

# FullyConnectedNets

2021 年 5 月 14 日

## 1 全连接神经网络

在前面的作业中，你在 CIFAR-10 上实现了一个两层的全连接神经网络。那个实现很简单，但不是很模块化，因为损失和梯度计算在一个函数内。对于一个简单的两层网络来说，还可以人为处理，但是当我们使用更大的模型时，人工处理损失和梯度就变得不切实际了。理想情况下，我们希望使用更加模块化的设计来构建网络，这样我们就可以独立地实现不同类型的层，然后将它们整合到不同架构的模型中。

在本练习中，我们将使用更模块化的方法实现全连接网络。对于每一层，我们将实现一个 `forward` 和一个 `backward` 的函数。`forward` 函数将接收输入、权重和其他参数，并返回一个输出和一个 `cache` 对象，存储反向传播所需的数据，如下所示：

```
def layer_forward(x, w):  
    """ Receive inputs x and weights w """  
    # Do some computations ...  
    z = wx# ... some intermediate value  
    # Do some more computations ...  
    out = z# the output  
  
    cache = (x, w, z, out) # Values we need to compute gradients  
  
    return out, cache
```

反向传播将接收上游的梯度和 `cache` 对象，并返回相对于输入和权重的梯度：

```
def layer_backward(dout, cache):  
    """  
    Receive dout (derivative of loss with respect to outputs) and cache,  
    and compute derivative with respect to inputs.
```

```

"""
# Unpack cache values
x, w, z, out = cache

# Use values in cache to compute derivatives
dx = dout.dot(w.T).reshape(x.shape) # Derivative of loss with respect to x
dw = x.reshape(x.shape[0], np.prod(x.shape[1:])).T.dot(dout) # Derivative of loss with respect to w

return dx, dw

```

以这种方式实现了一些层之后，我们能够轻松地将它们组合起来，以构建不同架构的分类器。

除了实现任意深度的全连接网络外，我们还将探索不同的优化更新规则，并引入 Dropout 作为正则化器和 Batch/Layer 归一化工具来更有效地优化网络。

```
[ ]: !ls
```

```
daseCV  drive  sample_data
```

```
[4]: # As usual, a bit of setup
from __future__ import print_function
import time
import numpy as np
import matplotlib.pyplot as plt
from daseCV.classifiers.fc_net import *
from daseCV.data_utils import get_CIFAR10_data
from daseCV.gradient_check import eval_numerical_gradient, \
    eval_numerical_gradient_array
from daseCV.solver import Solver

%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/
# autoreload-of-modules-in-ipython

```

```
# when changing other linked py files, this ipynb will auto-reload immediately
%load_ext autoreload
%autoreload 2

def rel_error(x, y):
    """ returns relative error """
    return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

```
[35]: # Load the (preprocessed) CIFAR10 data.
```

```
data = get_CIFAR10_data()
for k, v in list(data.items()):
    print('%s: ' % k, v.shape)
```

```
('X_train: ', (49000, 3, 32, 32))
('y_train: ', (49000,))
('X_val: ', (1000, 3, 32, 32))
('y_val: ', (1000,))
('X_test: ', (1000, 3, 32, 32))
('y_test: ', (1000,))
```

## 2 仿射层：前向传播

打开 `daseCV/layers.py` 并实现 `affine_forward` 函数。

当你完成上述函数后，你可以用下面的代码测试你的实现正确与否

```
[ ]: # Test the affine_forward function
```

```
num_inputs = 2
input_shape = (4, 5, 6)#120
output_dim = 3
 #(2*120) (120*3)+(3*1)
input_size = num_inputs * np.prod(input_shape)
weight_size = output_dim * np.prod(input_shape)

x = np.linspace(-0.1, 0.5, num=input_size).reshape(num_inputs, *input_shape)
```

```

w = np.linspace(-0.2, 0.3, num=weight_size).reshape(np.prod(input_shape),
↳output_dim)
b = np.linspace(-0.3, 0.1, num=output_dim)

out, _ = affine_forward(x, w, b)
correct_out = np.array([[ 1.49834967,  1.70660132,  1.91485297],
                        [ 3.25553199,  3.5141327,  3.77273342]])

# Compare your output with ours. The error should be around e-9 or less.
print('Testing affine_forward function:')
print('difference: ', rel_error(out, correct_out))

```

Testing affine\_forward function:

difference: 9.769849468192957e-10

```
[ ]: x.shape[1:]
```

```
[ ]: (4, 5, 6)
```

```
[ ]: list(x.shape)[1:]
```

```
[ ]: [4, 5, 6]
```

```
[ ]: np.ones([6,3,2,1]).shape[0]
```

```
[ ]: 6
```

```
[ ]: np.ones([2,3]).dot(np.ones([3,1]))
```

```
[ ]: array([[3.],
           [3.]])
```

```
[ ]: np.ones([3,2])+np.ones([2,])
```

```
[ ]: array([[2., 2.],
           [2., 2.],
           [2., 2.]])
```

### 3 仿射层：反向传播

实现 `affine_backwards` 函数，并使用数值梯度检查测试你的实现。

```
[51]: # Test the affine_backward function
np.random.seed(231)
x = np.random.randn(10, 2, 3) # N*D
w = np.random.randn(6, 5) # D*M
b = np.random.randn(5) # M*
dout = np.random.randn(10, 5) # N*M
'''
During the backward pass through the linear layer, we assume that  $\frac{\partial L}{\partial y}$ 
has already been computed. For example if the linear layer is part
of a linear classifier,
then the matrix  $y$  gives class scores; these scores are fed to a loss
function
(such as the softmax or multiclass SVM loss) which computes the scalar
 $\frac{\partial L}{\partial y}$ 
of the loss with respect to the scores.
'''

dx_num = eval_numerical_gradient_array(lambda x: affine_forward(x, w, b)[0], x,
                                       dout)
dw_num = eval_numerical_gradient_array(lambda w: affine_forward(x, w, b)[0], w,
                                       dout)
db_num = eval_numerical_gradient_array(lambda b: affine_forward(x, w, b)[0], b,
                                       dout)

_, cache = affine_forward(x, w, b)
dx, dw, db = affine_backward(dout, cache) # out: dout (x,w,b): cache

# The error should be around e-10 or less
print('Testing affine_backward function:')
print('dx error: ', rel_error(dx_num, dx))
print('dw error: ', rel_error(dw_num, dw))
```

```
print('db error: ', rel_error(db_num, db))
```

Testing affine\_backward function:

dx error: 5.399100368651805e-11

dw error: 9.904211865398145e-11

db error: 2.4122867568119087e-11

```
[48]: x.reshape(x.shape[0], *x.shape[1:]).T.shape
```

```
[48]: (3, 2, 10)
```

```
[43]: x.shape[1:]
```

```
[43]: (2, 3)
```

```
[38]: dout.dot(w.T).shape
```

```
[38]: (10, 6)
```

```
[ ]: x.shape
```

```
[39]: dout.dot(w.T).reshape(x.shape).shape
```

```
[39]: (10, 2, 3)
```

## 4 ReLU 激活函数：前向传播

在 `relu_forward` 函数中实现 ReLU 激活函数的前向传播，并使用以下代码测试您的实现：

```
[52]: # Test the relu_forward function

x = np.linspace(-0.5, 0.5, num=12).reshape(3, 4)
print(x)
out, _ = relu_forward(x)
correct_out = np.array([[ 0.,          0.,          0.,          0.,          ],
                        [ 0.,          0.,          0.04545455, 0.13636364,],
                        [ 0.22727273, 0.31818182, 0.40909091, 0.5,          ]])
print(correct_out)
```

```
# Compare your output with ours. The error should be on the order of e-8
print('Testing relu_forward function:')
print('difference: ', rel_error(out, correct_out))
```

```
[[ -0.5          -0.40909091 -0.31818182 -0.22727273]
 [ -0.13636364 -0.04545455  0.04545455  0.13636364]
 [ 0.22727273  0.31818182  0.40909091  0.5         ]]
[[0.          0.          0.          0.          ]
 [0.          0.          0.04545455 0.13636364]
 [0.22727273 0.31818182 0.40909091 0.5         ]]
Testing relu_forward function:
difference: 4.999999798022158e-08
```

## 5 ReLU 激活函数：反向传播

在 `relu_back` 函数中为 ReLU 激活函数实现反向传播，并使用数值梯度检查来测试你的实现

```
[53]: np.random.seed(231)
x = np.random.randn(10, 10)
dout = np.random.randn(*x.shape)

dx_num = eval_numerical_gradient_array(lambda x: relu_forward(x)[0], x, dout)

