

# 机器学习中的最优化方法

<https://github.com/w-gc/>

东北大学计算机系

版本:0.1

更新:July 21, 2020

## 1 引言

### 1.1 最优化与机器学习

在机器学习中用到的梯度下降类算法,实际上在最优化中早有提及,而且最优化中不少优化算法也是应用到了机器学习中。那么问题来了,最优化中的算法是不是可以直接拿过来用在机器学习中呢?

显然不是的,两者虽然都是“最优化”但是又有各自面临的问题。它们解决问题的思路是有区别,这导致两者的算法并不是简单的照搬照抄,而是需要做相应的“适应”。

举个例子,假设我们在拟合一条直线,现在有 6 个样本点,见表??。显然,我们需要拟合的模型是  $y = 3 * x + 5$ 。需要说明的是样本的  $y_i$  是不等于  $3 * x_i + 5$ , 因为采样过程中存在不可避免的误差。

表 1: 样本数据

样本序号 $i$	$x_i$	$y_i$
1	0.5000	6.9243
2	1.0000	8.4900
3	1.5000	10.0835
4	2.0000	11.0833
5	2.5000	13.1602
6	3.0000	14.0523

我们首先从解析求解的角度看待这个问题,建立模型  $y = ax + b$ 。求解参数  $a, b$ , 让  $y_i$  和  $ax_i + b$  尽量接近。这里如果我们采用  $l_2$  范式描述两者的接近程度(距离),那么得到的无约束最

优化模型为

$$\min_{a,b} \sum_{i=1}^6 \|ax_i + b - y_i\|_2^2. \quad (1)$$

为了书写的方便记  $L(a, b) = \sum_{i=1}^6 \|ax_i + b - y_i\|_2^2$ 。那么问题便是求解无约束的二元函数的极小值,我们可以用极值点的必要条件来求解。

$$\begin{cases} \frac{\partial L}{\partial a} = \sum_{i=1}^6 2x_i (ax_i + b - y_i) = 0 \\ \frac{\partial L}{\partial b} = \sum_{i=1}^6 2(ax_i + b - y_i) = 0. \end{cases} \quad (2)$$

显然上式可以求得解析解

$$\begin{cases} a = \frac{\sum_{i=1}^6 x_i y_i - \left(\frac{1}{6} \sum_{i=1}^6 y_i\right) \sum_{i=1}^6 x_i}{\sum_{i=1}^6 x_i^2 - \left(\frac{1}{6} \sum_{i=1}^6 x_i\right) \sum_{i=1}^6 y_i} \\ b = \left(\frac{1}{6} \sum_{i=1}^6 y_i\right) - a \left(\frac{1}{6} \sum_{i=1}^6 x_i\right). \end{cases} \quad (3)$$

带入样本点值得到  $y = 2.8943x + 5.5672$ , 见图??。

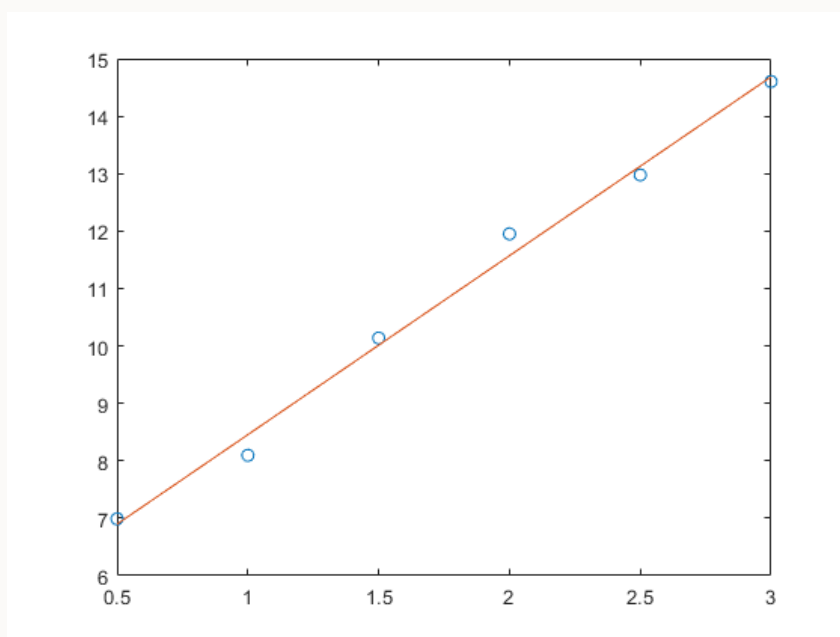


图 1: 采用解析法

显然(??)式解析解可以推广至有限的  $n$  个样本, 只需将式(??)的 6 改成  $n$  即可。上述过程也称作是做一元线性回归, 采用的具体解法称作最小二乘。

```
1 % 采用解析求解法示例代码
2 clear all, clc
3 x = (0.5 : 0.5 : 3)';
4 y = 3 * x + 5 + rand(length(x), 1);
```

```

5 syms a b real
6 y_ = (a*x + b);
7 L = (y_ - y)' * (y_ - y);
8 d_L_a = diff(L, a);
9 d_L_b = diff(L, b);
10 [slove_a, slove_b] = solve([d_L_a, d_L_b]);
11 disp([slove_a, slove_b]);
12 y_ = slove_a * x + slove_b;
13 plot(x, y, 'o', x, y_)

```

当问题变得更复杂,比如样本数目增加,或者参数增加,解析法将很能求解,涉及到了多元非线性方程组。这里我们再考虑最优化教材中常用的梯度下方法,这里以最速下降为例。对于参数  $\theta$  其通用算法为(??)。这里在确定步长因子  $t$  时,显然这个问题比较简单所以其实可以不用近似算法求解,直接整理得到关于  $t$  的一元函数,求解析解。但是当问题复杂时,求步长因子  $t$  的解析解并不简单,所以这里也同样使用近似求解。

$$\begin{cases} t = \arg \max_t f(\theta + t\nabla\theta) \\ \theta = \theta - t\nabla\theta. \end{cases} \quad (4)$$

具体到我们的这个问题其算法为

---

**Algorithm 1** 采用最速下降法

---

**Require:**  $(x_i, y_i), i = 1, \dots, 6$

**Ensure:**  $a, b$

随机初始化  $a, b$

**for**  $i = 1; i < 100; i++$  **do**

$\hat{y}_i = ax_i + b, i = 1, \dots, 6$

$L(a, b) = \sum_{i=1}^6 \|\hat{y}_i - y_i\|_2^2$

$\nabla a = \sum_{i=1}^6 2x_i (\hat{y}_i - y_i)$

$\nabla b = \sum_{i=1}^6 2(\hat{y}_i - y_i)$

寻找  $L(a + t\nabla a, b + t\nabla b)$  包含一个局部极值点的区间  $[l, r]$

在区间  $[l, r]$  选找  $L(a + t\nabla a, b + t\nabla b)$  的极值点  $t$

$a = a + t\nabla a$

$b = b + t\nabla b$

**end for**

---

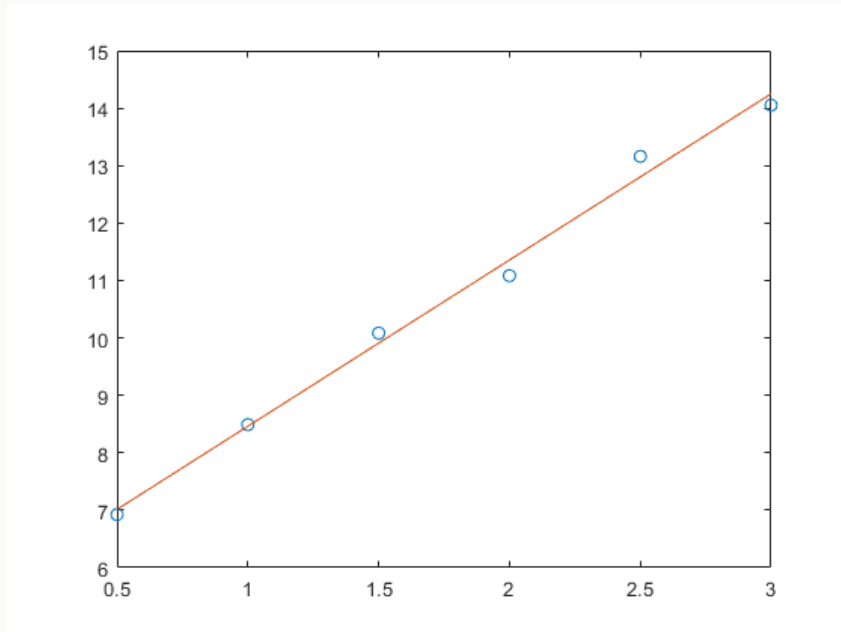


图 2: 采用最速下降法

同样的, 这个算法可以推广至有限个样本  $n$ , 只需要将上面的 6 改成  $n$  即可。参考代码中, 关于关于步长因子  $t$  的局部极值区间和确定步长因子  $t$  的局部极值部分实现, 参考 [?] 。

```

1 % 采用最速下降法示例代码
2 clear all,clc
3 x = (0.5 : 0.5 : 3)';
4 % y = 3 * x + 5 + rand(length(x),1);
5 y = [6.9243, 8.4900, 10.0835, 11.0833, 13.1602, 14.0523]';
6 a = rand(1);
7 b = rand(1);
8 for i = 1 : 100
9     y_ = (a*x + b);
10    L = (y_ - y)' * (y_ - y);
11    d_L_a = 45.5 * a + 21* b - 248.6028;
12    d_L_b = 21*a + 12*b - 127.5872;
13    %-----确定关于步长因子t的局部极值区间
14    ft = @(t)((a - t * d_L_a)*x + (b - t * d_L_b)) - y)' * ((a - t * d_L_a)*x + (b - t *
        d_L_b)) - y);
15    [InterL, InterR] = searchInterval(ft);
16    %-----确定步长因子t的局部极值
17    t = bestT(ft, InterL, InterR);
18    %-----更新参数a,b
19    a = (a - t * d_L_a);
20    b = (b - t * d_L_b);
21 end
22 disp([a, b]);
23 y_ = a * x + b;
24 plot(x, y, 'o', x, y_)

