

# Fully-Connected Neural Nets

In the previous homework you implemented a fully-connected two-layer neural network on CIFAR-10. The implementation was simple but not very modular since the loss and gradient were computed in a single monolithic function. This is manageable for a simple two-layer network, but would become impractical as we move to bigger models. Ideally we want to build networks using a more modular design so that we can implement different layer types in isolation and then snap them together into models with different architectures.

In this exercise we will implement fully-connected networks using a more modular approach. For each layer we will implement a `forward` and a `backward` function. The `forward` function will receive inputs, weights, and other parameters and will return both an output and a `cache` object storing data needed for the backward pass, like this:

```
def layer_forward(x, w):
    """ Receive inputs x and weights w """
    # Do some computations ...
    z = # ... some intermediate value
    # Do some more computations ...
    out = # the output

    cache = (x, w, z, out) # Values we need to compute gradients

    return out, cache
```

The backward pass will receive upstream derivatives and the `cache` object, and will return gradients with respect to the inputs and weights, like this:

```
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 = # Derivative of loss with respect to x
    dw = # Derivative of loss with respect to w

    return dx, dw
```

After implementing a bunch of layers this way, we will be able to easily combine them to build classifiers with different architectures.

In addition to implementing fully-connected networks of arbitrary depth, we will also explore different update rules for optimization, and introduce Dropout as a regularizer and Batch/Layer Normalization as a tool to more efficiently optimize deep networks.

```
In [1]: # As usual, a bit of setup
from __future__ import print_function
import time
import numpy as np
import matplotlib.pyplot as plt
from cs231n.classifiers.fc_net import *
from cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_grads
from cs231n.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
%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))))
```

```
In [2]: # 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,))
```

## Affine layer: forward

Open the file `cs231n/layers.py` and implement the `affine_forward` function.

Once you are done you can test your implementation by running the following:

```
In [3]: # Test the affine_forward function

num_inputs = 2
input_shape = (4, 5, 6)
output_dim = 3

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
```

## Affine layer: backward

Now implement the `affine_backward` function and test your implementation using numeric gradient checking.

```
In [4]: # Test the affine_backward function
np.random.seed(231)
x = np.random.randn(10, 2, 3)
w = np.random.randn(6, 5)
b = np.random.randn(5)
dout = np.random.randn(10, 5)

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)

# 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
```

## ReLU activation: forward

Implement the forward pass for the ReLU activation function in the `relu_forward` function and test your implementation using the following:

```
In [5]: # Test the relu_forward function

x = np.linspace(-0.5, 0.5, num=12).reshape(3, 4)

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,

# 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))

Testing relu_forward function:
difference:  4.999999798022158e-08
```

## ReLU activation: backward

Now implement the backward pass for the ReLU activation function in the `relu_backward` function and test your implementation using numeric gradient checking:

```
In [6]: 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
```

## Inline Question 1:

We've only asked you to implement ReLU, but there are a number of different activation functions that one could use in neural networks, each with its pros and cons. In particular, an issue commonly seen with activation functions is getting zero (or close to zero) gradient flow during backpropagation. Which of the following activation functions have this problem? If you consider these functions in the one dimensional case, what types of input would lead to this behaviour?

1. Sigmoid
2. ReLU

### 3. Leaky ReLU

## Answer:

[FILL THIS IN]

## "Sandwich" layers

There are some common patterns of layers that are frequently used in neural nets. For example, affine layers are frequently followed by a ReLU nonlinearity. To make these common patterns easy, we define several convenience layers in the file `cs231n/layer_utils.py`.

For now take a look at the `affine_relu_forward` and `affine_relu_backward` functions, and run the following to numerically gradient check the backward pass:

```
In [7]: from cs231n.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
```

## Loss layers: Softmax and SVM

You implemented these loss functions in the last assignment, so we'll give them to you for free here. You should still make sure you understand how they work by looking at the implementations in `cs231n/layers.py`.

You can make sure that the implementations are correct by running the following:

```

In [8]: 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 1e-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 1e-9
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

```

## Two-layer network

In the previous assignment you implemented a two-layer neural network in a single monolithic class. Now that you have implemented modular versions of the necessary layers, you will reimplement the two layer network using these modular implementations.

Open the file `cs231n/classifiers/fc_net.py` and complete the implementation of the `TwoLayerNet` class. This class will serve as a model for the other networks you will implement in this assignment, so read through it to make sure you understand the API. You can run the cell below to test your implementation.

```

In [9]: 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, 15.50004145, 16.18759162],
     [12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.05099822, 15.66781506, 16.28463189]])
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: 7.98e-08
b1 relative error: 1.35e-08
b2 relative error: 7.76e-10
```

In [ ]:

## Solver

In the previous assignment, the logic for training models was coupled to the models themselves. Following a more modular design, for this assignment we have split the logic for training models into a separate class.