_, cache = relu_forward(x)
dx = relu_backward(dout, cache)

# The error should be on the order of e-12
print('Testing relu_backward function:')
print('dx error: ', rel_error(dx_num, dx))
```

```
Testing relu_backward function:
dx error: 3.2756349136310288e-12
```

### 5.1 Inline Question 1:

作业中只要求你实现 ReLU, 但是神经网络可以使用很多不同的激活函数, 每个都有它的优点和缺点。但是, 激活函数的一个常见问题是在反向传播时出现零 (或接近零) 梯度流。下列哪个激活函数会有这个问题? 如果在一维情况下考虑这些函数, 什么样的输入将会发生这种现象? 1. Sigmoid 2. ReLU 3. Leaky ReLU

### 5.2 Answer:

ReLU 会在  $x$  是负数的时候出现 0 梯度流, 而 leaky ReLU 则避免了这样死亡 RELU 的问题

## 6 “三明治”层

在神经网络中有一些常用的层模式。例如, 仿射层后面经常跟一个 ReLU 层。为了简化这些常见模式, 我们在文件 `daseCV/layer_utils.py` 中定义了几个常用的层

请查看 `affine_relu_forward` 和 `affine_relu_backward` 函数, 并且运行下列代码进行数值梯度检查:

```
[5]: from daseCV.layer_utils import affine_relu_forward, affine_relu_backward
np.random.seed(231)
x = np.random.randn(2, 3, 4)
w = np.random.randn(12, 10)
b = np.random.randn(10)
dout = np.random.randn(2, 10)

out, cache = affine_relu_forward(x, w, b)
dx, dw, db = affine_relu_backward(dout, cache)

dx_num = eval_numerical_gradient_array(lambda x: affine_relu_forward(x, w, b), x, dout)
dw_num = eval_numerical_gradient_array(lambda w: affine_relu_forward(x, w, b), w, dout)
db_num = eval_numerical_gradient_array(lambda b: affine_relu_forward(x, w, b), b, dout)
```



```
# Relative error should be around e-10 or less
print('Testing affine_relu_forward and affine_relu_backward:')
print('dx error: ', rel_error(dx_num, dx))
print('dw error: ', rel_error(dw_num, dw))
print('db error: ', rel_error(db_num, db))
```

Testing affine\_relu\_forward and affine\_relu\_backward:

dx error: 2.299579177309368e-11

dw error: 8.162011105764925e-11

db error: 7.826724021458994e-12

## 7 损失层: Softmax and SVM

在上次作业中你已经实现了这些损失函数，所以这次作业就不用做了，免费送你了。当然，你仍然应该通过查看 `daseCV/layers.py` 其中的实现来确保理解它们是如何工作的。

你可以通过运行以下程序来确保实现是正确的:

```
[29]: np.random.seed(231)
num_classes, num_inputs = 10, 50
x = 0.001 * np.random.randn(num_inputs, num_classes)
y = np.random.randint(num_classes, size=num_inputs)

dx_num = eval_numerical_gradient(lambda x: svm_loss(x, y)[0], x, verbose=False)
loss, dx = svm_loss(x, y)

# Test svm_loss function. Loss should be around 9 and dx error should be around
→the order of e-9
print('Testing svm_loss:')
print('loss: ', loss)
print('dx error: ', rel_error(dx_num, dx))

dx_num = eval_numerical_gradient(lambda x: softmax_loss(x, y)[0], x,
→verbose=False)
loss, dx = softmax_loss(x, y)
```

```
# Test softmax_loss function. Loss should be close to 2.3 and dx error should
↳ be around e-8
print('\nTesting softmax_loss:')
print('loss: ', loss)
print('dx error: ', rel_error(dx_num, dx))
```

Testing svm\_loss:

loss: 8.999602749096233

dx error: 1.4021566006651672e-09

Testing softmax\_loss:

loss: 2.302545844500738

dx error: 9.384673161989355e-09

```
[6]: np.arange(3)
```

```
[6]: array([0, 1, 2])
```

```
[21]: np.zeros_like(np.ones((2,3)))
```

```
[21]: array([[0., 0., 0.],
            [0., 0., 0.]])
```

```
[14]: np.ones((2,3))[0,1,2].shape
```

```
[14]: (2,)
```

```
[17]: a=np.array([1,2,3,4,5])
      aa=a[:,np.newaxis]
      print(aa.shape)
      print (aa)
      # 增加一个维度
```

```
(5, 1)
```

```
[[1]
```

```
[2]
```

```
[3]
```

```
[4]
```

```
[5]]
```

## 8 两层网络

在之前的作业中，你已经实现了一个简单的两层神经网络。现在你已经模块化地实现了一些层，你将使用这些模块重新实现两层网络。

打开文件 `daseCV/classifiers/fc_net`。并完成 `TwoLayerNet` 类的实现。这个类将作为这个作业中其他网络的模块，所以请通读它以确保你理解了这个 API。你可以运行下面的单元来测试您的实现。

```
[32]: np.random.seed(231)
N, D, H, C = 3, 5, 50, 7
X = np.random.randn(N, D)
y = np.random.randint(C, size=N)

std = 1e-3
model = TwoLayerNet(input_dim=D, hidden_dim=H, num_classes=C, weight_scale=std)

print('Testing initialization ... ')
W1_std = abs(model.params['W1'].std() - std)
b1 = model.params['b1']
W2_std = abs(model.params['W2'].std() - std)
b2 = model.params['b2']
assert W1_std < std / 10, 'First layer weights do not seem right'
assert np.all(b1 == 0), 'First layer biases do not seem right'
assert W2_std < std / 10, 'Second layer weights do not seem right'
assert np.all(b2 == 0), 'Second layer biases do not seem right'

print('Testing test-time forward pass ... ')
model.params['W1'] = np.linspace(-0.7, 0.3, num=D*H).reshape(D, H)
model.params['b1'] = np.linspace(-0.1, 0.9, num=H)
model.params['W2'] = np.linspace(-0.3, 0.4, num=H*C).reshape(H, C)
model.params['b2'] = np.linspace(-0.9, 0.1, num=C)
X = np.linspace(-5.5, 4.5, num=N*D).reshape(D, N).T
scores = model.loss(X)
correct_scores = np.asarray(
    [[11.53165108, 12.2917344, 13.05181771, 13.81190102, 14.57198434, 15.
↪ 33206765, 16.09215096],
```

```

    [12.05769098, 12.74614105, 13.43459113, 14.1230412, 14.81149128, 15.
↪49994135, 16.18839143],
    [12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.05099822, 15.
↪66781506, 16.2846319 ]])
scores_diff = np.abs(scores - correct_scores).sum()
assert scores_diff < 1e-6, 'Problem with test-time forward pass'

print('Testing training loss (no regularization)')
y = np.asarray([0, 5, 1])
loss, grads = model.loss(X, y)
correct_loss = 3.4702243556
assert abs(loss - correct_loss) < 1e-10, 'Problem with training-time loss'

model.reg = 1.0
loss, grads = model.loss(X, y)
correct_loss = 26.5948426952
assert abs(loss - correct_loss) < 1e-10, 'Problem with regularization loss'

# Errors should be around e-7 or less
for reg in [0.0, 0.7]:
    print('Running numeric gradient check with reg = ', reg)
    model.reg = reg
    loss, grads = model.loss(X, y)

    for name in sorted(grads):
        f = lambda _: model.loss(X, y)[0]
        grad_num = eval_numerical_gradient(f, model.params[name], verbose=False)
        print('%s relative error: %.2e' % (name, rel_error(grad_num, grads[name])))

```

Testing initialization ...

Testing test-time forward pass ...

Testing training loss (no regularization)

Running numeric gradient check with reg = 0.0

W1 relative error: 1.83e-08

W2 relative error: 3.12e-10

b1 relative error: 9.83e-09

```

b2 relative error: 4.33e-10
Running numeric gradient check with reg = 0.7
W1 relative error: 2.53e-07
W2 relative error: 2.85e-08
b1 relative error: 1.56e-08
b2 relative error: 7.76e-10

```

## 9 Solver

在之前的作业中，模型的训练逻辑与模型本身是耦合的。在这次作业中，按照更加模块化的设计，我们将模型的训练逻辑划分为单独的类。

打开文件 `daseCV/solver`，通读一遍以熟悉 API。然后使用一个 `Sovler` 实例来训练一个 `TwoLayerNet`，它可以在验证集上达到至少 50% 的精度。