```

最后,我们来看看如果是从神经网络的角度怎么解决这个问题。将其看成一个非常简单的神经网络,见图??。神经网络的优化使用到了一个叫学习率,其实和最速下降法中的步长因子是同一事物,都是控制沿着梯度方向变化的程度。

```

1 % 采用机器学习方法示例代码
2 clear all,clc
3 x = (0.5 : 0.5 : 3)';
4 % y = 3 * x + 5 + rand(length(x),1);
5 y = [6.9243, 8.4900, 10.0835, 11.0833, 13.1602, 14.0523]';
6 a = rand(1);
7 b = rand(1);
8 lr = 0.0001;
9 for i = 1:1000
10     y_ = a * x - b;
11     e = y_ - y;
12     d_L_a = 2 * x' * e / length(x);
13     d_L_b = 2 * mean(e);
14     a = a - lr * d_L_a;
15     b = b - lr * d_L_b;
16 end
17 disp([a, b]);
18 y_ = a * x + b;
19 plot(x, y, 'o', x, y_)

```

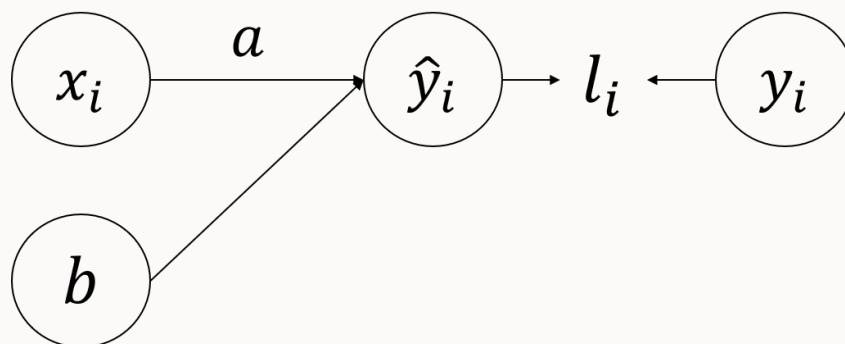


图 3: 看成一个简单的神经网络

观察图??,发现在这个简单的问题上,这个简单的神经网络的表现并没有好过传统的最优化方法,而且还非常不稳定。

总得来说,在以前的最优化问题中,可以采用解析法,或者梯度下降方法。虽然神经网络的优化算法也是梯度下降算法,但是两者还是有些区别。主要是关于步长因子 $t$ ,最优化教材中的算法往往是采用求解 $t$ 的最优步长,这里其实上是需要整理到关于 $t$ 的函数。但是在深度的网络中,虽然只有一个参数 $t$ 但要整理得到 $t$ 函数也是有困难的。而且还需要搜索一个包含 $t$ 的局部极小值区间,然后再在这个区间里搜索最优值。而深度网络的优化或者说机器学习中的优化,一般是将其设置为一个很小的数值,比如 $10^{-4}$ 。

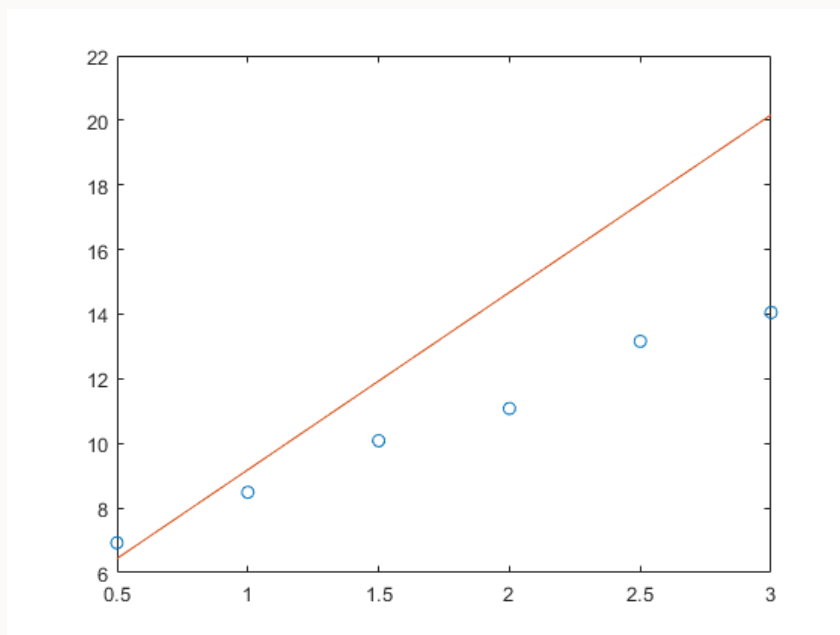


图 4: 采用 ML 的算法

## 1.2 一个简单的分类网络

CIFAR-10 数据集（加拿大高级研究所）是一种常用于训练图像的采集机器学习和计算机视觉算法。它是机器学习研究中使用最广泛的数据集之一。CIFAR-10 数据集包含 10 种不同类别的 60,000 张 32x32 彩色图像。10 个不同的类别分别代表飞机, 汽车, 鸟类, 猫, 鹿, 狗, 青蛙, 马, 轮船和卡车。每个类别有 6,000 张图像 [?]。

这里采用一个简单卷积神经网络：

### 网络结构

```

1 Net(
2   (conv1): Conv2d(3, 32, kernel_size=(5, 5), stride=(1, 1), padding=(2, 2))
3   (maxpool): MaxPool2d(kernel_size=(3, 3), stride=2, padding=(1, 1), dilation=1, ceil_mode=False)
4   (conv2): Conv2d(32, 32, kernel_size=(5, 5), stride=(1, 1), padding=(2, 2))
5   (avgpool): AvgPool2d(kernel_size=(3, 3), stride=2, padding=(1, 1))
6   (conv3): Conv2d(32, 64, kernel_size=(5, 5), stride=(1, 1), padding=(2, 2))
7   (conv4): Conv2d(64, 64, kernel_size=(4, 4), stride=(1, 1))
8   (conv5): Conv2d(64, 10, kernel_size=(1, 1), stride=(1, 1))
9 )

```

整个运行程序为：

```

1 # -*- coding: utf-8 -*-
2 import random
3 import torch
4 import torchvision
5 import torchvision.transforms as transforms
6 import torch.nn as nn
7 import torch.nn.functional as F

```

```

8 import torch.optim as optim
9 # 导入不同的优化器
10 # from SVRG import SVRG
11
12 batch_size = 100
13 EPOCH = 10
14
15 transform = transforms.Compose([transforms.ToTensor(),
16                                 transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))])
17
18 trainset = torchvision.datasets.CIFAR10(root='/sdb/wgc/cifar10', train=True,
19                                         download=True, transform=transform)
20 trainloader = torch.utils.data.DataLoader(trainset, batch_size = batch_size,
21                                           shuffle=True, num_workers=2)
22
23 testset = torchvision.datasets.CIFAR10(root='/sdb/wgc/cifar10', train=False,
24                                         download=True, transform=transform)
25 testloader = torch.utils.data.DataLoader(testset, batch_size = batch_size,
26                                           shuffle=False, num_workers=2)
27
28 classes = ('plane', 'car', 'bird', 'cat', 'deer',
29            'dog', 'frog', 'horse', 'ship', 'truck')
30
31 # dataiter = iter(trainloader)
32 # inputs, labels = dataiter.next()
33
34 class Net(nn.Module):
35     def __init__(self):
36         super(Net, self).__init__()
37         self.conv1 = nn.Conv2d(in_channels = 3, out_channels = 32, kernel_size = 5,
38                                 stride = 1, padding = (2, 2))
39         self.maxpool = nn.MaxPool2d(kernel_size = (3, 3), stride = 2, padding=(1, 1))
40         self.conv2 = nn.Conv2d(in_channels = 32, out_channels = 32, kernel_size = 5,
41                                 stride = 1, padding = (2, 2))
42         self.avgpool = nn.AvgPool2d(kernel_size = (3, 3), stride = 2, padding=(1, 1))
43         self.conv3 = nn.Conv2d(in_channels = 32, out_channels = 64, kernel_size = 5,
44                                 stride = 1, padding = (2, 2))
45         self.conv4 = nn.Conv2d(in_channels = 64, out_channels = 64, kernel_size = 4,
46                                 stride = 1)
47         self.conv5 = nn.Conv2d(in_channels = 64, out_channels = 10, kernel_size = 1,
48                                 stride = 1)
49
50     def forward(self, x):
51         batch_size = x.size()[0]
52         x = F.relu(self.maxpool(self.conv1(x)))
53         x = self.avgpool(F.relu(self.conv2(x)))
54         x = self.avgpool(F.relu(self.conv3(x)))
55         x = F.relu(self.conv4(x))
56         x = self.conv5(x)