Open the file `cs231n/solver.py` and read through it to familiarize yourself with the API. After doing so, use a `Solver` instance to train a `TwoLayerNet` that achieves at least 50% accuracy on the validation set.



```

In [10]: model = TwoLayerNet()
         #solver = None

#####
# TODO: Use a Solver instance to train a TwoLayerNet that achieves at least
# 50% accuracy on the validation set.
#####
#pass
data = get_CIFAR10_data()

solver = Solver(model, data, update_rule='sgd', optim_config={ 'learning_rate': 0.1,
                                                              lr_decay=0.95, num_epochs=10, batch_size=150, print_every=100 })
solver.train()

#####
#                               END OF YOUR CODE
#####

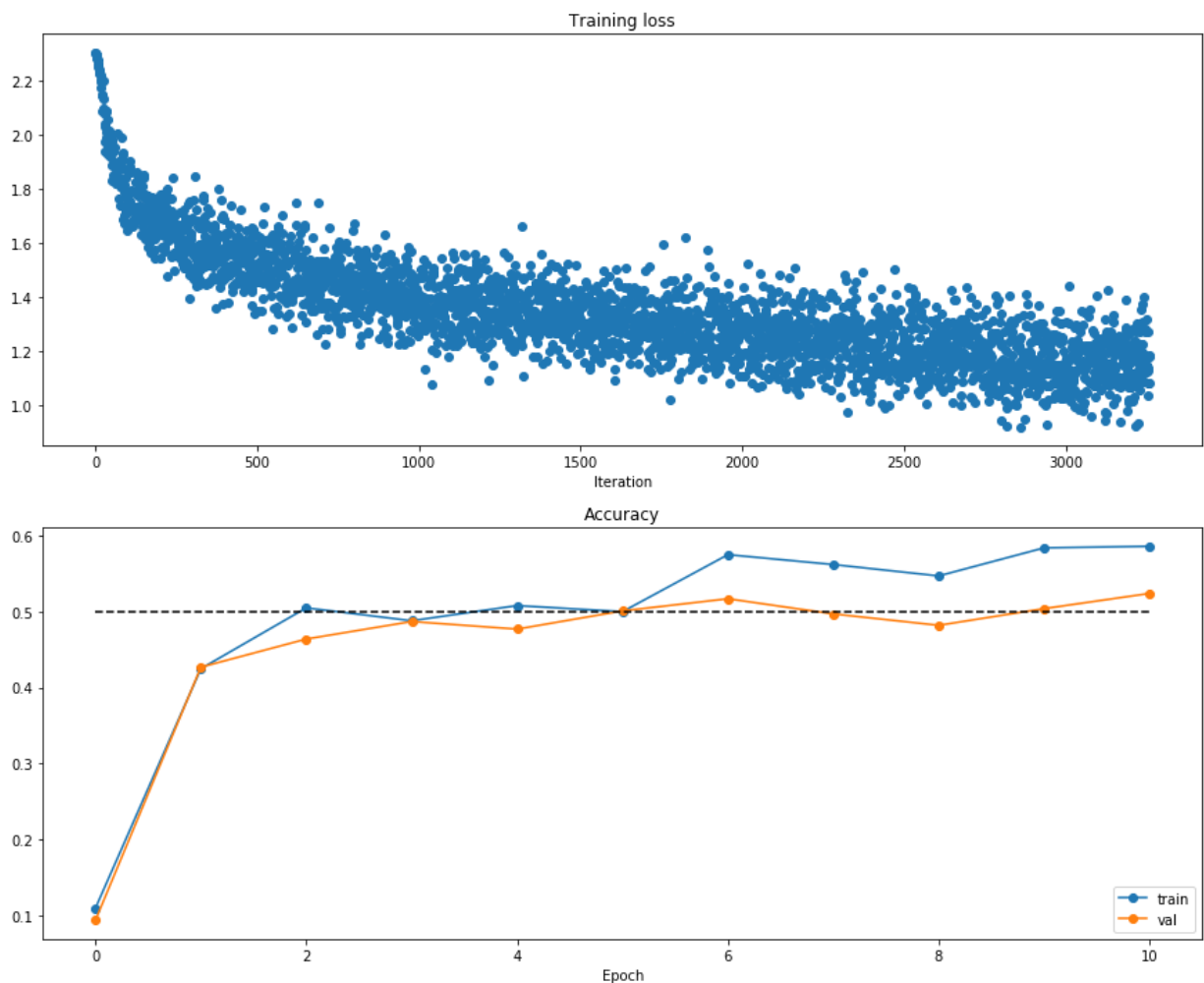
(Iteration 1 / 3260) loss: 2.304015
(Epoch 0 / 10) train acc: 0.110000; val_acc: 0.094000
(Iteration 11 / 3260) loss: 2.243340
(Iteration 21 / 3260) loss: 2.150927
(Iteration 31 / 3260) loss: 1.975819
(Iteration 41 / 3260) loss: 1.989264
(Iteration 51 / 3260) loss: 1.982700
(Iteration 61 / 3260) loss: 1.830876
(Iteration 71 / 3260) loss: 1.765306
(Iteration 81 / 3260) loss: 1.827145
(Iteration 91 / 3260) loss: 1.701437
(Iteration 101 / 3260) loss: 1.744181
(Iteration 111 / 3260) loss: 1.796182
(Iteration 121 / 3260) loss: 1.686499
(Iteration 131 / 3260) loss: 1.733715
(Iteration 141 / 3260) loss: 1.697911
(Iteration 151 / 3260) loss: 1.771549
(Iteration 161 / 3260) loss: 1.638439
(Iteration 171 / 3260) loss: 1.627800

```

In [11]: *# 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()
```



## Multilayer network

Next you will implement a fully-connected network with an arbitrary number of hidden layers.

Read through the `FullyConnectedNet` class in the file `cs231n/classifiers/fc_net.py`.

Implement the initialization, the forward pass, and the backward pass. For the moment don't worry about implementing dropout or batch/layer normalization; we will add those features soon.

## Initial loss and gradient check

As a sanity check, run the following to check the initial loss and to gradient check the network both with and without regularization. Do the initial losses seem reasonable?

For gradient checking, you should expect to see errors around  $1e-7$  or less.