```

[36]: model = TwoLayerNet()
      #solver = None

#####
# TODO: Use a Solver instance to train a TwoLayerNet that achieves at least #
# 50% accuracy on the validation set.                                     #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
solver = Solver(model, data,
                update_rule='sgd',
                lr_decay=0.98,
                batch_size=128,
                print_every=100,
                optim_config={'learning_rate': 1e-3},
                num_epochs = 9)
solver.train()

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
#####
#                               END OF YOUR CODE                               #
#####

```

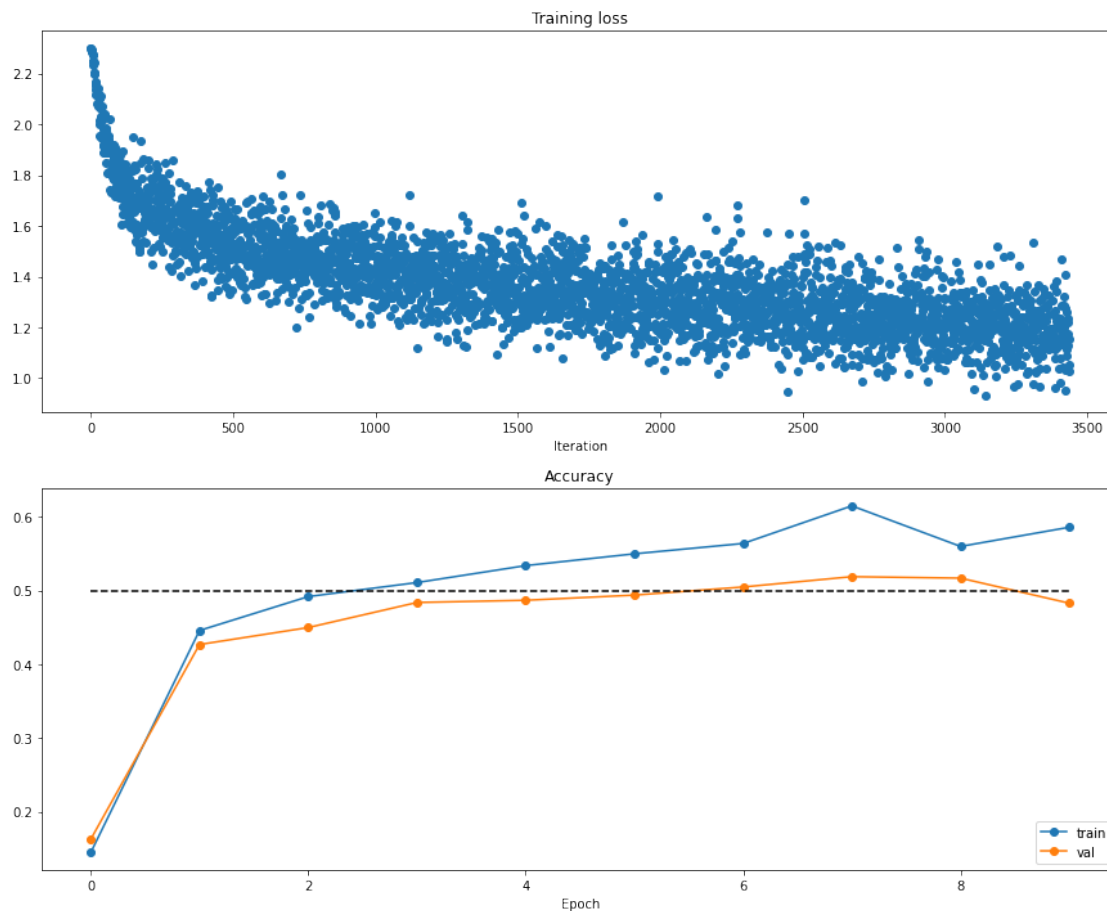
```
(Iteration 1 / 3438) loss: 2.299875
(Epoch 0 / 9) train acc: 0.145000; val_acc: 0.163000
(Iteration 101 / 3438) loss: 1.746332
(Iteration 201 / 3438) loss: 1.733316
(Iteration 301 / 3438) loss: 1.580383
(Epoch 1 / 9) train acc: 0.446000; val_acc: 0.427000
(Iteration 401 / 3438) loss: 1.499206
(Iteration 501 / 3438) loss: 1.479629
(Iteration 601 / 3438) loss: 1.460441
(Iteration 701 / 3438) loss: 1.359739
(Epoch 2 / 9) train acc: 0.492000; val_acc: 0.450000
(Iteration 801 / 3438) loss: 1.611551
(Iteration 901 / 3438) loss: 1.440300
(Iteration 1001 / 3438) loss: 1.324513
(Iteration 1101 / 3438) loss: 1.522455
(Epoch 3 / 9) train acc: 0.511000; val_acc: 0.484000
(Iteration 1201 / 3438) loss: 1.538411
(Iteration 1301 / 3438) loss: 1.293379
(Iteration 1401 / 3438) loss: 1.539784
(Iteration 1501 / 3438) loss: 1.230064
(Epoch 4 / 9) train acc: 0.534000; val_acc: 0.487000
(Iteration 1601 / 3438) loss: 1.599137
(Iteration 1701 / 3438) loss: 1.325792
(Iteration 1801 / 3438) loss: 1.435145
(Iteration 1901 / 3438) loss: 1.390626
(Epoch 5 / 9) train acc: 0.550000; val_acc: 0.494000
(Iteration 2001 / 3438) loss: 1.204316
(Iteration 2101 / 3438) loss: 1.253460
(Iteration 2201 / 3438) loss: 1.327926
(Epoch 6 / 9) train acc: 0.564000; val_acc: 0.505000
(Iteration 2301 / 3438) loss: 1.097801
(Iteration 2401 / 3438) loss: 1.226286
(Iteration 2501 / 3438) loss: 1.119927
(Iteration 2601 / 3438) loss: 1.269149
(Epoch 7 / 9) train acc: 0.615000; val_acc: 0.519000
(Iteration 2701 / 3438) loss: 1.355638
(Iteration 2801 / 3438) loss: 1.301551
```

```
(Iteration 2901 / 3438) loss: 1.326050
(Iteration 3001 / 3438) loss: 1.310682
(Epoch 8 / 9) train acc: 0.560000; val_acc: 0.517000
(Iteration 3101 / 3438) loss: 1.181888
(Iteration 3201 / 3438) loss: 1.113038
(Iteration 3301 / 3438) loss: 1.277730
(Iteration 3401 / 3438) loss: 1.230148
(Epoch 9 / 9) train acc: 0.586000; val_acc: 0.483000
```

[37]: *# Run this cell to visualize training loss and train / val accuracy*

```
plt.subplot(2, 1, 1)
plt.title('Training loss')
plt.plot(solver.loss_history, 'o')
plt.xlabel('Iteration')

plt.subplot(2, 1, 2)
plt.title('Accuracy')
plt.plot(solver.train_acc_history, '-o', label='train')
plt.plot(solver.val_acc_history, '-o', label='val')
plt.plot([0.5] * len(solver.val_acc_history), 'k--')
plt.xlabel('Epoch')
plt.legend(loc='lower right')
plt.gcf().set_size_inches(15, 12)
plt.show()
```



## 10 多层网络

接下来，请实现一个带有任意数量的隐层的全连接网络。

阅读 `daseCV/classifiers/fc_net.py` 中的 `FullyConnectedNet` 类。

实现初始化、前向传播和反向传播的函数，暂时不要考虑实现 dropout 或 batch/layer normalization，我们将在后面添加上去。

### 10.1 初始化 loss 和梯度检查

刚开始要做完整性检查，运行以下代码来检查初始 loss，并对有正则化和无正则化的网络进行梯度检查。请问初始的 loss 合理吗？



在梯度检查中，你应该期望得到  $1e-7$  或更少的 errors。

```
[41]: np.random.seed(231)
N, D, H1, H2, C = 2, 15, 20, 30, 10
X = np.random.randn(N, D)
y = np.random.randint(C, size=(N,))

for reg in [0, 3.14]:
    print('Running check with reg = ', reg)
    model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
                              reg=reg, weight_scale=5e-2, dtype=np.float64)

    loss, grads = model.loss(X, y)
    print('Initial loss: ', loss)

    # Most of the errors should be on the order of e-7 or smaller.
    # NOTE: It is fine however to see an error for W2 on the order of e-5
    # for the check when reg = 0.0
    for name in sorted(grads):
        f = lambda _: model.loss(X, y)[0]
        grad_num = eval_numerical_gradient(f, model.params[name], verbose=False,
        ↪h=1e-5)
        print('%s relative error: %.2e' % (name, rel_error(grad_num, grads[name])))
```