```

```

57         x = x.view(batch_size, -1)
58         return x
59
60 net = Net()
61
62 use_gpu = torch.cuda.is_available()
63 if(use_gpu):
64     net = net.cuda()
65
66 criterion = nn.CrossEntropyLoss()
67
68 # 选择不同优化器
69 # optimizer = optim.SGD(net.parameters(), lr = 1e-2)
70
71 for epoch in range(EPOCH): # loop over the dataset multiple times
72     running_loss = 0.0
73     for i, data in enumerate(trainloader, 0):
74         # get the inputs
75         inputs, labels = data
76         if(use_gpu):
77             inputs = inputs.cuda()
78             labels = labels.cuda()
79         outputs = net(inputs)
80         loss = criterion(outputs, labels)
81         optimizer.zero_grad()
82
83         loss.backward()
84         optimizer.step()
85
86         # print statistics
87         running_loss += loss.item()
88         if i % 100 == 99: # print every 100 mini-batches
89             print('[%d, %5d] loss: %.3f' % (epoch + 1, i + 1, running_loss / 2000))
90             running_loss = 0.0
91
92 # 一个EPOCH训练的准确率
93 correct = 0
94 total = 0
95 with torch.no_grad():
96     for data in testloader:
97         inputs, labels = data
98         if(use_gpu):
99             inputs = inputs.cuda()
100             labels = labels.cuda()
101         outputs = net(inputs)
102         # outputs = F.softmax(outputs)
103         _, predicted = torch.max(outputs.data, 1)
104         total += labels.size(0)
105         correct += (predicted == labels).sum().item()

```



```

106     print('%d Accuracy of the network on the 10000 test images: %d %%' % (epoch + 1, 100 *
107         correct / total))
108
109
110 print('-----Finished Training-----')
111
112 # 训练结束后在测试集上测试，并输出各个类别的准确率
113 class_correct = list(0. for i in range(10))
114 class_total = list(0. for i in range(10))
115 with torch.no_grad():
116     for data in testloader:
117         inputs, labels = data
118         if(use_gpu):
119             inputs = inputs.cuda()
120             labels = labels.cuda()
121         outputs = net(inputs)
122         _, predicted = torch.max(outputs, 1)
123         c = (predicted == labels).squeeze()
124         for i in range(4):
125             label = labels[i]
126             class_correct[label] += c[i].item()
127             class_total[label] += 1
128
129 for i in range(10):
130     print('Accuracy of %5s : %2d %%' % (classes[i], 100 * class_correct[i] / class_total[i]))

```

### 1.3 优化方法主要类别

机器学习中的优化方法主要按导数类型,分为三大类:

- 一阶导信息 GD,SGD,NAG,AdaDelta/RMSProp,Adam,SAG,SVRG,ADMM,Frank-Wolfe
- 高阶导信息共轭梯度法、(拟)牛顿、Natural Gradient
- 非导数启发式,采样法

## 2 机器学习中常见优化模型

### 2.1 有监督学习

有一组样本  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ 。定义  $f(x, \theta)$  为机器学习模型（可以是回归，也可以是神经网络）。那么定义损失函数为  $l(y, f)$ ，这里的损失可以是 log-softmax，也可以是某种范式。那么建立的优化模型为

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l_i(y_i, f(x_i, \theta)). \quad (5)$$

我们得到的这个目标函数(??)对于确定的样本，其实是只跟参数  $\theta$  有关。而参数  $\theta$  是我们需要”学习“的，它的值确定了我们的机器学习模型。 $\theta$  可以是回归模型里的参数  $a, b$ ，有可以是神经网络里的连接权重  $w$ ，当然也可以是卷积神经(CNN)中的卷积核。

为了防止过拟合，我们通常会在目标函数后面增加关于参数  $\theta$  的惩罚项。比如在高维统计学中常用的  $L_1, L_2$  范式惩罚项，当然，在机器学习中也是常用  $L_1, L_2$  范式惩罚项。这里为了描述问题的更一般性，假设采用了  $L_p$  范式。则

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N l_i(y_i, f(x_i, \theta)) + \lambda \|\theta\|_{L_p}^2. \quad (6)$$

### 2.2 半监督学习

半监督学习是，部分数据有标签，部分数据没有标签。这里我们仅考虑 S3VM。假设有带标签数据  $(x_1, y_1), (x_1, y_1), \dots, (x_l, y_l)$ ，和无标签数据  $(x_{l+1}, y_{l+1}), (x_{l+2}, y_{l+2}), \dots, (x_{l+u}, y_{l+u})$ ，且  $N = l + u$ 。

$$\begin{cases} \min \|w\| + C \left[ \sum_{i=1}^l \zeta^i + \sum_{j=l+1}^N \min(\epsilon_j^0, \epsilon_j^1) \right] \\ s.t. y_i (wx_i + b) + \zeta^i \geq 1, \quad i = 1, 2, \dots, l \\ wx_j + b + \epsilon_j^0 \geq 1, \quad \epsilon_j^0 \geq 0, \quad j = l+1, l+2, \dots, N \\ -(wx_j + b) + \epsilon_j^1 \geq 1, \quad \epsilon_j^1 \geq 0, \quad j = l+1, l+2, \dots, N \end{cases} \quad (7)$$

其中  $\zeta^i$  是松弛变量， $C$  是一个错误率的惩罚因子，而  $\epsilon_j^0, \epsilon_j^1$  分别是样本  $x_j$  分到 0、1 类的错分率。

### 2.3 无监督学习

这里主要介绍下 k-means、主成分分析(PCA)、和极大似然估计(MLE)的最优化模型。首先关于 k-means，假设有  $K$  个类别，第  $k$  类的样本集为  $S_k$ ，类中心为  $\mu_k$ 。那么我们的目标就是找到

$K$ , 让每个样本集中的样本  $x \in S_k$ , 尽量接近类中心  $\mu_k$ 。

$$\min_S \sum_{k=1}^K \sum_{x \in S_k} \|x - \mu_k\|_{L_p}^2. \quad (8)$$

PCA, 将原来的样本从原来的样本空间投影至一个更小的空间里, 目标是寻找能够使得投影前后两者相差最小的投影方式(映射)。

$$\min \sum_{i=1}^N \left\| \sum_{j=1}^{D'} z_{i,j} e_j - x_i \right\|_{L_p}^2, \quad D \gg D'. \quad (9)$$

其中,  $e_j$  是小空间的基, 而  $z_i$  是样本  $x_i$  在小空间上的投影坐标, 而  $\sum_{j=1}^{D'} z_{i,j} e_j$  是对样本  $x_i$  的重构。跟自编码器一样, 先编码, 再解码, 解码结果与真实结果尽量接近。

极大似然估计, 思想是让寻找参数让出现现在有的样本的情况最大。也就是说, 样本联合概率最大。

$$\max \sum_{i=1}^N \ln p(x_i; \theta); \quad (10)$$

## 2.4 强化学习

此领域不是我关心的重点, 这里只简单介绍下。

$$\max_{\pi} V_{\pi}(s) = \mathbb{E} \left( \sum_{k=0}^{\infty} \gamma^k r_{t+k} | S_t = s \right). \quad (11)$$

(??)的意思是, 在时刻  $t$ , 所在状态为  $s$ , 采用策略  $\pi$  获得的回报是  $r$ ,  $\gamma \in [0, 1]$  是折扣率, 目标是选择最佳的策略  $\pi$  使得获得的回报是  $r$  最大。

## 3 一阶导信息

能使函数值减小的都可以称之为下降方向，而负梯度方向则是下降最快的方向。当然，需要说明的是，这里的所谓的“下降最快”是指在当前点的局部领域内下降最快。否则会以为直接用负梯度就是最好的，最速下降是最快的，实则不是。

### 3.1 梯度下降

这整个章节使用一阶导数信息的下降方法都可以说是梯度下降方法，但这里单独的这一小章节主要是讨论随机梯度、小批量、批量梯度下降。这三者的区别是由样本数目决定的，依次是一个样本、小批量样本、整个样本集合，也就是批量的大小 `batch-size`。`batch-size` 的取值在这里可以看出对参数更新策略是有影响的，但它的其他影响已经如何选择 `batch-size` 我们在之后的章节专讨论。

#### 3.1.1 随机梯度下降

随机梯度下降是指将样本集打乱后，每个样本  $(x_i, y_i)$  有优化模型

$$\min_{\theta} l_i(y_i, f(x_i, \theta)). \quad (12)$$

更新策略为

$$\theta \leftarrow \theta - \eta \frac{\partial l_i}{\partial \theta}. \quad (13)$$

```
1 import torch
2 from torch.optim.optimizer import Optimizer, required
3 class SGD(Optimizer):
4     """
5     (batch, mini-batch)Stochastic gradient descent:
6     $$
7     \backslash\theta \backslash\text{gets} \backslash\theta - \backslash\eta \backslash\text{sum}_{\{i=1\}^{\{B\}} \backslash\frac{\{\partial\} l_i}{\{\partial\} \backslash\theta}.
8     $$
9     """
10    def __init__(self, params, lr=required, weight_decay=0):
11        if lr is not required and lr < 0.0:
12            raise ValueError("Invalid learning rate: {}".format(lr))
13        if weight_decay < 0.0:
14            raise ValueError("Invalid weight_decay value: {}".format(weight_decay))
15        defaults = dict(lr=lr, weight_decay=weight_decay)
16        super(SGD, self).__init__(params, defaults)
17
18    def __setstate__(self, state):
19        super(SGD, self).__setstate__(state)
20        # for group in self.param_groups:
21            # group.setdefault('nesterov', False)
```

```

22
23 def step(self, closure=None):
24     """Performs a single optimization step.
25     Arguments:
26         closure (callable, optional): A closure that reevaluates the model
27         and returns the loss.
28     """
29     loss = None
30     if closure is not None:
31         loss = closure()
32     for group in self.param_groups:
33         weight_decay = group['weight_decay']
34         for p in group['params']:
35             if p.grad is None:
36                 continue
37             d_p = p.grad.data
38             if weight_decay != 0:
39                 d_p.add_(weight_decay, p.data)
40             p.data.add_(-group['lr'], d_p)
41     return loss