```
In [12]: 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)
        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: 5.940411485412347
W1 relative error: 7.36e-09
W2 relative error: 6.87e-08
W3 relative error: 3.80e-07
b1 relative error: 1.48e-08
b2 relative error: 1.72e-09
b3 relative error: 1.80e-10
```

As another sanity check, make sure you can overfit a small dataset of 50 images. First we will try a three-layer network with 100 units in each hidden layer. In the following cell, tweak the learning rate and initialization scale to overfit and achieve 100% training accuracy within 20 epochs.

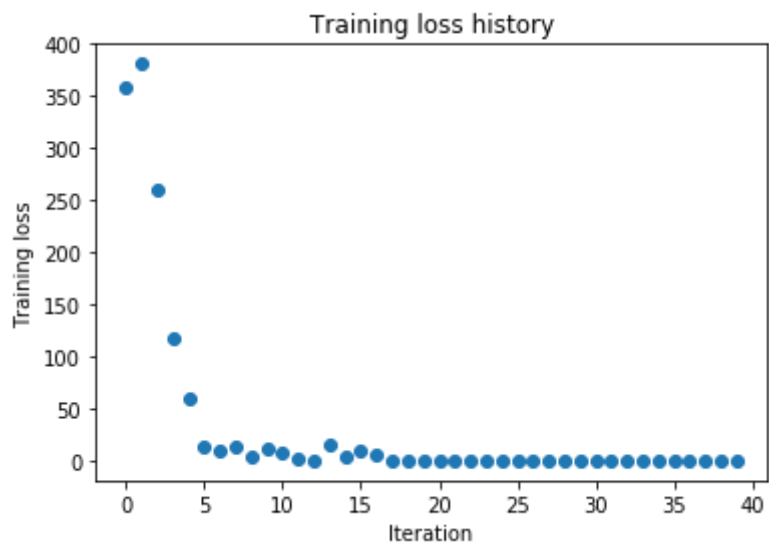
In [13]: *# TODO: Use a three-layer Net to overfit 50 training examples by  
# tweaking just the learning rate and initialization scale.*

```
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-1
learning_rate = 1e-3
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: 357.428290
(Epoch 0 / 20) train acc: 0.220000; val_acc: 0.111000
(Epoch 1 / 20) train acc: 0.380000; val_acc: 0.141000
(Epoch 2 / 20) train acc: 0.520000; val_acc: 0.138000
(Epoch 3 / 20) train acc: 0.740000; val_acc: 0.130000
(Epoch 4 / 20) train acc: 0.820000; val_acc: 0.153000
(Epoch 5 / 20) train acc: 0.860000; val_acc: 0.175000
(Iteration 11 / 40) loss: 6.726589
(Epoch 6 / 20) train acc: 0.940000; val_acc: 0.163000
(Epoch 7 / 20) train acc: 0.960000; val_acc: 0.166000
(Epoch 8 / 20) train acc: 0.960000; val_acc: 0.164000
(Epoch 9 / 20) train acc: 0.980000; val_acc: 0.162000
(Epoch 10 / 20) train acc: 0.980000; val_acc: 0.162000
(Iteration 21 / 40) loss: 0.800243
(Epoch 11 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 12 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 13 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 14 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.158000
(Iteration 31 / 40) loss: 0.000000
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.158000
```



Now try to use a five-layer network with 100 units on each layer to overfit 50 training examples. Again you will have to adjust the learning rate and weight initialization, but you should be able to achieve 100% training accuracy within 20 epochs.

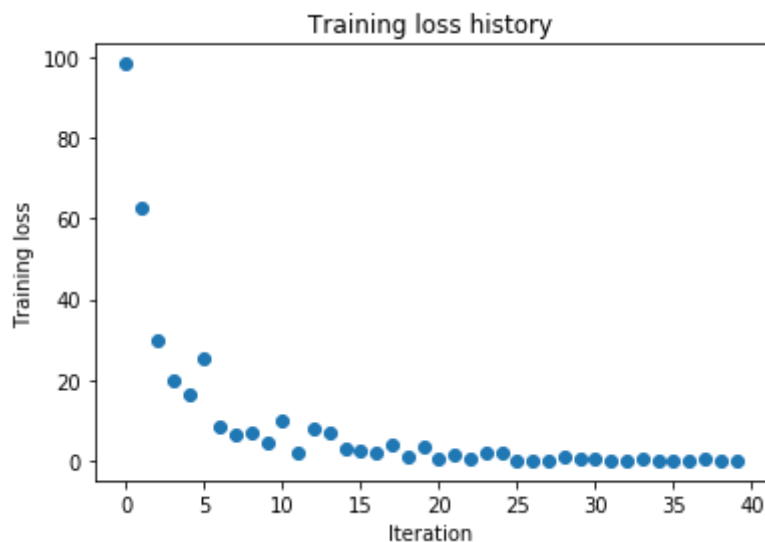
In [14]: *# TODO: Use a five-layer Net to overfit 50 training examples by  
# tweaking just the learning rate and initialization scale.*