```
Running check with reg = 0
Initial loss: 2.3004790897684924
W1 relative error: 1.48e-07
W2 relative error: 2.21e-05
W3 relative error: 3.53e-07
b1 relative error: 5.38e-09
b2 relative error: 2.09e-09
b3 relative error: 5.80e-11
Running check with reg = 3.14
Initial loss: 7.052114776533016
W1 relative error: 7.36e-09
W2 relative error: 6.87e-08
W3 relative error: 3.48e-08
```

```

b1 relative error: 1.48e-08
b2 relative error: 1.72e-09
b3 relative error: 1.80e-10

```

实现另一个完整性检查，请确保你可以过拟合 50 个图像的小数据集。首先，我们将尝试一个三层网络，每个隐藏层有 100 个单元。在接下来的代码中，调整 **learning rate** 和 **weight initialization scale** 以达到过拟合，在 20 epoch 内达到 100% 的训练精度。

[42]: *# TODO: Use a three-layer Net to overfit 50 training examples by  
# tweaking just the learning rate and initialization scale.*

```

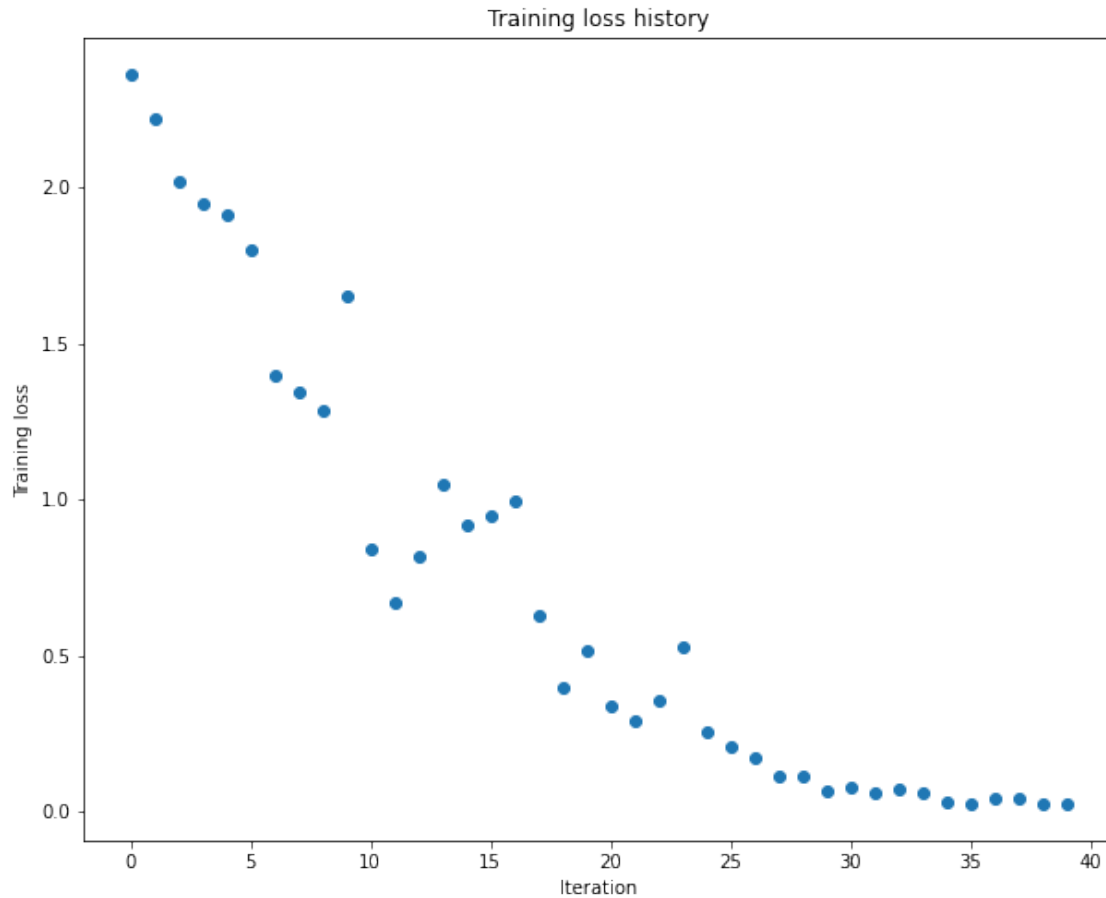
num_train = 50
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

weight_scale = 1e-2    # Experiment with this!
learning_rate = 1e-2   # Experiment with this!
model = FullyConnectedNet([100, 100],
                           weight_scale=weight_scale, dtype=np.float64)
solver = Solver(model, small_data,
                 print_every=10, num_epochs=20, batch_size=25,
                 update_rule='sgd',
                 optim_config={
                     'learning_rate': learning_rate,
                 })
solver.train()

plt.plot(solver.loss_history, 'o')
plt.title('Training loss history')
plt.xlabel('Iteration')
plt.ylabel('Training loss')
plt.show()

```

```
(Iteration 1 / 40) loss: 2.363364
(Epoch 0 / 20) train acc: 0.180000; val_acc: 0.108000
(Epoch 1 / 20) train acc: 0.320000; val_acc: 0.127000
(Epoch 2 / 20) train acc: 0.440000; val_acc: 0.172000
(Epoch 3 / 20) train acc: 0.500000; val_acc: 0.184000
(Epoch 4 / 20) train acc: 0.540000; val_acc: 0.181000
(Epoch 5 / 20) train acc: 0.740000; val_acc: 0.190000
(Iteration 11 / 40) loss: 0.839976
(Epoch 6 / 20) train acc: 0.740000; val_acc: 0.187000
(Epoch 7 / 20) train acc: 0.740000; val_acc: 0.183000
(Epoch 8 / 20) train acc: 0.820000; val_acc: 0.177000
(Epoch 9 / 20) train acc: 0.860000; val_acc: 0.200000
(Epoch 10 / 20) train acc: 0.920000; val_acc: 0.191000
(Iteration 21 / 40) loss: 0.337174
(Epoch 11 / 20) train acc: 0.960000; val_acc: 0.189000
(Epoch 12 / 20) train acc: 0.940000; val_acc: 0.180000
(Epoch 13 / 20) train acc: 1.000000; val_acc: 0.199000
(Epoch 14 / 20) train acc: 1.000000; val_acc: 0.199000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.195000
(Iteration 31 / 40) loss: 0.075911
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.182000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.201000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.207000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.185000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.192000
```



现在尝试使用一个五层的网络，每层 100 个单元，对 50 张图片进行训练。同样，你将调整 learning rate 和 weight initialization scale 比例，你应该能够在 20 个 epoch 内实现 100% 的训练精度。

[43]: *# TODO: Use a five-layer Net to overfit 50 training examples by  
# tweaking just the learning rate and initialization scale.*

```
num_train = 50
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}
```

```

weight_scale = 1e-1 # Experiment with this!
learning_rate = 2e-3 # Experiment with this!
model = FullyConnectedNet([100, 100, 100, 100],
                           weight_scale=weight_scale, dtype=np.float64)
solver = Solver(model, small_data,
                 print_every=10, num_epochs=20, batch_size=25,
                 update_rule='sgd',
                 optim_config={
                     'learning_rate': learning_rate,
                 })
solver.train()

plt.plot(solver.loss_history, 'o')
plt.title('Training loss history')
plt.xlabel('Iteration')
plt.ylabel('Training loss')
plt.show()