```

### 3.1.2 批量梯度下降

对整个样本集  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ 。优化模型??便是批量梯度下降的优化模型, 而它的更新策略为

$$\theta \leftarrow \theta - \eta \sum_{i=1}^N \frac{\partial l_i}{\partial \theta}. \quad (14)$$

### 3.1.3 小批量梯度下降

小批量梯度下降是介于随机梯度下降和批量梯度下降中间的一种方法。对于一小批样本  $(x_1, y_1), (x_2, y_2), \dots, (x_B, y_B)$ 。  $B$  是批量的大小 (样本的个数), 一般  $N \gg B$ 。建立的最优模型为

$$\min_{\theta} \frac{1}{B} \sum_{i=1}^B l_i(y_i, f(x_i, \theta)). \quad (15)$$

更新策略为

$$\theta \leftarrow \theta - \eta \sum_{i=1}^B \frac{\partial l_i}{\partial \theta}. \quad (16)$$

## 3.2 动量法

动量法的启发源自物理, 想象一个小球从山顶上滚落下来, 当它到达谷底时, 从数学的角度来看此时的梯度时为 0 的, 但是小球有惯性会继续想前。受这种启发, 有了动量法, 它的好处是能

够更好的摆脱局部最小值。

$$\begin{cases} v \leftarrow \gamma v + \eta \frac{\partial L}{\partial \theta} \\ \theta \leftarrow \theta - v. \end{cases} \quad (17)$$

其中  $L$  是模型的损失函数,  $\gamma$  通常取值为 0.9。

可以看到,所谓的“动量”其实是梯度的历史值的累计,这个算法说明了梯度的历史信息也是有用的,有助于收敛。当然,记录梯度的历史信息是需要额外的存储空间,这个空间的大小跟参数的大小一致。

下面的代码实现本质来说也是一种小批量梯度下降方法,也是使用了批量的技术。

```
1 import torch
2 from torch.optim.optimizer import Optimizer, required
3
4 class MSGD(Optimizer):
5     """
6     Momentum Stochastic gradient descent:
7     $$
8     \left\{ \begin{aligned} & v \leftarrow \gamma v + \eta \frac{\partial L}{\partial \theta} \\ & \theta \leftarrow \theta - v. \end{aligned} \right.
9
10
11     \end{aligned}
12
13     \right.
14     $$
15     """
16
17     def __init__(self, params, lr=required, momentum=0.9, dampening=0, weight_decay=0):
18         if lr is not required and lr < 0.0:
19             raise ValueError("Invalid learning rate: {}".format(lr))
20         if momentum < 0.0:
21             raise ValueError("Invalid momentum value: {}".format(momentum))
22         if weight_decay < 0.0:
23             raise ValueError("Invalid weight_decay value: {}".format(weight_decay))
24
25         defaults = dict(lr=lr, momentum=momentum, dampening=dampening, weight_decay=
26             weight_decay)
27         super(MSGD, self).__init__(params, defaults)
28
29     def __setstate__(self, state):
30         super(MSGD, self).__setstate__(state)
31
32     def step(self, closure=None):
33         """Performs a single optimization step.
34         Arguments:
35             closure (callable, optional): A closure that reevaluates the model
36             and returns the loss.
37         """
```

```

37     loss = None
38     if closure is not None:
39         loss = closure()
40
41     for group in self.param_groups:
42         for p in group['params']:
43             if p.grad is None:
44                 continue
45             d_p = p.grad.data
46             if group['weight_decay'] != 0:
47                 d_p.add_(group['weight_decay'], p.data)
48             if group['momentum'] != 0:
49                 param_state = self.state[p]
50                 if 'momentum_buffer' not in param_state:
51                     param_state['momentum_buffer'] = torch.clone(d_p).detach()
52                 else:
53                     param_state['momentum_buffer'].mul_(group['momentum']).add_(1 - group['
                        dampening'], d_p)
54                 d_p = param_state['momentum_buffer']
55             p.data.add_(-group['lr'], d_p)
56     return loss

```

### 3.3 NAG(Nesterov Accelerated Gradient Descent)

NAG 方比动量法具有更快的收敛速度,更新策略为

$$\begin{cases} v \leftarrow \gamma v + \eta \frac{\partial L}{\partial (\theta - \gamma v)} \\ \theta \leftarrow \theta - v. \end{cases} \quad (18)$$

观察图??和图??,两个最重要的差别是  $\frac{\partial L}{\partial \theta}$  和  $\frac{\partial L}{\partial (\theta - \gamma v)}$ 。用 Nesterov[?] 本人的解释是用了下一步  $(\theta - \gamma v)$  的梯度,见图。当然,后来有人[?]证明上式(??)与下式等价:

$$\begin{cases} v_t \leftarrow \gamma v_{t-1} + \eta \frac{\partial L}{\partial \theta_{t-1}} + \gamma \left( \frac{\partial L}{\partial \theta_{t-1}} - \frac{\partial L}{\partial \theta_{t-2}} \right) \\ \theta_t \leftarrow \theta_{t-1} - \gamma v_t. \end{cases} \quad (19)$$

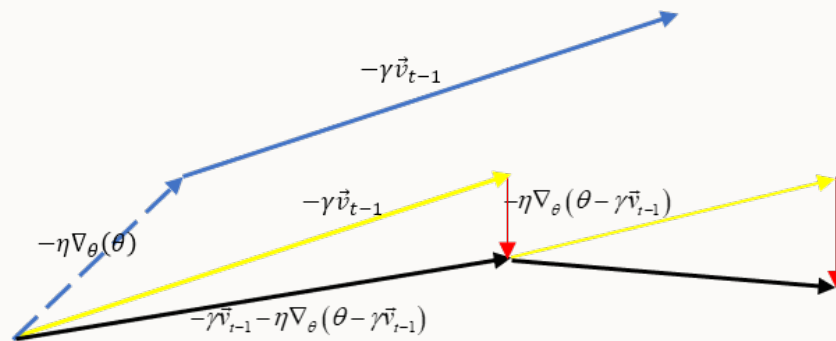


图 5: NAG 法更新

NAG 等价性证明.

$$\begin{aligned}
\theta_t - \gamma v_t &= (\theta_{t-1} - v_t) - \gamma v_t \\
&= \theta_{t-1} - (1 + \gamma) v_t \\
&= \theta_{t-1} - (1 + \gamma) [\gamma v_{t-1} + \eta \nabla_{\theta} L(\theta - \gamma v_{t-1})] \\
&= \theta_{t-1} - \gamma v_{t-1} - \gamma^2 v_{t-1} - (1 + \gamma) \eta \nabla_{\theta} L(\theta - \gamma v_{t-1}).
\end{aligned} \tag{20}$$

记  $\hat{\theta}_{t-1} = \theta_{t-1} - \gamma v_{t-1}$ ,  $\hat{v}_t = \gamma^2 v_{t-1} + (1 + \gamma) \eta \nabla_{\theta} L(\hat{\theta}_{t-1})$ . 另外

$$\begin{aligned}
\hat{v}_t &= \gamma^2 v_{t-1} + (1 + \gamma) \eta \nabla_{\theta} L(\hat{\theta}_{t-1}) \\
&= \gamma^2 [\gamma v_{t-2} + \eta \nabla_{\theta} L(\hat{\theta}_{t-2})] + (1 + \gamma) \eta \nabla_{\theta} L(\hat{\theta}_{t-1}) \\
&= \gamma^3 v_{t-2} + \eta \gamma^2 \nabla_{\theta} L(\hat{\theta}_{t-2}) + \eta \gamma \nabla_{\theta} L(\hat{\theta}_{t-1}) + \eta \nabla_{\theta} L(\hat{\theta}_{t-1}) \\
&= \gamma^t v_1 + \eta \sum_{i=1}^{t-1} \gamma^i \nabla_{\theta} L(\hat{\theta}_{t-i}) + \eta \nabla_{\theta} L(\hat{\theta}_{t-1})
\end{aligned} \tag{21}$$

因此得到  $\hat{v}_t - \gamma \hat{v}_{t-1} = \eta \nabla_{\theta} L(\hat{\theta}_{t-1}) + \gamma [\nabla_{\theta} L(\hat{\theta}_{t-1}) - \nabla_{\theta} L(\hat{\theta}_{t-2})]$ . 整理后便有

$$\begin{aligned}
\hat{v}_t &= \gamma \hat{v}_{t-1} + \eta \nabla_{\theta} L(\hat{\theta}_{t-1}) + \gamma [\nabla_{\theta} L(\hat{\theta}_{t-1}) - \nabla_{\theta} L(\hat{\theta}_{t-2})] \\
\hat{\theta}_t &= \hat{\theta}_{t-1} - \gamma \hat{v}_t
\end{aligned}$$

```

1 import torch
2 from torch.optim.optimizer import Optimizer, required
3
4 class Nesterov(Optimizer):
5     """
6     Nesterov Accelerated Gradient Descent
7     $$
8     \left\{
9     \begin{aligned}
10     v &\leftarrow \gamma v + \eta \frac{\partial L}{\partial \left( \theta - \gamma v \right)} \\
11     \theta &\leftarrow \theta - v.
12     \end{aligned}
13     \right.
14     $$
15     """
16
17     def __init__(self, params, lr=required, momentum=0.9, weight_decay=0):
18         if lr is not required and lr < 0.0:
19             raise ValueError("Invalid learning rate: {}".format(lr))
20         if momentum <= 0.0:
21             raise ValueError("Invalid momentum value: {}".format(momentum))
22         if weight_decay < 0.0:
23             raise ValueError("Invalid weight_decay value: {}".format(weight_decay))
24
25         defaults = dict(lr=lr, momentum=momentum, weight_decay=weight_decay)