```
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'],
}

learning_rate = 2e-4
weight_scale = 9e-2
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: 98.325850
(Epoch 0 / 20) train acc: 0.160000; val_acc: 0.125000
(Epoch 1 / 20) train acc: 0.220000; val_acc: 0.127000
(Epoch 2 / 20) train acc: 0.300000; val_acc: 0.114000
(Epoch 3 / 20) train acc: 0.380000; val_acc: 0.126000
(Epoch 4 / 20) train acc: 0.480000; val_acc: 0.120000
(Epoch 5 / 20) train acc: 0.600000; val_acc: 0.099000
(Iteration 11 / 40) loss: 10.025033
(Epoch 6 / 20) train acc: 0.680000; val_acc: 0.128000
(Epoch 7 / 20) train acc: 0.720000; val_acc: 0.123000
(Epoch 8 / 20) train acc: 0.820000; val_acc: 0.125000
(Epoch 9 / 20) train acc: 0.820000; val_acc: 0.114000
(Epoch 10 / 20) train acc: 0.880000; val_acc: 0.126000
(Iteration 21 / 40) loss: 0.481338
(Epoch 11 / 20) train acc: 0.900000; val_acc: 0.121000
(Epoch 12 / 20) train acc: 0.920000; val_acc: 0.122000
(Epoch 13 / 20) train acc: 0.940000; val_acc: 0.121000
(Epoch 14 / 20) train acc: 0.960000; val_acc: 0.121000
(Epoch 15 / 20) train acc: 0.960000; val_acc: 0.123000
(Iteration 31 / 40) loss: 0.699036
(Epoch 16 / 20) train acc: 0.940000; val_acc: 0.118000
(Epoch 17 / 20) train acc: 0.960000; val_acc: 0.125000
(Epoch 18 / 20) train acc: 0.980000; val_acc: 0.122000
(Epoch 19 / 20) train acc: 0.980000; val_acc: 0.119000
(Epoch 20 / 20) train acc: 0.980000; val_acc: 0.119000
```



## Inline Question 2:

Did you notice anything about the comparative difficulty of training the three-layer net vs training the five layer net? In particular, based on your experience, which network seemed more sensitive to the initialization scale? Why do you think that is the case?

## Answer:

[FILL THIS IN]

## Update rules

So far we have used vanilla stochastic gradient descent (SGD) as our update rule. More sophisticated update rules can make it easier to train deep networks. We will implement a few of the most commonly used update rules and compare them to vanilla SGD.

## SGD+Momentum

Stochastic gradient descent with momentum is a widely used update rule that tends to make deep networks converge faster than vanilla stochastic gradient descent. See the Momentum Update section at <http://cs231n.github.io/neural-networks-3/#sgd> (<http://cs231n.github.io/neural-networks-3/#sgd>) for more information.

Open the file `cs231n/optim.py` and read the documentation at the top of the file to make sure you understand the API. Implement the SGD+momentum update rule in the function `sgd_momentum` and run the following to check your implementation. You should see errors less than  $e-8$ .

```

In [15]: from cs231n.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

```

Once you have done so, run the following to train a six-layer network with both SGD and SGD+momentum. You should see the SGD+momentum update rule converge faster.



```

In [16]: 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': 1e-2,
                    },
                    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 sgd
(Iteration 1 / 200) loss: 2.559978
(Epoch 0 / 5) train acc: 0.103000; val_acc: 0.108000
(Iteration 11 / 200) loss: 2.291086

```

```
(Iteration 21 / 200) loss: 2.153591
(Iteration 31 / 200) loss: 2.082693
(Epoch 1 / 5) train acc: 0.277000; val_acc: 0.242000
(Iteration 41 / 200) loss: 2.004171
(Iteration 51 / 200) loss: 2.010409
(Iteration 61 / 200) loss: 2.023753
(Iteration 71 / 200) loss: 2.026621
(Epoch 2 / 5) train acc: 0.352000; val_acc: 0.312000
(Iteration 81 / 200) loss: 1.807164
(Iteration 91 / 200) loss: 1.915726
(Iteration 101 / 200) loss: 1.921843
(Iteration 111 / 200) loss: 1.708426
(Epoch 3 / 5) train acc: 0.399000; val_acc: 0.318000
(Iteration 121 / 200) loss: 1.699806
(Iteration 131 / 200) loss: 1.771857
(Iteration 141 / 200) loss: 1.790064
(Iteration 151 / 200) loss: 1.820861
(Epoch 4 / 5) train acc: 0.430000; val_acc: 0.324000
(Iteration 161 / 200) loss: 1.627657
(Iteration 171 / 200) loss: 1.899449
(Iteration 181 / 200) loss: 1.551447
(Iteration 191 / 200) loss: 1.698943
(Epoch 5 / 5) train acc: 0.437000; val_acc: 0.330000
```

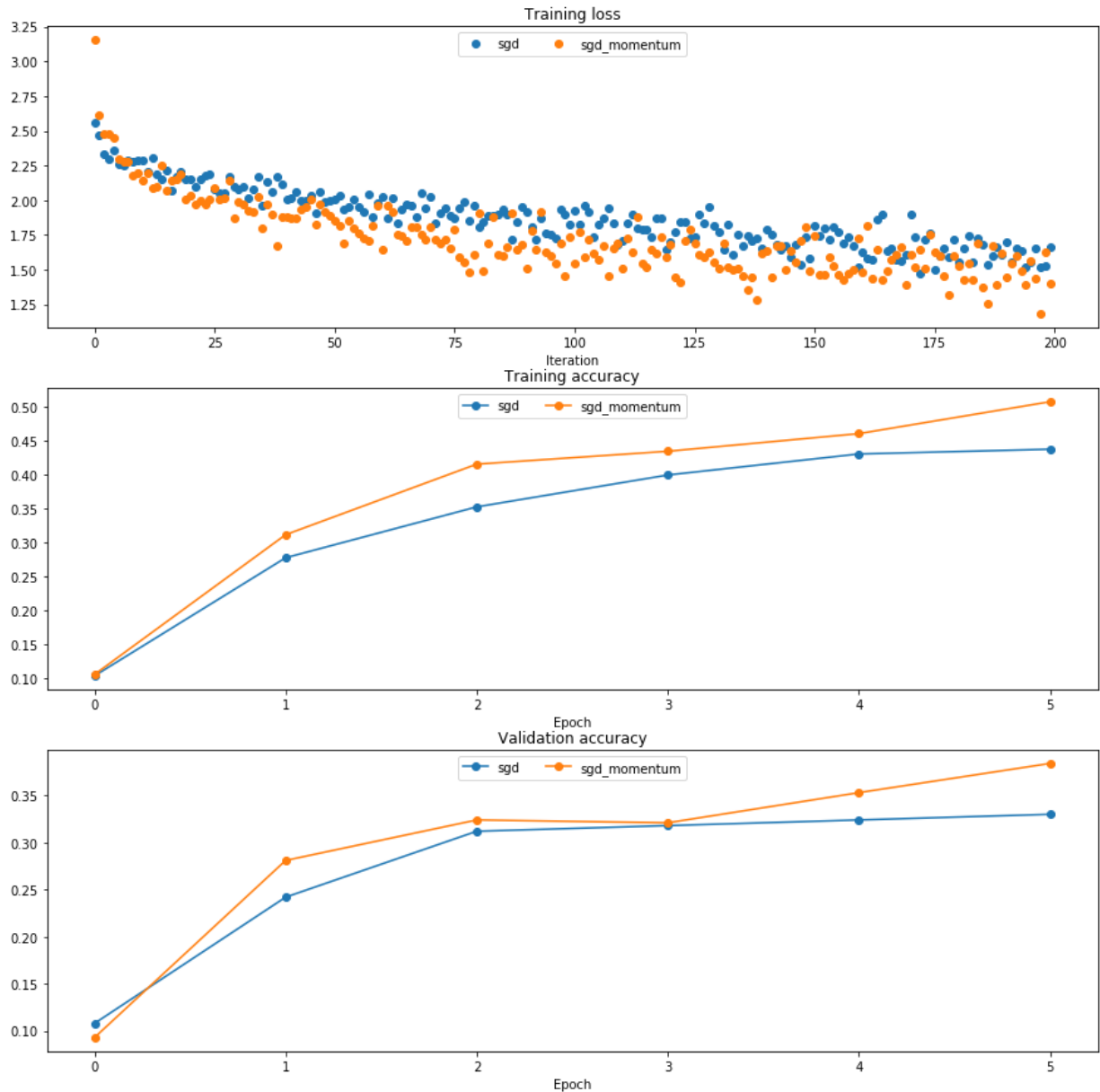
running with `sgd_momentum`