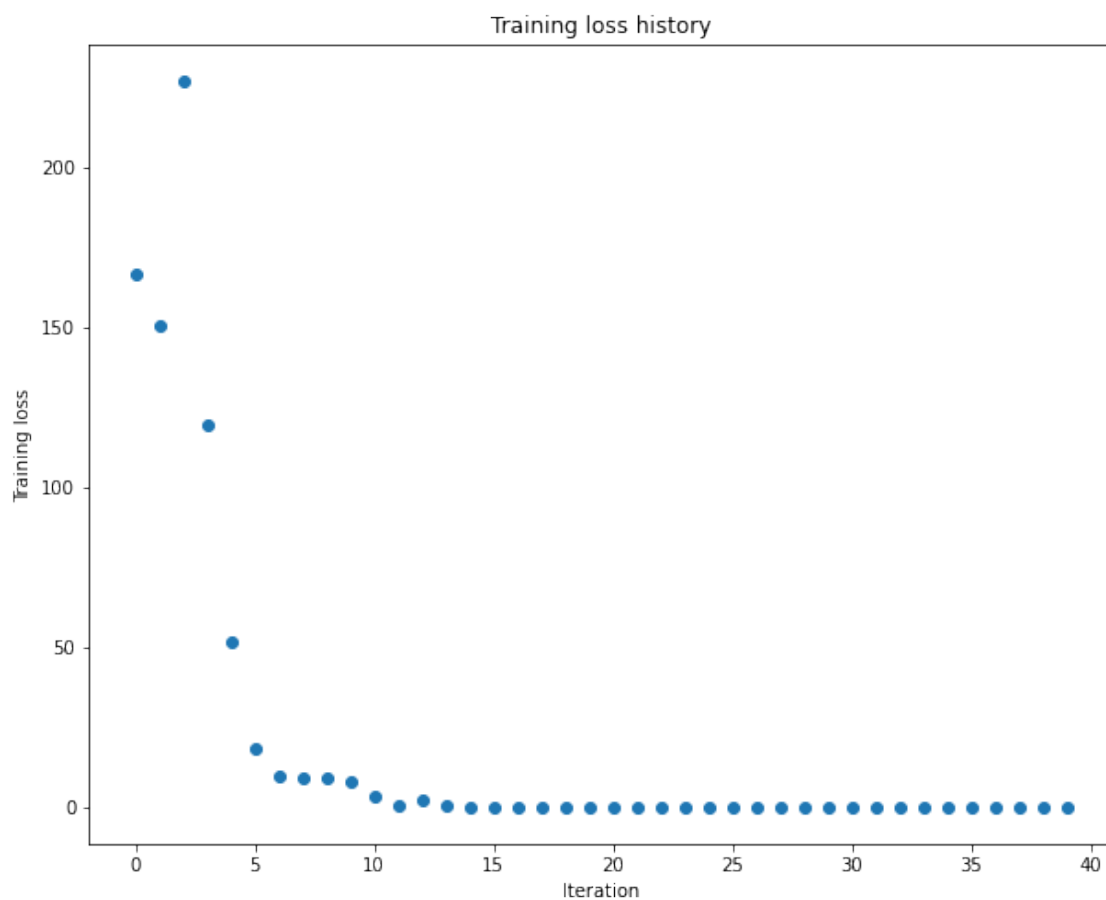
```

```

(Iteration 1 / 40) loss: 166.501707
(Epoch 0 / 20) train acc: 0.100000; val_acc: 0.107000
(Epoch 1 / 20) train acc: 0.320000; val_acc: 0.101000
(Epoch 2 / 20) train acc: 0.160000; val_acc: 0.122000
(Epoch 3 / 20) train acc: 0.380000; val_acc: 0.106000
(Epoch 4 / 20) train acc: 0.520000; val_acc: 0.111000
(Epoch 5 / 20) train acc: 0.760000; val_acc: 0.113000
(Iteration 11 / 40) loss: 3.343141
(Epoch 6 / 20) train acc: 0.840000; val_acc: 0.122000
(Epoch 7 / 20) train acc: 0.920000; val_acc: 0.113000
(Epoch 8 / 20) train acc: 0.940000; val_acc: 0.125000
(Epoch 9 / 20) train acc: 0.960000; val_acc: 0.125000
(Epoch 10 / 20) train acc: 0.980000; val_acc: 0.121000
(Iteration 21 / 40) loss: 0.039138
(Epoch 11 / 20) train acc: 0.980000; val_acc: 0.123000
(Epoch 12 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 13 / 20) train acc: 1.000000; val_acc: 0.121000

```

```
(Epoch 14 / 20) train acc: 1.000000; val_acc: 0.121000  
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.121000  
(Iteration 31 / 40) loss: 0.000644  
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.121000  
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.121000  
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.121000  
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.121000  
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.121000
```



**Inline Question 2:** 你注意到训练三层网和训练五层网难度的区别了吗? 根据你的经验, 哪个网络对 initialization scale 更敏感? 为什么会这样呢?

## 10.2 Answer:

五层网，初始网络规模更大时，为了学到更多特征，学习会变得更复杂困难，所以网络越深，对更大的网络规模会更敏感。

## 11 更新规则

到目前为止，我们使用了普通的随机梯度下降法 (SGD) 作为我们的更新规则。更复杂的更新规则可以更容易地训练深度网络。我们将实现一些最常用的更新规则，并将它们与普通的 SGD 进行比较。

## 12 SGD+Momentum

带动量的随机梯度下降法是一种广泛使用的更新规则，它使深度网络的收敛速度快于普通的随机梯度下降法。更多信息参见 <http://cs231n.github.io/neural-networks-3/#sgd> 动量更新部分。

打开文件 `daseCV/optim`，并阅读该文件顶部的文档，以确保你理解了该 API。在函数 `sgd_momentum` 中实现 SGD+ 动量更新规则，并运行以下代码检查你的实现。你会看到 `errors` 小于  $e-8$ 。

```
[45]: from daseCV.optim import sgd_momentum

N, D = 4, 5
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
v = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-3, 'velocity': v}
next_w, _ = sgd_momentum(w, dw, config=config)

expected_next_w = np.asarray([
    [ 0.1406,      0.20738947,  0.27417895,  0.34096842,  0.40775789],
    [ 0.47454737,  0.54133684,  0.60812632,  0.67491579,  0.74170526],
    [ 0.80849474,  0.87528421,  0.94207368,  1.00886316,  1.07565263],
    [ 1.14244211,  1.20923158,  1.27602105,  1.34281053,  1.4096    ]])
expected_velocity = np.asarray([
```

```

[ 0.5406,      0.55475789,  0.56891579, 0.58307368,  0.59723158],
[ 0.61138947,  0.62554737,  0.63970526,  0.65386316,  0.66802105],
[ 0.68217895,  0.69633684,  0.71049474,  0.72465263,  0.73881053],
[ 0.75296842,  0.76712632,  0.78128421,  0.79544211,  0.8096    ]]

# Should see relative errors around e-8 or less
print('next_w error: ', rel_error(next_w, expected_next_w))
print('velocity error: ', rel_error(expected_velocity, config['velocity']))

```

```
next_w error:  8.882347033505819e-09
```

```
velocity error:  4.269287743278663e-09
```

当你完成了上面的步骤，运行以下代码来训练一个具有 SGD 和 SGD+momentum 的六层网络。你应该看到 SGD+momentum 更新规则收敛得更快。

```

[46]: num_train = 4000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

solvers = {}

for update_rule in ['sgd', 'sgd_momentum']:
    print('running with ', update_rule)
    model = FullyConnectedNet([100, 100, 100, 100, 100], weight_scale=5e-2)

    solver = Solver(model, small_data,
                    num_epochs=5, batch_size=100,
                    update_rule=update_rule,
                    optim_config={
                        'learning_rate': 5e-3,
                    },
                    verbose=True)
    solvers[update_rule] = solver

```



```

    solver.train()
    print()

plt.subplot(3, 1, 1)
plt.title('Training loss')
plt.xlabel('Iteration')

plt.subplot(3, 1, 2)
plt.title('Training accuracy')
plt.xlabel('Epoch')

plt.subplot(3, 1, 3)
plt.title('Validation accuracy')
plt.xlabel('Epoch')

for update_rule, solver in solvers.items():
    plt.subplot(3, 1, 1)
    plt.plot(solver.loss_history, 'o', label="loss_%s" % update_rule)

    plt.subplot(3, 1, 2)
    plt.plot(solver.train_acc_history, '-o', label="train_acc_%s" % update_rule)

    plt.subplot(3, 1, 3)
    plt.plot(solver.val_acc_history, '-o', label="val_acc_%s" % update_rule)

for i in [1, 2, 3]:
    plt.subplot(3, 1, i)
    plt.legend(loc='upper center', ncol=4)
plt.gcf().set_size_inches(15, 15)
plt.show()

```

running with `sgd`

(Iteration 1 / 200) loss: 2.559978

(Epoch 0 / 5) train acc: 0.104000; val\_acc: 0.107000

(Iteration 11 / 200) loss: 2.356069

(Iteration 21 / 200) loss: 2.214091