```



```

26
27     super(Nesterov, self).__init__(params, defaults)
28
29 def __setstate__(self, state):
30     super(Nesterov, self).__setstate__(state)
31
32 def step(self, closure=None):
33     """Performs a single optimization step.
34     Arguments:
35         closure (callable, optional): A closure that reevaluates the model
36         and returns the loss.
37     """
38     loss = None
39     if closure is not None:
40         loss = closure()
41
42     for group in self.param_groups:
43         weight_decay = group['weight_decay']
44         momentum = group['momentum']
45
46         for p in group['params']:
47             if p.grad is None:
48                 continue
49             d_p = p.grad.data
50             if weight_decay != 0:
51                 d_p.add_(weight_decay, p.data)
52             if momentum != 0:
53                 param_state = self.state[p]
54                 if 'momentum_buffer' not in param_state:
55                     param_state['momentum_buffer'] = torch.clone(d_p).detach()
56                 else:
57                     param_state['momentum_buffer'].mul_(momentum).add_(d_p)
58                 d_p = d_p.add(momentum, param_state['momentum_buffer'])
59             p.data.add_(-group['lr'], d_p)
60     return loss

```

### 3.4 AdaGrad

AdaGrad 较少使用了,因为本来是为了用历史梯度信息来获得自适应的学习率,但是  $\frac{\eta}{V_t}$  是单调递减的。

$$\begin{cases} g_t \leftarrow \frac{\partial L(\theta_t)}{\partial \theta} \\ V_t \leftarrow \sqrt{\sum_{i=0}^t (g_i)^2} + \epsilon \\ \theta \leftarrow \theta - \frac{\eta}{V_t} g_t \end{cases} \quad (22)$$

需要存储一个和参数一样大小的  $V_t$  和临时梯度  $\frac{\partial L(\theta_t)}{\partial \theta}$ 。

因为  $\frac{\eta}{V_t}$  是单调递减的, 这里比之其他方法需要比较大的学习率  $\eta$ , 否则会过早梯度消失。

```

1 import torch
2 from torch.optim.optimizer import Optimizer, required
3
4 class AdaGrad(Optimizer):
5     """AdaGrad algorithm.
6     $$
7     \left\{
8     \begin{aligned}
9         g_t &\leftarrow \frac{\partial L(\theta_t)}{\partial \theta} \\
10        V_t &\leftarrow \sqrt{\sum_{i=0}^t \left\| g_t \right\|^2} + \epsilon \\
11        \theta &\leftarrow \theta - \frac{\eta}{V_t} g_t
12    \end{aligned}
13    \right.
14    $$
15    """
16
17    def __init__(self, params, lr=1e-2, lr_decay=0, weight_decay=0, initial_accumulator_value
18                =0, eps=1e-10):
19        if not 0.0 <= lr:
20            raise ValueError("Invalid learning rate: {}".format(lr))
21        if not 0.0 <= lr_decay:
22            raise ValueError("Invalid lr_decay value: {}".format(lr_decay))
23        if not 0.0 <= weight_decay:
24            raise ValueError("Invalid weight_decay value: {}".format(weight_decay))
25        if not 0.0 <= initial_accumulator_value:
26            raise ValueError("Invalid initial_accumulator_value value: {}".format(
27                initial_accumulator_value))
28        if not 0.0 <= eps:
29            raise ValueError("Invalid epsilon value: {}".format(eps))
30
31        defaults = dict(lr=lr, lr_decay=lr_decay, eps=eps, weight_decay=weight_decay,
32                        initial_accumulator_value=initial_accumulator_value)
33        super(AdaGrad, self).__init__(params, defaults)
34
35        for group in self.param_groups:
36            for p in group['params']:
37                state = self.state[p]
38                state['step'] = 0
39                state['sum'] = torch.full_like(p.data, initial_accumulator_value)
40
41    def __setstate__(self, state):
42        super(AdaGrad, self).__setstate__(state)
43
44    def step(self, closure=None):
45        """Performs a single optimization step.
46        Arguments:

```

```

45         closure (callable, optional): A closure that reevaluates the model
46         and returns the loss.
47
48     """
49     loss = None
50     if closure is not None:
51         loss = closure()
52     for group in self.param_groups:
53         for p in group['params']:
54             if p.grad is None:
55                 continue
56             grad = p.grad.data
57             if grad.is_sparse:
58                 raise RuntimeError('Adadelata does not support sparse gradients')
59             state = self.state[p]
60             state['step'] += 1
61             if group['weight_decay'] != 0:
62                 grad = grad.add(group['weight_decay'], p.data)
63             clr = group['lr'] / ( 1 + (state['step'] - 1) * group['lr_decay'] )
64             state['sum'].addcmul_(1, grad, grad)
65             std = state['sum'].sqrt().add_(group['eps'])
66             p.data.addcdiv_(-clr, grad, std)
67     return loss

```

### 3.5 AdaDelta/RMSProp

$$\left\{ \begin{array}{l} g_t \leftarrow \frac{\partial L(\theta_t)}{\partial \theta} \\ E[g^2]_t \leftarrow \gamma E[g^2]_{t-1} + (1 - \gamma)(g_t)^2 \\ \Delta \theta_t \leftarrow \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t \\ \theta_{t+1} \leftarrow \theta_t - \Delta \theta_t \end{array} \right. \quad (23)$$

需要存储一个和参数一样大小的  $E[g^2]_t$ ,  $\Delta \theta_t$  和临时梯度  $\frac{\partial L(\theta_t)}{\partial \theta}$ 。超参数  $\gamma$  通常取值为 0.9。

上式(??)中的学习率  $\eta$  也可以改成自适应的。总共需要 4 倍的和参数  $\theta$  大小的内存。

$$\left\{ \begin{array}{l} g_t \leftarrow \frac{\partial L(\theta_t)}{\partial \theta} \\ E[g^2]_t \leftarrow \gamma E[g^2]_{t-1} + (1 - \gamma)(g_t)^2 \\ \Delta \theta_t \leftarrow \frac{\sqrt{E[(\Delta \theta_t)^2]_{t-1} + \epsilon}}{\sqrt{E[g^2]_t + \epsilon}} g_t \\ E[(\Delta \theta_t)^2]_t \leftarrow \gamma E[(\Delta \theta_t)^2]_{t-1} + (1 - \gamma)(\Delta \theta_t)^2 \\ \theta_{t+1} \leftarrow \theta_t - \Delta \theta_t \end{array} \right. \quad (24)$$

式(??)中的  $\gamma$  取值为 0.9, 则为 RMSProp 算法。

$$\left\{ \begin{array}{l} g_t \leftarrow \frac{\partial L(\theta_t)}{\partial \theta} \\ E[g^2]_t \leftarrow 0.9E[g^2]_{t-1} + 0.1(g_t)^2 \\ \Delta\theta_t \leftarrow \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t \\ \theta_{t+1} \leftarrow \theta_t - \Delta\theta_t \end{array} \right. \quad (25)$$

这里, 我们使用的是学习率自适应法。学习率不自适应的时候, 表现非常出色。而学习率自适应的时候, 表现得发挥不稳定, 有时候表现极好非常快就收敛了, 有时候却始终不更新权重(也算是提前收敛了)。这是跟  $E[g^2]_t, E[(\Delta\theta_t)^2]_t$  的初始化有关系。

```

1 import torch
2 from torch.optim.optimizer import Optimizer, required
3
4 class AdaDelta(Optimizer):
5     """Adadelta algorithm.
6     $$
7     \left\{ \begin{array}{l}
8     \begin{aligned}
9         g_t &\leftarrow \frac{\partial L(\theta_t)}{\partial \theta} \\
10        E[g^2]_t &\leftarrow \gamma E[g^2]_{t-1} + (1 - \gamma) (g_t)^2 \\
11        \Delta\theta_t &\leftarrow \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t \\
12        \theta_{t+1} &\leftarrow \theta_t - \Delta\theta_t
13      \end{aligned}
14    \end{array} \right.
15    $$
16    """
17
18    def __init__(self, params, lr=1.0, rho=0.9, eps=1e-6, weight_decay=0):
19        if not 0.0 <= lr:
20            raise ValueError("Invalid learning rate: {}".format(lr))
21        if not 0.0 <= rho <= 1.0:
22            raise ValueError("Invalid rho value: {}".format(rho))
23        if not 0.0 <= eps:
24            raise ValueError("Invalid epsilon value: {}".format(eps))
25        if not 0.0 <= weight_decay:
26            raise ValueError("Invalid weight_decay value: {}".format(weight_decay))
27
28        defaults = dict(lr=lr, rho=rho, eps=eps, weight_decay=weight_decay)
29        super(AdaDelta, self).__init__(params, defaults)
30
31        for group in self.param_groups:
32            for p in group['params']:
33                state = self.state[p]
34                state['square_avg'] = torch.zeros_like(p.data)
35                if group['lr'] == 1:
36                    state['acc_delta'] = torch.zeros_like(p.data)