```
(Iteration 1 / 200) loss: 3.153777
(Epoch 0 / 5) train acc: 0.105000; val_acc: 0.093000
(Iteration 11 / 200) loss: 2.145874
(Iteration 21 / 200) loss: 2.032563
(Iteration 31 / 200) loss: 1.985848
(Epoch 1 / 5) train acc: 0.311000; val_acc: 0.281000
(Iteration 41 / 200) loss: 1.882353
(Iteration 51 / 200) loss: 1.855372
(Iteration 61 / 200) loss: 1.649133
(Iteration 71 / 200) loss: 1.806432
(Epoch 2 / 5) train acc: 0.415000; val_acc: 0.324000
(Iteration 81 / 200) loss: 1.907840
(Iteration 91 / 200) loss: 1.510681
(Iteration 101 / 200) loss: 1.546872
(Iteration 111 / 200) loss: 1.512047
(Epoch 3 / 5) train acc: 0.434000; val_acc: 0.321000
(Iteration 121 / 200) loss: 1.677301
(Iteration 131 / 200) loss: 1.504686
(Iteration 141 / 200) loss: 1.633253
(Iteration 151 / 200) loss: 1.745081
(Epoch 4 / 5) train acc: 0.460000; val_acc: 0.353000
(Iteration 161 / 200) loss: 1.485411
(Iteration 171 / 200) loss: 1.610416
(Iteration 181 / 200) loss: 1.528331
(Iteration 191 / 200) loss: 1.449159
(Epoch 5 / 5) train acc: 0.507000; val_acc: 0.384000
```

/home/shared/anaconda3/lib/python3.6/site-packages/matplotlib/cbook/deprecation.py:106: 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. Me

anwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

```
warnings.warn(message, mplDeprecation, stacklevel=1)
```



## RMSProp and Adam

RMSProp [1] and Adam [2] are update rules that set per-parameter learning rates by using a running average of the second moments of gradients.

In the file `cs231n/optim.py`, implement the RMSProp update rule in the `rmsprop` function and implement the Adam update rule in the `adam` function, and check your implementations using the tests below.

**NOTE:** Please implement the *complete* Adam update rule (with the bias correction mechanism), not the first simplified version mentioned in the course notes.

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

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

```
In [17]: # Test RMSProp implementation
from cs231n.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: 9.502645229894295e-08
cache error: 2.6477955807156126e-09
```

```

In [18]: # Test Adam implementation
from cs231n.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: 1.139887467333134e-07
v error: 4.208314038113071e-09
m error: 4.214963193114416e-09

```

Once you have debugged your RMSProp and Adam implementations, run the following to train a pair of deep networks using these new update rules:

```

In [19]: 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: 0.126000; val_acc: 0.110000
(Iteration 11 / 200) loss: 2.015214
(Iteration 21 / 200) loss: 2.190157
(Iteration 31 / 200) loss: 1.785418
(Epoch 1 / 5) train acc: 0.387000; val_acc: 0.335000
(Iteration 41 / 200) loss: 1.743079
(Iteration 51 / 200) loss: 1.731332
(Iteration 61 / 200) loss: 1.992775
(Iteration 71 / 200) loss: 1.593985
(Epoch 2 / 5) train acc: 0.436000; val_acc: 0.360000
(Iteration 81 / 200) loss: 1.626604

```

```
(Iteration 91 / 200) loss: 1.494471
(Iteration 101 / 200) loss: 1.445403
(Iteration 111 / 200) loss: 1.398111
(Epoch 3 / 5) train acc: 0.499000; val_acc: 0.367000
(Iteration 121 / 200) loss: 1.222326
(Iteration 131 / 200) loss: 1.414736
```

### Inline Question 3:

AdaGrad, like Adam, is a per-parameter optimization method that uses the following update rule:

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

John notices that when he was training a network with AdaGrad that the updates became very small, and that his network was learning slowly. Using your knowledge of the AdaGrad update rule, why do you think the updates would become very small? Would Adam have the same issue?

### Answer:

## Train a good model!

Train the best fully-connected model that you can on CIFAR-10, storing your best model in the `best_model` variable. We require you to get at least 50% accuracy on the validation set using a fully-connected net.

If you are careful it should be possible to get accuracies above 55%, but we don't require it for this part and won't assign extra credit for doing so. Later in the assignment we will ask you to train the best convolutional network that you can on CIFAR-10, and we would prefer that you spend your effort working on convolutional nets rather than fully-connected nets.

You might find it useful to complete the `BatchNormalization.ipynb` and `Dropout.ipynb` notebooks before completing this part, since those techniques can help you train powerful models.