```
(Iteration 31 / 200) loss: 2.205928
(Epoch 1 / 5) train acc: 0.225000; val_acc: 0.193000
(Iteration 41 / 200) loss: 2.132095
(Iteration 51 / 200) loss: 2.118950
(Iteration 61 / 200) loss: 2.116443
(Iteration 71 / 200) loss: 2.132549
(Epoch 2 / 5) train acc: 0.298000; val_acc: 0.260000
(Iteration 81 / 200) loss: 1.977227
(Iteration 91 / 200) loss: 2.007528
(Iteration 101 / 200) loss: 2.004762
(Iteration 111 / 200) loss: 1.885342
(Epoch 3 / 5) train acc: 0.343000; val_acc: 0.287000
(Iteration 121 / 200) loss: 1.891517
(Iteration 131 / 200) loss: 1.923677
(Iteration 141 / 200) loss: 1.957743
(Iteration 151 / 200) loss: 1.966736
(Epoch 4 / 5) train acc: 0.322000; val_acc: 0.305000
(Iteration 161 / 200) loss: 1.801483
(Iteration 171 / 200) loss: 1.973780
(Iteration 181 / 200) loss: 1.666572
(Iteration 191 / 200) loss: 1.909494
(Epoch 5 / 5) train acc: 0.372000; val_acc: 0.319000
```

running with `sgd_momentum`

```
(Iteration 1 / 200) loss: 3.153778
(Epoch 0 / 5) train acc: 0.099000; val_acc: 0.088000
(Iteration 11 / 200) loss: 2.227203
(Iteration 21 / 200) loss: 2.125706
(Iteration 31 / 200) loss: 1.932695
(Epoch 1 / 5) train acc: 0.307000; val_acc: 0.260000
(Iteration 41 / 200) loss: 1.946488
(Iteration 51 / 200) loss: 1.778584
(Iteration 61 / 200) loss: 1.758119
(Iteration 71 / 200) loss: 1.849137
(Epoch 2 / 5) train acc: 0.382000; val_acc: 0.322000
(Iteration 81 / 200) loss: 2.048671
(Iteration 91 / 200) loss: 1.693223
```

```
(Iteration 101 / 200) loss: 1.511693
(Iteration 111 / 200) loss: 1.390754
(Epoch 3 / 5) train acc: 0.458000; val_acc: 0.338000
(Iteration 121 / 200) loss: 1.670614
(Iteration 131 / 200) loss: 1.540271
(Iteration 141 / 200) loss: 1.597365
(Iteration 151 / 200) loss: 1.609851
(Epoch 4 / 5) train acc: 0.490000; val_acc: 0.327000
(Iteration 161 / 200) loss: 1.472687
(Iteration 171 / 200) loss: 1.378620
(Iteration 181 / 200) loss: 1.378175
(Iteration 191 / 200) loss: 1.306439
(Epoch 5 / 5) train acc: 0.529000; val_acc: 0.369000
```

```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:39:
```

```
MatplotlibDeprecationWarning: Adding an axes using the same arguments as a
previous axes currently reuses the earlier instance. In a future version, a new
instance will always be created and returned. Meanwhile, this warning can be
suppressed, and the future behavior ensured, by passing a unique label to each
axes instance.
```

```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:42:
```

```
MatplotlibDeprecationWarning: Adding an axes using the same arguments as a
previous axes currently reuses the earlier instance. In a future version, a new
instance will always be created and returned. Meanwhile, this warning can be
suppressed, and the future behavior ensured, by passing a unique label to each
axes instance.
```

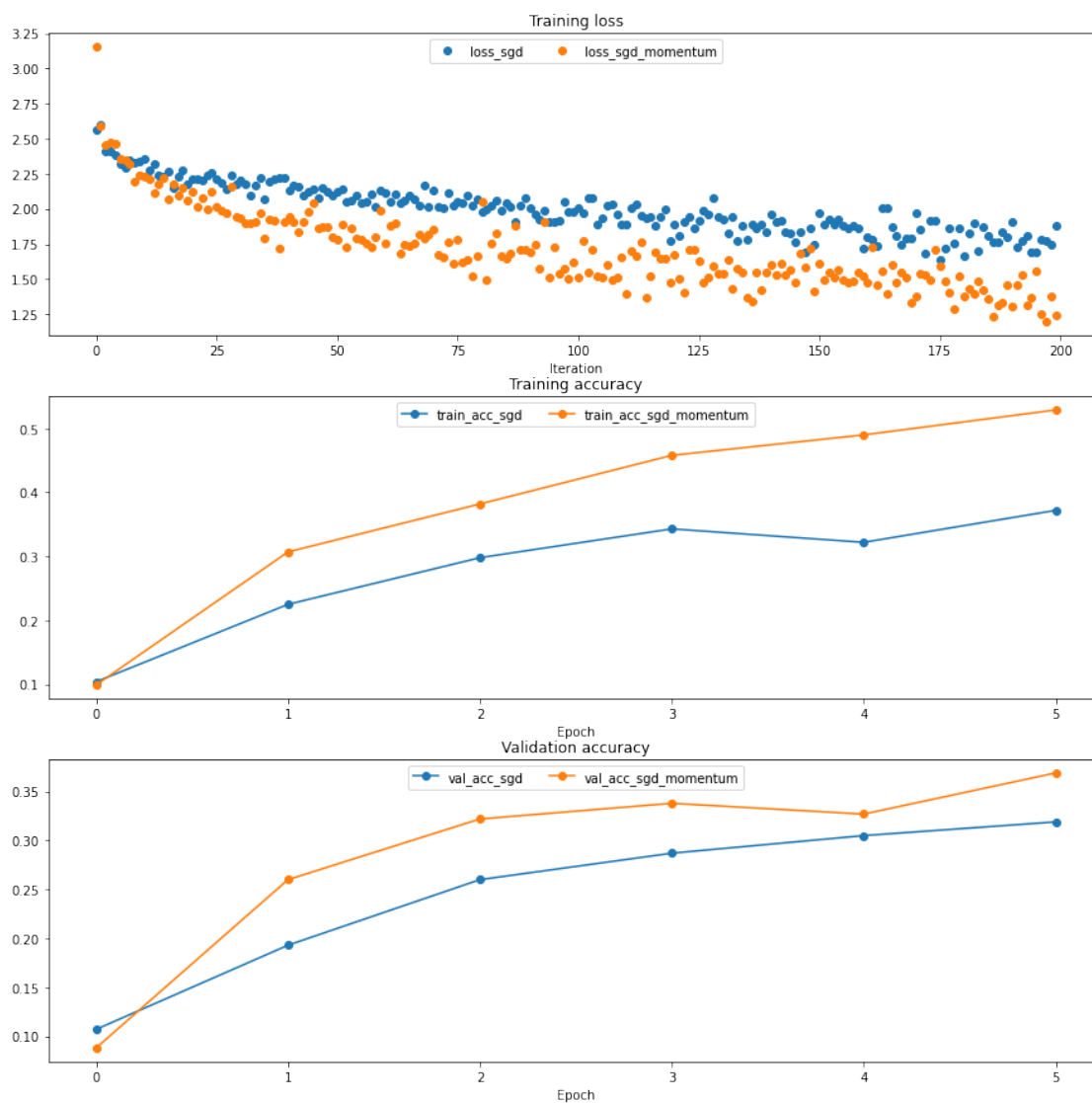
```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:45:
```

```
MatplotlibDeprecationWarning: Adding an axes using the same arguments as a
previous axes currently reuses the earlier instance. In a future version, a new
instance will always be created and returned. Meanwhile, this warning can be
suppressed, and the future behavior ensured, by passing a unique label to each
axes instance.
```

```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:49:
```

```
MatplotlibDeprecationWarning: Adding an axes using the same arguments as a
previous axes currently reuses the earlier instance. In a future version, a new
instance will always be created and returned. Meanwhile, this warning can be
```

suppressed, and the future behavior ensured, by passing a unique label to each axes instance.



## 13 RMSProp and Adam

RMSProp [1] 和 Adam [2] 是另外两个更新规则，它们通过使用梯度的二阶矩平均值来设置每个参数的学习速率。

在文件 `daseCV/optim` 中实现 RMSProp 函数和 Adam 函数，并使用下面的代码来检查您的实现。

[1] Tijmen Tieleman and Geoffrey Hinton. “Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude.” COURSERA: Neural Networks for Machine Learning 4 (2012).

[2] Diederik Kingma and Jimmy Ba, “Adam: A Method for Stochastic Optimization”, ICLR 2015.