```

```

37
38 def step(self, closure=None):
39     """Performs a single optimization step.
40     Arguments:
41         closure (callable, optional): A closure that reevaluates the model
42         and returns the loss.
43     """
44     loss = None
45     if closure is not None:
46         loss = closure()
47
48     for group in self.param_groups:
49         for p in group['params']:
50             if p.grad is None:
51                 continue
52             grad = p.grad.data
53             if grad.is_sparse:
54                 raise RuntimeError('Adadelta does not support sparse gradients')
55             state = self.state[p]
56
57
58             rho, eps = group['rho'], group['eps']
59
60             if group['weight_decay'] != 0:
61                 grad = grad.add(group['weight_decay'], p.data)
62
63             state['square_avg'].mul_(rho).addcmul_(1 - rho, grad, grad)
64             std = state['square_avg'].add(eps).sqrt_()
65
66             if group['lr'] == 1:
67                 delta = state['acc_delta'].add(eps).sqrt_().div_(std).mul_(grad)
68                 state['acc_delta'].mul_(rho).addcmul_(1 - rho, delta, delta)
69             else:
70                 delta = 1 / std * grad #1.div_(std).mul_(grad)
71
72             p.data.add_(-group['lr'], delta)
73
74     return loss

```

### 3.6 Adam

Adam 是目前最为常用的算法,但是它需要 4 倍的和参数  $\theta$  大小的内存,同时拥有 3 个超参。通常取  $\beta_1$  为 0.9,  $\beta_2$  为 0.999, 而  $\epsilon$  取值为  $10^{-8}$ 。

$$\left\{ \begin{array}{l} g_t \leftarrow \frac{\partial L(\theta_t)}{\partial \theta} \\ m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t \\ v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) (g_t)^2 \\ \hat{m}_t \leftarrow \frac{m_t}{1 - (\beta_1)^t} \\ \hat{v}_t \leftarrow \frac{v_t}{1 - (\beta_2)^t} \\ \theta \leftarrow \theta - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t. \end{array} \right. \quad (26)$$

```

1 import math
2 import torch
3 from torch.optim.optimizer import Optimizer, required
4
5 class Adam(Optimizer):
6     """Adam algorithm.
7     $$
8     \left\{ \begin{array}{l}
9         \begin{aligned}
10             g_t &\gets \frac{\partial L(\theta_t)}{\partial \theta} \\
11             m_t &\gets \beta_1 m_{t-1} + (1 - \beta_1) g_t \\
12             v_t &\gets \beta_2 v_{t-1} + (1 - \beta_2) (g_t)^2 \\
13             \hat{m}_t &\gets \frac{m_t}{1 - (\beta_1)^t} \\
14             \hat{v}_t &\gets \frac{v_t}{1 - (\beta_2)^t} \\
15             \theta &\gets \theta - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t.
16         \end{aligned}
17     \end{array} \right.
18     $$
19     """
20
21     def __init__(self, params, lr=1e-3, betas=(0.9, 0.999), eps=1e-8, weight_decay=0, amsgrad=False):
22         if not 0.0 <= lr:
23             raise ValueError("Invalid learning rate: {}".format(lr))
24         if not 0.0 <= eps:
25             raise ValueError("Invalid epsilon value: {}".format(eps))
26         if not 0.0 <= betas[0] < 1.0:
27             raise ValueError("Invalid beta parameter at index 0: {}".format(betas[0]))
28         if not 0.0 <= betas[1] < 1.0:
29             raise ValueError("Invalid beta parameter at index 1: {}".format(betas[1]))
30         defaults = dict(lr=lr, betas=betas, eps=eps, weight_decay=weight_decay, amsgrad=amsgrad)
31         super(Adam, self).__init__(params, defaults)
32
33         # State initialization
34         for group in self.param_groups:
35             for p in group['params']:
36                 state = self.state[p]

```

```

37         state['step'] = 0
38         state['exp_avg'] = torch.zeros_like(p.data)
39         state['exp_avg_sq'] = torch.zeros_like(p.data)
40
41     def __setstate__(self, state):
42         super(Adam, self).__setstate__(state)
43
44     def step(self, closure=None):
45         """Performs a single optimization step.
46         Arguments:
47             closure (callable, optional): A closure that reevaluates the model
48             and returns the loss.
49         """
50         loss = None
51         if closure is not None:
52             loss = closure()
53
54         for group in self.param_groups:
55             for p in group['params']:
56                 if p.grad is None:
57                     continue
58                 grad = p.grad.data
59                 if grad.is_sparse:
60                     raise RuntimeError('Adam does not support sparse gradients')
61                 amsgrad = group['amsgrad']
62
63                 state = self.state[p]
64
65                 beta1, beta2 = group['betas']
66                 state['step'] += 1
67                 bias_correction1 = 1 - beta1 ** state['step']
68                 bias_correction2 = 1 - beta2 ** state['step']
69
70                 if group['weight_decay'] != 0:
71                     grad.add_(group['weight_decay'], p.data)
72
73                 state['exp_avg'].mul_(beta1).add_(1 - beta1, grad)
74                 state['exp_avg_sq'].mul_(beta2).addcmul_(1 - beta2, grad, grad)
75                 denom = (state['exp_avg_sq'].sqrt() / math.sqrt(bias_correction2)).add_(group['eps'])
76                 exp_avg_hat = state['exp_avg'] / bias_correction1
77
78                 p.data.addcdiv_(-group['lr'], exp_avg_hat, denom)
79         return loss

```

### 3.7 Nadam

Nadam 是采用了 NAG 的思想,对式(??)进行了调整,将  $\hat{m}_{t-1}$  改成了“下一步”的。

$$\begin{cases} g_t \leftarrow \frac{\partial L(\theta_t)}{\partial \theta} \\ m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t \\ v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) (g_t)^2 \\ \hat{m}_t \leftarrow \frac{m_t}{1 - (\beta_1)^t} \\ \hat{v}_t \leftarrow \frac{v_t}{1 - (\beta_2)^t} \\ \theta \leftarrow \theta - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \left( \beta_1 \hat{m}_t + \frac{1 - \beta_1}{1 - (\beta_1)^t} g_t \right). \end{cases} \quad (27)$$

简要推导下式(??)和式(??)的关系

$$\begin{aligned} \theta_{t+1} &= \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t \\ &= \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \left( \frac{m_t}{1 - (\beta_1)^t} \right) \\ &= \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \left( \frac{\beta_1 m_{t-1} + (1 - \beta_1) g_t}{1 - (\beta_1)^t} \right) \\ &= \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \left( \frac{\beta_1 m_{t-1}}{1 - (\beta_1)^t} + \frac{(1 - \beta_1) g_t}{1 - (\beta_1)^t} \right) \\ &= \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \left( \beta_1 \hat{m}_{t-1} + \frac{(1 - \beta_1) g_t}{1 - (\beta_1)^t} \right) \end{aligned} \quad (28)$$

将  $\hat{m}_{t-1}$  改最新的  $\hat{m}_t$ , 得到

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \left( \beta_1 \hat{m}_t + \frac{(1 - \beta_1) g_t}{1 - (\beta_1)^t} \right) \quad (29)$$

```

1 import math
2 import torch
3 from torch.optim.optimizer import Optimizer, required
4
5 class Nadam(Optimizer):
6     """Nadam algorithm.
7     $$
8     \left\{ \begin{aligned}
9         g_t &\gets \frac{\partial L(\theta_t)}{\partial \theta} \\
10        m_t &\gets \beta_1 m_{t-1} + (1 - \beta_1) g_t \\
11        v_t &\gets \beta_2 v_{t-1} + (1 - \beta_2) (g_t)^2 \\
12        \hat{m}_t &\gets \frac{m_t}{1 - (\beta_1)^t} \\
13        \hat{v}_t &\gets \frac{v_t}{1 - (\beta_2)^t} \\
14        \theta &\gets \theta - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \left( \beta_1 \hat{m}_t + \frac{1 - \beta_1}{1 - (\beta_1)^t} g_t \right).
15    \end{aligned} \right.
16
17 \end{aligned}

```



```

18 \right.
19 $$
20 """
21
22 def __init__(self, params, lr=1e-3, betas=(0.9, 0.999), eps=1e-8, weight_decay=0, amsgrad=
    False):
23     if not 0.0 <= lr:
24         raise ValueError("Invalid learning rate: {}".format(lr))
25     if not 0.0 <= eps:
26         raise ValueError("Invalid epsilon value: {}".format(eps))
27     if not 0.0 <= betas[0] < 1.0:
28         raise ValueError("Invalid beta parameter at index 0: {}".format(betas[0]))
29     if not 0.0 <= betas[1] < 1.0:
30         raise ValueError("Invalid beta parameter at index 1: {}".format(betas[1]))
31     defaults = dict(lr=lr, betas=betas, eps=eps, weight_decay=weight_decay, amsgrad=amsgrad
        )
32     super(Nadam, self).__init__(params, defaults)
33
34     # State initialization
35     for group in self.param_groups:
36         for p in group['params']:
37             state = self.state[p]
38             state['step'] = 0
39             state['exp_avg'] = torch.zeros_like(p.data)
40             state['exp_avg_sq'] = torch.zeros_like(p.data)
41
42     def __setstate__(self, state):
43         super(Nadam, self).__setstate__(state)
44
45     def step(self, closure=None):
46         """Performs a single optimization step.
47         Arguments:
48             closure (callable, optional): A closure that reevaluates the model
49             and returns the loss.
50         """
51         loss = None
52         if closure is not None:
53             loss = closure()
54
55         for group in self.param_groups:
56             for p in group['params']:
57                 if p.grad is None:
58                     continue
59                 grad = p.grad.data
60                 if grad.is_sparse:
61                     raise RuntimeError('Nadam does not support sparse gradients.')
62                 amsgrad = group['amsgrad']
63
64                 state = self.state[p]