```

In [20]: 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
# the best_model variable.
#####
#pass

data = get_CIFAR10_data()

best_model = FullyConnectedNet([100, 100, 100, 100], dropout=0.8, normalizat

solver = Solver(best_model, data,
                 num_epochs=10, batch_size=100,
                 update_rule='adam',
                 optim_config={
                     'learning_rate': 4e-4
                 },
                 verbose=True)
solver.train()
print()

#####
#                               END OF YOUR CODE
#####

```

```

(Iteration 1 / 4900) loss: 2.312783
(Epoch 0 / 10) train acc: 0.082000; val_acc: 0.103000
(Iteration 11 / 4900) loss: 2.281196
(Iteration 21 / 4900) loss: 2.236326
(Iteration 31 / 4900) loss: 2.195081
(Iteration 41 / 4900) loss: 2.168172
(Iteration 51 / 4900) loss: 2.125695
(Iteration 61 / 4900) loss: 2.143014
(Iteration 71 / 4900) loss: 1.957847
(Iteration 81 / 4900) loss: 2.133950
(Iteration 91 / 4900) loss: 2.027694
(Iteration 101 / 4900) loss: 1.998782
(Iteration 111 / 4900) loss: 1.940760
(Iteration 121 / 4900) loss: 1.954139
(Iteration 131 / 4900) loss: 1.975175
(Iteration 141 / 4900) loss: 1.930862
(Iteration 151 / 4900) loss: 1.986191
(Iteration 161 / 4900) loss: 1.930245
(Iteration 171 / 4900) loss: 1.854939
(Iteration 181 / 4900) loss: 1.940928
(Iteration 191 / 4900) loss: 1.906350
(Iteration 201 / 4900) loss: 1.793624
(Iteration 211 / 4900) loss: 1.908489
(Iteration 221 / 4900) loss: 1.985113
(Iteration 231 / 4900) loss: 1.809485
(Iteration 241 / 4900) loss: 1.868577
(Iteration 251 / 4900) loss: 1.686672
(Iteration 261 / 4900) loss: 1.914262
(Iteration 271 / 4900) loss: 1.800328
(Iteration 281 / 4900) loss: 1.733879
(Iteration 291 / 4900) loss: 1.786188

```



```
(Iteration 301 / 4900) loss: 1.808115
(Iteration 311 / 4900) loss: 1.654289
(Iteration 321 / 4900) loss: 1.747490
(Iteration 331 / 4900) loss: 1.882386
(Iteration 341 / 4900) loss: 1.771834
(Iteration 351 / 4900) loss: 1.859555
(Iteration 361 / 4900) loss: 1.804786
(Iteration 371 / 4900) loss: 1.653695
(Iteration 381 / 4900) loss: 1.816778
(Iteration 391 / 4900) loss: 1.673257
(Iteration 401 / 4900) loss: 1.643267
(Iteration 411 / 4900) loss: 1.695853
(Iteration 421 / 4900) loss: 1.633278
(Iteration 431 / 4900) loss: 1.871397
(Iteration 441 / 4900) loss: 1.640895
(Iteration 451 / 4900) loss: 1.700059
(Iteration 461 / 4900) loss: 1.842141
(Iteration 471 / 4900) loss: 1.673366
(Iteration 481 / 4900) loss: 1.688936
(Epoch 1 / 10) train acc: 0.418000; val_acc: 0.413000
(Iteration 491 / 4900) loss: 1.635011
(Iteration 501 / 4900) loss: 1.657575
(Iteration 511 / 4900) loss: 1.947451
(Iteration 521 / 4900) loss: 1.827936
(Iteration 531 / 4900) loss: 1.685695
(Iteration 541 / 4900) loss: 1.555308
(Iteration 551 / 4900) loss: 1.752016
(Iteration 561 / 4900) loss: 1.573506
(Iteration 571 / 4900) loss: 1.759551
(Iteration 581 / 4900) loss: 1.705281
(Iteration 591 / 4900) loss: 1.536525
(Iteration 601 / 4900) loss: 1.871777
(Iteration 611 / 4900) loss: 1.653273
(Iteration 621 / 4900) loss: 1.680849
(Iteration 631 / 4900) loss: 1.715795
(Iteration 641 / 4900) loss: 1.600358
(Iteration 651 / 4900) loss: 1.594983
(Iteration 661 / 4900) loss: 1.608028
(Iteration 671 / 4900) loss: 1.640727
(Iteration 681 / 4900) loss: 1.743578
(Iteration 691 / 4900) loss: 1.705601
(Iteration 701 / 4900) loss: 1.671087
(Iteration 711 / 4900) loss: 1.738299
(Iteration 721 / 4900) loss: 1.586689
(Iteration 731 / 4900) loss: 1.562898
(Iteration 741 / 4900) loss: 1.688164
(Iteration 751 / 4900) loss: 1.726936
(Iteration 761 / 4900) loss: 1.783545
(Iteration 771 / 4900) loss: 1.767448
(Iteration 781 / 4900) loss: 1.612397
(Iteration 791 / 4900) loss: 1.656693
(Iteration 801 / 4900) loss: 1.692088
(Iteration 811 / 4900) loss: 1.825692
(Iteration 821 / 4900) loss: 1.582759
(Iteration 831 / 4900) loss: 1.530102
(Iteration 841 / 4900) loss: 1.725787
(Iteration 851 / 4900) loss: 1.648151
```