```
[48]: # Test RMSProp implementation
from daseCV.optim import rmsprop

N, D = 4, 5
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
cache = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-2, 'cache': cache}
next_w, _ = rmsprop(w, dw, config=config)

expected_next_w = np.asarray([
    [-0.39223849, -0.34037513, -0.28849239, -0.23659121, -0.18467247],
    [-0.132737,   -0.08078555, -0.02881884,  0.02316247,  0.07515774],
    [ 0.12716641,  0.17918792,  0.23122175,  0.28326742,  0.33532447],
    [ 0.38739248,  0.43947102,  0.49155973,  0.54365823,  0.59576619]])
expected_cache = np.asarray([
    [ 0.5976,      0.6126277,  0.6277108,  0.64284931,  0.65804321],
    [ 0.67329252,  0.68859723,  0.70395734,  0.71937285,  0.73484377],
    [ 0.75037008,  0.7659518,   0.78158892,  0.79728144,  0.81302936],
    [ 0.82883269,  0.84469141,  0.86060554,  0.87657507,  0.8926   ]])

# You should see relative errors around e-7 or less
print('next_w error: ', rel_error(expected_next_w, next_w))
print('cache error: ', rel_error(expected_cache, config['cache']))
```

```
next_w error:  9.524687511038133e-08
```

```
cache error:  2.6477955807156126e-09
```

```
[50]: # Test Adam implementation
from daseCV.optim import adam
```

```
N, D = 4, 5
```

```

w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
m = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)
v = np.linspace(0.7, 0.5, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-2, 'm': m, 'v': v, 't': 5}
next_w, _ = adam(w, dw, config=config)

expected_next_w = np.asarray([
    [-0.40094747, -0.34836187, -0.29577703, -0.24319299, -0.19060977],
    [-0.1380274, -0.08544591, -0.03286534, 0.01971428, 0.0722929],
    [ 0.1248705, 0.17744702, 0.23002243, 0.28259667, 0.33516969],
    [ 0.38774145, 0.44031188, 0.49288093, 0.54544852, 0.59801459]])
expected_v = np.asarray([
    [ 0.69966, 0.68908382, 0.67851319, 0.66794809, 0.65738853,],
    [ 0.64683452, 0.63628604, 0.6257431, 0.61520571, 0.60467385,],
    [ 0.59414753, 0.58362676, 0.57311152, 0.56260183, 0.55209767,],
    [ 0.54159906, 0.53110598, 0.52061845, 0.51013645, 0.49966, ]])
expected_m = np.asarray([
    [ 0.48, 0.49947368, 0.51894737, 0.53842105, 0.55789474],
    [ 0.57736842, 0.59684211, 0.61631579, 0.63578947, 0.65526316],
    [ 0.67473684, 0.69421053, 0.71368421, 0.73315789, 0.75263158],
    [ 0.77210526, 0.79157895, 0.81105263, 0.83052632, 0.85 ]])

# You should see relative errors around e-7 or less
print('next_w error: ', rel_error(expected_next_w, next_w))
print('v error: ', rel_error(expected_v, config['v']))
print('m error: ', rel_error(expected_m, config['m']))

```

```
next_w error: 0.20720703668629928
```

```
v error: 4.208314038113071e-09
```

```
m error: 4.214963193114416e-09
```

当你完成了上面 RMSProp 和 Adam 函数后，运行下面的代码训练一对网络，其中分别使用了上述两个方法

```

[51]: learning_rates = {'rmsprop': 1e-4, 'adam': 1e-3}
for update_rule in ['adam', 'rmsprop']:
    print('running with ', update_rule)
    model = FullyConnectedNet([100, 100, 100, 100, 100], weight_scale=5e-2)

    solver = Solver(model, small_data,
                    num_epochs=5, batch_size=100,
                    update_rule=update_rule,
                    optim_config={
                        'learning_rate': learning_rates[update_rule]
                    },
                    verbose=True)
    solvers[update_rule] = solver
    solver.train()
    print()

plt.subplot(3, 1, 1)
plt.title('Training loss')
plt.xlabel('Iteration')

plt.subplot(3, 1, 2)
plt.title('Training accuracy')
plt.xlabel('Epoch')

plt.subplot(3, 1, 3)
plt.title('Validation accuracy')
plt.xlabel('Epoch')

for update_rule, solver in list(solvers.items()):
    plt.subplot(3, 1, 1)
    plt.plot(solver.loss_history, 'o', label=update_rule)

    plt.subplot(3, 1, 2)
    plt.plot(solver.train_acc_history, '-o', label=update_rule)

    plt.subplot(3, 1, 3)

```

```
plt.plot(solver.val_acc_history, '-o', label=update_rule)

for i in [1, 2, 3]:
    plt.subplot(3, 1, i)
    plt.legend(loc='upper center', ncol=4)
plt.gcf().set_size_inches(15, 15)
plt.show()
```

running with adam

```
(Iteration 1 / 200) loss: 3.476928
(Epoch 0 / 5) train acc: 0.105000; val_acc: 0.100000
(Iteration 11 / 200) loss: 2.036672
(Iteration 21 / 200) loss: 2.231720
(Iteration 31 / 200) loss: 1.986361
(Epoch 1 / 5) train acc: 0.290000; val_acc: 0.234000
(Iteration 41 / 200) loss: 1.965305
(Iteration 51 / 200) loss: 1.855719
(Iteration 61 / 200) loss: 1.983865
(Iteration 71 / 200) loss: 1.688213
(Epoch 2 / 5) train acc: 0.319000; val_acc: 0.295000
(Iteration 81 / 200) loss: 1.735082
(Iteration 91 / 200) loss: 1.653525
(Iteration 101 / 200) loss: 1.709237
(Iteration 111 / 200) loss: 1.905604
(Epoch 3 / 5) train acc: 0.335000; val_acc: 0.307000
(Iteration 121 / 200) loss: 1.661850
(Iteration 131 / 200) loss: 1.769786
(Iteration 141 / 200) loss: 1.668342
(Iteration 151 / 200) loss: 1.719504
(Epoch 4 / 5) train acc: 0.426000; val_acc: 0.358000
(Iteration 161 / 200) loss: 1.691765
(Iteration 171 / 200) loss: 1.416469
(Iteration 181 / 200) loss: 1.620160
(Iteration 191 / 200) loss: 1.578580
(Epoch 5 / 5) train acc: 0.405000; val_acc: 0.347000
```



```
running with  rmsprop
(Iteration 1 / 200) loss: 2.589166
(Epoch 0 / 5) train acc: 0.119000; val_acc: 0.146000
(Iteration 11 / 200) loss: 2.032921
(Iteration 21 / 200) loss: 1.897278
(Iteration 31 / 200) loss: 1.770793
(Epoch 1 / 5) train acc: 0.381000; val_acc: 0.320000
(Iteration 41 / 200) loss: 1.895731
(Iteration 51 / 200) loss: 1.681091
(Iteration 61 / 200) loss: 1.487204
(Iteration 71 / 200) loss: 1.629973
(Epoch 2 / 5) train acc: 0.429000; val_acc: 0.350000
(Iteration 81 / 200) loss: 1.506686
(Iteration 91 / 200) loss: 1.610742
(Iteration 101 / 200) loss: 1.486124
(Iteration 111 / 200) loss: 1.559454
(Epoch 3 / 5) train acc: 0.492000; val_acc: 0.361000
(Iteration 121 / 200) loss: 1.497406
(Iteration 131 / 200) loss: 1.530736
(Iteration 141 / 200) loss: 1.550957
(Iteration 151 / 200) loss: 1.652046
(Epoch 4 / 5) train acc: 0.530000; val_acc: 0.361000
(Iteration 161 / 200) loss: 1.599574
(Iteration 171 / 200) loss: 1.401073
(Iteration 181 / 200) loss: 1.509365
(Iteration 191 / 200) loss: 1.365773
(Epoch 5 / 5) train acc: 0.531000; val_acc: 0.369000
```

```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:30:
```

```
MatplotlibDeprecationWarning: Adding an axes using the same arguments as a
previous axes currently reuses the earlier instance.  In a future version, a new
instance will always be created and returned.  Meanwhile, this warning can be
suppressed, and the future behavior ensured, by passing a unique label to each
axes instance.
```

