个样本，但实验中随机梯度下降的一个迭代周期耗时更多。这是因为随机梯度下降在一个迭代周期里做了更多次的自变量迭代，而且单样本的梯度计算难以有效利用矢量计算。

当批量大小为10时，优化使用的是小批量随机梯度下降。它在每个迭代周期的耗时介于梯度下降和随机梯度下降的耗时之间。

```
train_sgd(0.05, 10)
```

输出：

```
loss: 0.242805, 0.078792 sec per epoch
```



7.3.3 简洁实现

在PyTorch里可以通过创建 `optimizer` 实例来调用优化算法。这能让实现更简洁。下面实现一个通用的训练函数，它通过优化算法的函数 `optimizer_fn` 和超参数 `optimizer_hyperparams` 来创建 `optimizer` 实例。

```
# 本函数与原书不同的是这里第一个参数优化器函数而不是优化器的名字
# 例如: optimizer_fn=torch.optim.SGD, optimizer_hyperparams={"lr": 0.05}
def train_pytorch_ch7(optimizer_fn, optimizer_hyperparams, features, labels,
                      batch_size=10, num_epochs=2):
    # 初始化模型
    net = nn.Sequential(
        nn.Linear(features.shape[-1], 1)
    )
    loss = nn.MSELoss()
    optimizer = optimizer_fn(net.parameters(), **optimizer_hyperparams)

    def eval_loss():
        return loss(net(features).view(-1), labels).item() / 2
```

```

ls = [eval_loss()]
data_iter = torch.utils.data.DataLoader(
    torch.utils.data.TensorDataset(features, labels), batch_size, shuffle=True)

for _ in range(num_epochs):
    start = time.time()
    for batch_i, (X, y) in enumerate(data_iter):
        # 除以2是为了和train_ch7保持一致，因为squared_loss中除了2
        l = loss(net(X).view(-1), y) / 2

        optimizer.zero_grad()
        l.backward()
        optimizer.step()
        if (batch_i + 1) * batch_size % 100 == 0:
            ls.append(eval_loss())
# 打印结果和作图
print('loss: %f, %f sec per epoch' % (ls[-1], time.time() - start))
d2l.set_figsize()
d2l.plt.plot(np.linspace(0, num_epochs, len(ls)), ls)
d2l.plt.xlabel('epoch')
d2l.plt.ylabel('loss')

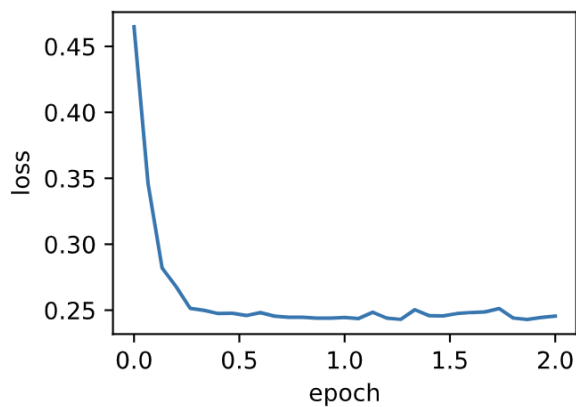
```

使用PyTorch重复上一个实验。

```
train_pytorch_ch7(optim.SGD, {"lr": 0.05}, features, labels, 10)
```

输出:

```
loss: 0.245491, 0.044150 sec per epoch
```



小结

- 小批量随机梯度每次随机均匀采样一个小批量的训练样本来计算梯度。
- 在实际中，（小批量）随机梯度下降的学习率可以在迭代过程中自我衰减。
- 通常，小批量随机梯度在每个迭代周期的耗时介于梯度下降和随机梯度下降的耗时之间。