```
(Iteration 861 / 4900) loss: 1.632727
(Iteration 871 / 4900) loss: 1.722813
(Iteration 881 / 4900) loss: 1.521523
(Iteration 891 / 4900) loss: 1.576306
(Iteration 901 / 4900) loss: 1.770597
(Iteration 911 / 4900) loss: 1.831992
(Iteration 921 / 4900) loss: 1.575039
(Iteration 931 / 4900) loss: 1.392259
(Iteration 941 / 4900) loss: 1.772898
(Iteration 951 / 4900) loss: 1.537475
(Iteration 961 / 4900) loss: 1.514005
(Iteration 971 / 4900) loss: 1.712311
(Epoch 2 / 10) train acc: 0.459000; val_acc: 0.458000
(Iteration 981 / 4900) loss: 1.747870
(Iteration 991 / 4900) loss: 1.908510
(Iteration 1001 / 4900) loss: 1.625330
(Iteration 1011 / 4900) loss: 1.649811
(Iteration 1021 / 4900) loss: 1.648202
(Iteration 1031 / 4900) loss: 1.585096
(Iteration 1041 / 4900) loss: 1.532263
(Iteration 1051 / 4900) loss: 1.713584
(Iteration 1061 / 4900) loss: 1.788891
(Iteration 1071 / 4900) loss: 1.824980
(Iteration 1081 / 4900) loss: 1.699505
(Iteration 1091 / 4900) loss: 1.588897
(Iteration 1101 / 4900) loss: 1.679769
(Iteration 1111 / 4900) loss: 1.565877
(Iteration 1121 / 4900) loss: 1.408617
(Iteration 1131 / 4900) loss: 1.556175
(Iteration 1141 / 4900) loss: 1.490858
(Iteration 1151 / 4900) loss: 1.680202
(Iteration 1161 / 4900) loss: 1.610204
(Iteration 1171 / 4900) loss: 1.673611
(Iteration 1181 / 4900) loss: 1.515057
(Iteration 1191 / 4900) loss: 1.710998
(Iteration 1201 / 4900) loss: 1.663781
(Iteration 1211 / 4900) loss: 1.463776
(Iteration 1221 / 4900) loss: 1.575967
(Iteration 1231 / 4900) loss: 1.715743
(Iteration 1241 / 4900) loss: 1.383559
(Iteration 1251 / 4900) loss: 1.591899
(Iteration 1261 / 4900) loss: 1.705413
(Iteration 1271 / 4900) loss: 1.427769
(Iteration 1281 / 4900) loss: 1.615685
(Iteration 1291 / 4900) loss: 1.723740
(Iteration 1301 / 4900) loss: 1.528576
(Iteration 1311 / 4900) loss: 1.591839
(Iteration 1321 / 4900) loss: 1.614699
(Iteration 1331 / 4900) loss: 1.633935
(Iteration 1341 / 4900) loss: 1.597452
(Iteration 1351 / 4900) loss: 1.744962
(Iteration 1361 / 4900) loss: 1.426230
(Iteration 1371 / 4900) loss: 1.570380
(Iteration 1381 / 4900) loss: 1.458848
(Iteration 1391 / 4900) loss: 1.481594
(Iteration 1401 / 4900) loss: 1.680717
(Iteration 1411 / 4900) loss: 1.609862
```

```
(Iteration 1421 / 4900) loss: 1.671969
(Iteration 1431 / 4900) loss: 1.651903
(Iteration 1441 / 4900) loss: 1.554943
(Iteration 1451 / 4900) loss: 1.500102
(Iteration 1461 / 4900) loss: 1.591759
(Epoch 3 / 10) train acc: 0.508000; val_acc: 0.484000
(Iteration 1471 / 4900) loss: 1.640040
(Iteration 1481 / 4900) loss: 1.541870
(Iteration 1491 / 4900) loss: 1.689369
(Iteration 1501 / 4900) loss: 1.446454
(Iteration 1511 / 4900) loss: 1.545516
(Iteration 1521 / 4900) loss: 1.644365
(Iteration 1531 / 4900) loss: 1.464418
(Iteration 1541 / 4900) loss: 1.489686
(Iteration 1551 / 4900) loss: 1.472097
(Iteration 1561 / 4900) loss: 1.566960
(Iteration 1571 / 4900) loss: 1.710042
(Iteration 1581 / 4900) loss: 1.574410
(Iteration 1591 / 4900) loss: 1.572487
(Iteration 1601 / 4900) loss: 1.438016
(Iteration 1611 / 4900) loss: 1.621301
(Iteration 1621 / 4900) loss: 1.529721
(Iteration 1631 / 4900) loss: 1.438909
(Iteration 1641 / 4900) loss: 1.389319
(Iteration 1651 / 4900) loss: 1.699098
(Iteration 1661 / 4900) loss: 1.616598
(Iteration 1671 / 4900) loss: 1.561009
(Iteration 1681 / 4900) loss: 1.530746
(Iteration 1691 / 4900) loss: 1.561213
(Iteration 1701 / 4900) loss: 1.547092
(Iteration 1711 / 4900) loss: 1.421806
(Iteration 1721 / 4900) loss: 1.484774
(Iteration 1731 / 4900) loss: 1.553614
(Iteration 1741 / 4900) loss: 1.615050
(Iteration 1751 / 4900) loss: 1.615691
(Iteration 1761 / 4900) loss: 1.544079
(Iteration 1771 / 4900) loss: 1.505500
(Iteration 1781 / 4900) loss: 1.439436
(Iteration 1791 / 4900) loss: 1.579443
(Iteration 1801 / 4900) loss: 1.634896
(Iteration 1811 / 4900) loss: 1.458149
(Iteration 1821 / 4900) loss: 1.447337
(Iteration 1831 / 4900) loss: 1.544582
(Iteration 1841 / 4900) loss: 1.543989
(Iteration 1851 / 4900) loss: 1.684127
(Iteration 1861 / 4900) loss: 1.415639
(Iteration 1871 / 4900) loss: 1.552438
(Iteration 1881 / 4900) loss: 1.567776
(Iteration 1891 / 4900) loss: 1.594695
(Iteration 1901 / 4900) loss: 1.693023
(Iteration 1911 / 4900) loss: 1.659982
(Iteration 1921 / 4900) loss: 1.547112
(Iteration 1931 / 4900) loss: 1.540765
(Iteration 1941 / 4900) loss: 1.663216
(Iteration 1951 / 4900) loss: 1.498218
(Epoch 4 / 10) train acc: 0.503000; val_acc: 0.500000
(Iteration 1961 / 4900) loss: 1.329362
```