```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:33:
```

```
MatplotlibDeprecationWarning: Adding an axes using the same arguments as a
```

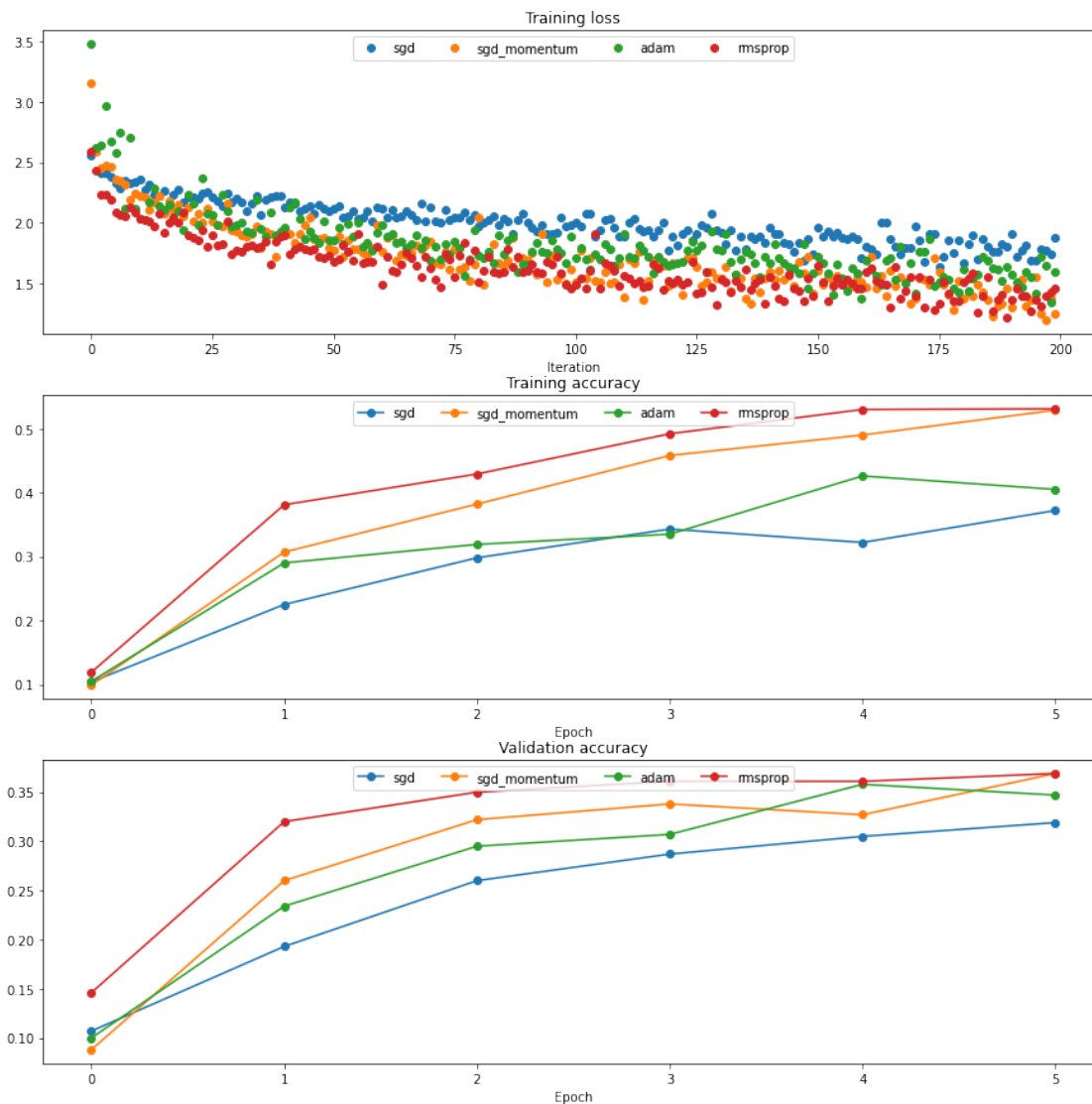
previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:36:
```

```
MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.
```

```
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:40:
```

```
MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.
```



### 13.1 Inline Question 3:

AdaGrad, 类似于 Adam, 是一个 per-parameter 优化方法, 它使用以下更新规则:

```
cache += dw**2
w += - learning_rate * dw / (np.sqrt(cache) + eps)
```

当使用 AdaGrad 训练一个网络时, 更新的值会变得非常小, 而且他的网络学习的非常慢。利用你对 AdaGrad 更新规则的了解, 解释为什么更新的值会变得非常小? Adam 会有同样的问题吗?

### 13.2 Answer:

根据上面的公式, adagrad 每步更新会越除越小。adam 结合了 adagrad 和动量的方法, 避免了这个问题。

## 14 训练一个效果足够好的模型!

在 CIFAR-10 上尽可能训练最好的全连接模型, 将最好的模型存储在 `best_model` 变量中。我们要求你在验证集上获得至少 50% 的准确性。

如果你细心的话, 应该是有可能得到 55% 以上精度的, 但我们不苛求你达到这么高的精度。在后面的作业上, 我们会要求你们在 CIFAR-10 上训练最好的卷积神经网络, 我们希望你们把精力放在卷积网络上, 而不是全连接网络上。

在做这部分之前完成 `BatchNormalization.ipynb` 和 `Dropout.ipynb` 可能会对你有帮助, 因为这些技术可以帮助你训练强大的模型。

```
[52]: best_model = None

#####
# TODO: Train the best FullyConnectedNet that you can on CIFAR-10. You might #
# find batch/layer normalization and dropout useful. Store your best model in #
# the best_model variable.                                                    #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

weight_scale = 4e-2
learning_rate = 6e-4
model = FullyConnectedNet([100, 100, 100, 100],
                           weight_scale=weight_scale, dtype=np.float64, reg=0.001)
solver = Solver(model, data,
                 print_every=100, num_epochs=6, batch_size=200,
                 update_rule='adam',
                 optim_config={
                     'learning_rate': learning_rate,
                 })
solver.train()
```

```
# accuracy is about 51%.
best_model = model

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
#####
#                                     END OF YOUR CODE                                     #
#####
```

```
(Iteration 1 / 1470) loss: 2.760619
(Epoch 0 / 6) train acc: 0.123000; val_acc: 0.132000
(Iteration 101 / 1470) loss: 1.940116
(Iteration 201 / 1470) loss: 1.814940
(Epoch 1 / 6) train acc: 0.449000; val_acc: 0.443000
(Iteration 301 / 1470) loss: 1.856808
(Iteration 401 / 1470) loss: 1.563786
(Epoch 2 / 6) train acc: 0.480000; val_acc: 0.473000
(Iteration 501 / 1470) loss: 1.600052
(Iteration 601 / 1470) loss: 1.504214
(Iteration 701 / 1470) loss: 1.457760
(Epoch 3 / 6) train acc: 0.507000; val_acc: 0.484000
(Iteration 801 / 1470) loss: 1.351943
(Iteration 901 / 1470) loss: 1.468692
(Epoch 4 / 6) train acc: 0.535000; val_acc: 0.493000
(Iteration 1001 / 1470) loss: 1.346348
(Iteration 1101 / 1470) loss: 1.329311
(Iteration 1201 / 1470) loss: 1.322137
(Epoch 5 / 6) train acc: 0.546000; val_acc: 0.495000
(Iteration 1301 / 1470) loss: 1.276124
(Iteration 1401 / 1470) loss: 1.386978
(Epoch 6 / 6) train acc: 0.583000; val_acc: 0.499000
```

## 15 测试你的模型!

在验证和测试集上运行您的最佳模型。验证集的准确率应达到 50% 以上。

```
[53]: y_test_pred = np.argmax(best_model.loss(data['X_test']), axis=1)
y_val_pred = np.argmax(best_model.loss(data['X_val']), axis=1)
print('Validation set accuracy: ', (y_val_pred == data['y_val']).mean())
print('Test set accuracy: ', (y_test_pred == data['y_test']).mean())
```

Validation set accuracy: 0.499

Test set accuracy: 0.49

## 16 重要

这里是作业的结尾处，请执行以下步骤：

1. 点击 File -> Save 或者用 control+s 组合键，确保你最新的 notebook 的作业已经保存到谷歌云。
2. 执行以下代码确保 .py 文件保存回你的谷歌云。

```
[56]: import os

FOLDER_TO_SAVE = os.path.join('drive/My Drive/', FOLDERNAME)
FILES_TO_SAVE = ['daseCV/classifiers/cnn.py', 'daseCV/classifiers/fc_net.py', \
                  'daseCV/layers.py', 'daseCV/optim.py', \
                  'daseCV/solver.py', 'daseCV/layer_utils.py', \
                  'daseCV/data_utils.py', 'daseCV/fast_layers.py']

for files in FILES_TO_SAVE:
    with open(os.path.join(FOLDER_TO_SAVE, '/'.join(files.split('/')[1:])), 'w') as f:
        f.write(''.join(open(files).readlines()))
```