```

```

65
66         beta1, beta2 = group['betas']
67         state['step'] += 1
68         bias_correction1 = 1 - beta1 ** state['step']
69         bias_correction2 = 1 - beta2 ** state['step']
70
71         if group['weight_decay'] != 0:
72             grad.add_(group['weight_decay'], p.data)
73
74         state['exp_avg'].mul_(beta1).add_(1 - beta1, grad)
75         state['exp_avg_sq'].mul_(beta2).addcmul_(1 - beta2, grad, grad)
76         denom = (state['exp_avg_sq'].sqrt() / math.sqrt(bias_correction2)).add_(group['
77             eps'])
78         exp_avg_hat = state['exp_avg'] / bias_correction1
79         exp_avg_hat = beta1 * exp_avg_hat + (1 - beta1) / bias_correction1 * grad
80         p.data.addcdiv_(-group['lr'], exp_avg_hat, denom)
81
82     return loss

```

### 3.8 SAG(Stochastic Average Gradient)

每次迭代只用一个样本,在第  $t$  次迭代中, SAG[?] 需要记录过去  $n$  个 ( $n \leq t$ ) 的梯度。同时在每次迭代中都随机选择一个  $i_t \in \{1, 2, \dots, n\}$ , 用当前迭代点的梯度  $\frac{\partial L(\theta, x_k)}{\partial \theta}$  替换  $n$  个记录中的第  $i_t$  个梯度。即

$$d_i^t = \begin{cases} \frac{\partial L(\theta, x_k)}{\partial \theta} & \text{if } i = i_t, \\ d_i^{t-1} & \text{otherwise.} \end{cases} \quad (30)$$

$$\theta_{t+1} \leftarrow \theta_t - \eta_t \left( \frac{1}{n} \sum_{i=1}^n d_i^t \right). \quad (31)$$

---

#### Algorithm 2 SAG

---

**Require:** 目标函数  $\min \frac{1}{n} \sum_{i=0}^n l_i(\theta)$ , 学习率  $\eta$

**Ensure:**  $\theta$

$d = 0$

**for**  $i = 1; i \leq n; i++$  **do**

$\hat{g}_i = 0$

**end for**

**for**  $t = 0; t < EPOCH; t++$  **do**

    sample  $i$  from  $\{1, 2, \dots, n\}$

$d = d - \hat{g}_i + \frac{\partial L}{\partial \theta}$

$\hat{g}_i = \frac{\partial L}{\partial \theta}$

$\theta = \theta - \frac{\eta}{n} d$

**end for**

---

SAG 比随机梯度下降快得多,是线性收敛算法。但是它要求损失函数是光滑的,而且目标函数是凸函数。而且需要存储不少历史梯度信息,用在非凸的神经网络上不太合适。

```
1 import math
2 import random
3 import torch
4 from torch.optim.optimizer import Optimizer, required
5
6 class SAG(Optimizer):
7     """Stochastic Average Gradient.
8     """
9
10    def __init__(self, params, lr=1e-2, old_grad_N = 10, weight_decay=0):
11        if not 0.0 <= lr:
12            raise ValueError("Invalid learning rate: {}".format(lr))
13        if not 0.0 <= weight_decay:
14            raise ValueError("Invalid weight_decay value: {}".format(weight_decay))
15        self.lold_grad_N = old_grad_N
16        defaults = dict(lr = lr, weight_decay = weight_decay)
17        super(SAG, self).__init__(params, defaults)
18
19        # State initialization
20        for group in self.param_groups:
21            for p in group['params']:
22                state = self.state[p]
23                state['old_grad'] = torch.zeros(torch.Size([self.lold_grad_N]) + p.data.size())
24                state['d'] = torch.zeros_like(p.data)
25
26    def step(self, closure=None):
27        """Performs a single optimization step.
28        Arguments:
29            closure (callable, optional): A closure that reevaluates the model
30            and returns the loss.
31        """
32        loss = None
33        if closure is not None:
34            loss = closure()
35
36        r = random.randrange(0, self.lold_grad_N, 1)
37
38        for group in self.param_groups:
39            for p in group['params']:
40                if p.grad is None:
41                    continue
42                grad = p.grad.data
43
44                if grad.is_sparse:
45                    raise RuntimeError('ASGD does not support sparse gradients')
46                state = self.state[p]
```

```

47
48         if group['weight_decay'] != 0:
49             grad = grad.add(group['weight_decay'], p.data)
50
51         if(torch.cuda.is_available()):
52             state['d'] = state['d'] - state['old_grad'][r].cuda() + grad
53         else:
54             state['d'] = state['d'] - state['old_grad'][r] + grad
55
56         state['old_grad'][r] = grad
57         p.data = p.data - group['lr'] / self.lold_grad_N * state['d']
58
59     return loss

```

### 3.9 SVRG(Stochastic Variance Reduction Gradient)

SVRG[?] 是对 SAG 的改进。同样需要保存前面的  $n$  次梯度, 在第  $t$  次迭代中, 随机取一个  $i_t \in \{1, 2, \dots, n\}$ , 取代  $n$  个中的第  $i_t$  个。除此之外, 每迭代  $m$  次, 更新  $\tilde{\theta}$ , 可以用最后一次的  $\theta$ , 也可以是  $m$  次中随机的一次。

$$\begin{cases} \tilde{\mu} \leftarrow \frac{1}{n} \sum_{i=0}^n g_i(\tilde{\theta}) \\ \theta_t \leftarrow \theta_{t-1} - \eta (g_{i_t}(\theta_{t-1}) - g_{i_t}(\tilde{\theta}) + \tilde{\mu}) \end{cases} \quad (32)$$

---

#### Algorithm 3 SVRG

---

**Require:** 目标函数  $\min \frac{1}{n} \sum_{i=0}^n l_i(\theta)$ , 更新频率  $m$ , 学习率  $\eta$

**Ensure:**  $\theta$

Initialize  $\hat{\theta}$

**for**  $s = 1; s < EPOCH; s++$  **do**

$\hat{\theta} = \hat{\theta}_{s-1}$

$\hat{\mu} = \frac{1}{n} \sum_{i=0}^n \frac{\partial l_i}{\partial \hat{\theta}}$

$\theta_0 = \hat{\theta}$

**for**  $t = 1; t \leq m; t++$  **do**

Randomly pick  $i_t \in \{1, 2, \dots, n\}$

$\theta_t = \theta_{t-1} - \eta \left( \frac{\partial l_{i_t}}{\partial \theta_{t-1}} - \frac{\partial l_{i_t}}{\partial \hat{\mu}} + \hat{\mu} \right)$

**end for**

OPTION I: set  $\hat{\theta}_s = \theta_m$

OPTION II: set  $\hat{\theta}_s = \theta_t$ , randomly chosen  $t \in \{1, 2, \dots, m-1\}$

**end for**

---

```

1 import math

```

```

2 import torch
3 from torch.optim.optimizer import Optimizer, required
4
5 class SVRG(Optimizer):
6     """Stochastic Variance Reduced Gradient.
7     """
8
9     def __init__(self, params, lr=1e-2, weight_decay=0, epoch = 10, batch_size = None):
10         if not 0.0 <= lr:
11             raise ValueError("Invalid learning rate: {}".format(lr))
12         if not 0.0 <= weight_decay:
13             raise ValueError("Invalid weight_decay value: {}".format(weight_decay))
14
15         self.epoch = epoch
16         self.batch_size = batch_size
17         defaults = dict(lr = lr, weight_decay = weight_decay)
18         super(SVRG, self).__init__(params, defaults)
19
20         # State initialization
21         for group in self.param_groups:
22             for p in group['params']:
23                 state = self.state[p]
24                 state['old_grad'] = torch.zeros_like(p.data)
25                 state['mu'] = torch.zeros_like(p.data)
26
27     def save_grad(self):
28         for group in self.param_groups:
29             for p in group['params']:
30                 state = self.state[p]
31                 state['old_grads'] = p.grad.data.clone()
32
33     def step(self, closure=None):
34         """Performs a single optimization step.
35         Arguments:
36             closure (callable, optional): A closure that reevaluates the model
37                 and returns the loss.
38         """
39         loss = None
40         if closure is None:
41             raise ValueError("Invalid closure function")
42
43         for group in self.param_groups:
44             for p in group['params']:
45                 state = self.state[p]
46                 state['mu'] = p.grad.data.clone()
47
48         for t in range(self.epoch):
49             closure()
50             for group in self.param_groups:

```

```

51         for p in group['params']:
52             if p.grad is None:
53                 continue
54             grad = p.grad.data
55             if grad.is_sparse:
56                 raise RuntimeError('SVRD does not support sparse gradients')
57             state = self.state[p]
58             if group['weight_decay'] != 0:
59                 grad = grad.add(group['weight_decay'], p.data)
60             p.data = p.data - group['lr'] * ( grad - state['old_grad'] + state['mu'])
61 # return loss

```

### 3.10 小结

## 4 二阶导信息

### 4.1 共轭梯度法

---

**Algorithm 4** 用于正定二次函数的共轭梯度法 (Conjugate Gradient)

---

**Require:**  $f(\vec{x}) = \frac{1}{2}\vec{x}^T Q \vec{x} + \vec{b}^T \vec{x} + \vec{c}$ , H 终止准则

**Ensure:** 极小值点  $\vec{x}^*$

set  $\vec{x}_0$ ,  $\vec{p}_0 = -\nabla f(\vec{x}_0)$ ,  $k = 0$

**for**  $k = 0; k < n; k++$  **do**

$$t_k = \frac{\nabla f(\vec{x}_k)^T \nabla f(\vec{x}_k)}{\vec{p}_k^T Q \vec{p}_k}$$

$$\vec{x}_{k+1} = \vec{x}_k + t_k \vec{p}_k$$

**if**  $\|\nabla f(\vec{x}_{k+1})\| < \epsilon$  **then**

**return**  $\vec{x}^* = \vec{x}_{k+1}$

**else**

$$a_k = \frac{\nabla f(\vec{x}_{k+1})^T Q \vec{p}_k}{\vec{p}_k^T Q \vec{p}_k}$$

$$\vec{p}_{k+1} = -\nabla f(\vec{x}_{k+1}) + a_k \vec{p}_k$$

**end if**

$k = k + 1$

**end for**

---

---

**Algorithm 5** Fletcher-Reeves 共轭梯度法

---

**Require:**  $L = \sum_{i=0}^n l_i$ , 初始  $\theta_0$

**Ensure:**  $\theta$

```
for  $t = 1; t \leq m; t++$  do
     $L_0 = L(\theta_0), \quad g_0 = g(\theta_0), \quad p_0 = -g_0, \quad k = 0$ 
    for  $k = 0; k < n; k++$  do
         $\theta_{k+1} = ls(\theta_k, p_k), \quad L_{k+1} = L(\theta_{k+1}), \quad g_{k+1} = g(\theta_{k+1})$ 
        if 满足 H 准则 then
            return  $\theta = \theta_{k+1}$ 
        else
            if  $k == n$  then
                 $\theta_0 = \theta_{k+1}$  break;
            else
                 $\alpha_k = \frac{\|g_{k+1}\|^2}{\|g_k\|^2}, \quad p_{k+1} = -g_{k+1} + \alpha_k p_k$ 
                if  $|p_{k+1}^T g_{k+1}| \leq \epsilon$  then
                     $\theta_0 = \theta_{k+1}$  break;
                else
                    if  $p_{k+1}^T g_{k+1} \geq \epsilon$  then
                         $p_{k+1} = -p_{k+1}$ 
                    end if
                end if
            end if
        end if
         $k = k + 1$ 
    end for
end for
```

---



## 4.2 牛顿法

---

### Algorithm 6 Newton 法 (Newton Method)

---

**Require:** 目标函数  $f(\vec{x})$  及其梯度  $G(\vec{x})$ , Hesse 矩阵度  $G(\vec{x})$ , H 终止准则

**Ensure:** 极小值点  $\vec{x}^*$

```

1: set  $\vec{x}_0$ ,  $k = 0$ ,  $f_0 = f(\vec{x}_0)$ ,  $G_0 = G(\vec{x}_0)$ 
2: for  $k = 0; k < n; k++$  do
3:    $G_k = G(\vec{x}_k)$ 
4:   求解  $G_k \vec{p}_k = -G_k$ 
5:    $\vec{x}_{k+1} = \vec{x}_k + \vec{p}_k$ 
6:    $f_{k+1} = f(\vec{x}_{k+1})$ ,  $G_{k+1} = G(\vec{x}_{k+1})$ 
7:   if H 终止准则 then
8:     return  $\vec{x}^* = \vec{x}_{k+1}$ 
9:   else
10:     $k = k + 1$ 
11:   end if
12: end for

```

---

## 4.3 拟牛顿法

---

### Algorithm 7 拟 Newton 法 (Quasi-Newton Method)

---

**Require:** 目标函数  $f(\vec{x})$  及其梯度  $\vec{g}(\vec{x})$ , H 终止准则

**Ensure:** 极小值点  $\vec{x}^*$

```

set  $\vec{x}_0$  和初始对称正定矩阵  $H_0$ ,  $k = 0$ ,  $f_0 = f(\vec{x}_0)$ ,  $\vec{g}_0 = \vec{g}(\vec{x}_0)$ 
for  $k = 0; k < n; k++$  do
   $\vec{x}_{k+1} = ls(\vec{x}_k, -H_k \vec{g}_k)$ 
   $f_{k+1} = f(\vec{x}_{k+1})$ ,  $\vec{g}_{k+1} = \vec{g}(\vec{x}_{k+1})$ 
   $\vec{s}_k = \vec{x}_{k+1} - \vec{x}_k$ ,  $\vec{y}_k = \vec{g}_{k+1} - \vec{g}_k$ 
  if H 终止准则 then
    return  $\vec{x}^* = \vec{x}_{k+1}$ 
  else
     $H_{k+1} = H_k + E_k$ 
     $k = k + 1$ 
  end if
end for

```

---

DFP 校正公式

$$H_{k+1} = H_k + \frac{\vec{s}_k \vec{s}_k^T}{\vec{s}_k^T \vec{y}_k} - \frac{H_k \vec{y}_k \vec{y}_k^T H_k}{\vec{y}_k^T H_k \vec{y}_k}. \quad (33)$$

BFGS 校正公式

$$H_{k+1} = H_k + \frac{1}{\vec{s}_k^T \vec{y}_k} \left[ \left( 1 + \frac{\vec{y}_k^T H_k \vec{y}_k}{\vec{s}_k^T \vec{y}_k} \right) \vec{s}_k \vec{s}_k^T - H_k \vec{y}_k \vec{s}_k^T - \vec{s}_k \vec{y}_k^T H_k \right]. \quad (34)$$

Broyden 校正公式

$$\begin{cases} H_{k+1} = H_k + \frac{\vec{s}_k \vec{s}_k^T}{\vec{y}_k^T \vec{s}_k} - \frac{H_k \vec{y}_k \vec{y}_k^T H_k}{\vec{y}_k^T H_k \vec{y}_k} + \beta \left( \vec{y}_k^T \vec{s}_k \right) \left( \vec{y}_k^T H_k \vec{y}_k \right) \vec{w}_k \vec{w}_k^T \\ \vec{w}_k = \frac{\vec{s}_k}{\vec{y}_k^T \vec{s}_k} - \frac{H_k \vec{y}_k}{\vec{y}_k^T H_k \vec{y}_k}. \end{cases} \quad (35)$$

Broyden 校正公式(??)中的  $\beta$  取值为 0 时得到 DFP 校正公式(??), 取值为  $\frac{1}{\vec{s}^T \vec{y}}$  则是 BFGS 校正公式(??)。

另外, 当  $\beta \geq 0$ , 且  $H_0$  为正定阵时, 则  $H_k$  每一个都是正定的。

#### 4.4 随机拟牛顿法

---

##### Algorithm 8 SQN Framework

---

**Require:**  $\theta_0, V, m, \eta_t$

**Ensure:**  $\theta$

```

for  $t = 0; t < EPOCH; t++$  do
   $s'_t = H_t g_t$  using the two-loop recursion.
   $s_t = -\eta_t s'_t$ 
   $\theta_{t+1} = \theta_t + s'_t$ 
  if update pairs then
    Compute  $s_t$  and  $u_t$ 
    Add a new displacement pair  $s_t, u_t$  to  $V$ 
    if  $|V| > m$  then
      Remove the eldest pair from  $V$ 
    end if
  end if
end for

```

---

---

**Algorithm 9** Two-Lopp Recursion for  $H_t g_t$ 

---

**Require:**  $\nabla f_t, u_t, s_t$ **Ensure:**  $H_{t+1} g_{t+1}$ 

```
g_t = ∇f_t
H_t^0 = (s_t^T u_t / ||u_t||^2) I
for l = t - 1; l ≥ t - p; l ++ do
    η_l = ρ_l s_l^T g_{l+1}
    g_l = g_{l+1} - η_l u_l
end for
r_{t-p-1} = H_t^0 g_{t-p}
for l = t - p; l ≥ t - 1; l -- do
    β_l = ρ_l u_l^T ρ_{l-1}
    ρ_l = ρ_{l-1} + s_l(η_t - β_l)
end for
H_{t+1} g_{t+1} = ρ
```

---

## 4.5 Hessian Free Method

---

**Algorithm 10** Hessian-Free Optimization Method

---

**Require:**  $\theta_0, \nabla f(\theta_0), \lambda$ **Ensure:**  $\theta$ 

```
for t = 0; t < EPOCH; t ++ do
    g_t = ∇f(θ_t)
    Compute λ by some methods
    B_t(v) = H(θ_t)v + λv
    Compute the step size η_t
    d_t = CG(B_t, -g_t)
    θ_{t+1} = θ_t + η_t d_t
end for
```

---

## 4.6 Sub-sampled Hessian Free Method

$S$  是整个样本的小样本集, 在第  $t$  次迭代中, 随机梯度可以估计为

$$\nabla L_{S_t}(\theta_t) = \frac{1}{|S_t|} \sum_{i \in S_t} \nabla l_i(\theta_t). \quad (36)$$

随机 Hessian 估计

$$\nabla^2 L_{S_t^H}(\theta_t) = \frac{1}{|S_t^H|} \sum_{i \in S_t^H} \nabla^2 l_i(\theta_t). \quad (37)$$

求解下降方向

$$\nabla^2 L_{S_t^H}(\theta_t) d_t = -\nabla L_{S_t}(\theta_t). \quad (38)$$

更新策略为

$$\theta \leftarrow \theta + \eta d \quad (39)$$

## 4.7 Natural Gradient

$$\begin{cases} \nabla_N L \leftarrow \left( \mathbb{E}_\theta \left[ \frac{\partial^2 L}{\partial \theta^2} \right] \right)^{-1} \frac{\partial L}{\partial \theta} \\ \theta \leftarrow \theta - \eta \nabla_N L. \end{cases} \quad (40)$$

## 5 其他类型优化方法

## 6 学习率 $\eta$ 问题

## 7 批量大小问题