```
(Iteration 1971 / 4900) loss: 1.372321
(Iteration 1981 / 4900) loss: 1.522879
(Iteration 1991 / 4900) loss: 1.494375
(Iteration 2001 / 4900) loss: 1.460619
(Iteration 2011 / 4900) loss: 1.654206
(Iteration 2021 / 4900) loss: 1.410848
(Iteration 2031 / 4900) loss: 1.417846
(Iteration 2041 / 4900) loss: 1.525596
(Iteration 2051 / 4900) loss: 1.716102

(Iteration 2061 / 4900) loss: 1.701369
(Iteration 2071 / 4900) loss: 1.509935
(Iteration 2081 / 4900) loss: 1.709408
(Iteration 2091 / 4900) loss: 1.554359
(Iteration 2101 / 4900) loss: 1.581738
(Iteration 2111 / 4900) loss: 1.695055
(Iteration 2121 / 4900) loss: 1.738292
(Iteration 2131 / 4900) loss: 1.547154
(Iteration 2141 / 4900) loss: 1.539329
(Iteration 2151 / 4900) loss: 1.564831
(Iteration 2161 / 4900) loss: 1.557065
(Iteration 2171 / 4900) loss: 1.400740
(Iteration 2181 / 4900) loss: 1.465492
(Iteration 2191 / 4900) loss: 1.486901
(Iteration 2201 / 4900) loss: 1.502205
(Iteration 2211 / 4900) loss: 1.504770
(Iteration 2221 / 4900) loss: 1.459327
(Iteration 2231 / 4900) loss: 1.537847
(Iteration 2241 / 4900) loss: 1.673207
(Iteration 2251 / 4900) loss: 1.301671
(Iteration 2261 / 4900) loss: 1.423535
(Iteration 2271 / 4900) loss: 1.703118
(Iteration 2281 / 4900) loss: 1.655642
(Iteration 2291 / 4900) loss: 1.505084
(Iteration 2301 / 4900) loss: 1.462715
(Iteration 2311 / 4900) loss: 1.419914
(Iteration 2321 / 4900) loss: 1.341203
(Iteration 2331 / 4900) loss: 1.344862
(Iteration 2341 / 4900) loss: 1.466331
(Iteration 2351 / 4900) loss: 1.341917
(Iteration 2361 / 4900) loss: 1.516055
(Iteration 2371 / 4900) loss: 1.353075
(Iteration 2381 / 4900) loss: 1.384555
(Iteration 2391 / 4900) loss: 1.546158
(Iteration 2401 / 4900) loss: 1.410270
(Iteration 2411 / 4900) loss: 1.553441
(Iteration 2421 / 4900) loss: 1.336065
(Iteration 2431 / 4900) loss: 1.488361
(Iteration 2441 / 4900) loss: 1.570992
(Epoch 5 / 10) train acc: 0.520000; val_acc: 0.506000
(Iteration 2451 / 4900) loss: 1.365809
(Iteration 2461 / 4900) loss: 1.407332
(Iteration 2471 / 4900) loss: 1.422069
(Iteration 2481 / 4900) loss: 1.479580
(Iteration 2491 / 4900) loss: 1.457647
(Iteration 2501 / 4900) loss: 1.628578
(Iteration 2511 / 4900) loss: 1.542198
```

```
(Iteration 2521 / 4900) loss: 1.509316
(Iteration 2531 / 4900) loss: 1.577862
(Iteration 2541 / 4900) loss: 1.551175
(Iteration 2551 / 4900) loss: 1.488764
(Iteration 2561 / 4900) loss: 1.540468
(Iteration 2571 / 4900) loss: 1.461957
(Iteration 2581 / 4900) loss: 1.404701
(Iteration 2591 / 4900) loss: 1.417520
(Iteration 2601 / 4900) loss: 1.514831
(Iteration 2611 / 4900) loss: 1.490023
(Iteration 2621 / 4900) loss: 1.451329
(Iteration 2631 / 4900) loss: 1.485876
(Iteration 2641 / 4900) loss: 1.588990
(Iteration 2651 / 4900) loss: 1.369174
(Iteration 2661 / 4900) loss: 1.465471
(Iteration 2671 / 4900) loss: 1.458743
(Iteration 2681 / 4900) loss: 1.488607
(Iteration 2691 / 4900) loss: 1.454031
(Iteration 2701 / 4900) loss: 1.543523
(Iteration 2711 / 4900) loss: 1.581099
(Iteration 2721 / 4900) loss: 1.609772
(Iteration 2731 / 4900) loss: 1.346582
(Iteration 2741 / 4900) loss: 1.492978
(Iteration 2751 / 4900) loss: 1.434017
(Iteration 2761 / 4900) loss: 1.384382
(Iteration 2771 / 4900) loss: 1.560223
(Iteration 2781 / 4900) loss: 1.630149
(Iteration 2791 / 4900) loss: 1.524025
(Iteration 2801 / 4900) loss: 1.392865
(Iteration 2811 / 4900) loss: 1.468193
(Iteration 2821 / 4900) loss: 1.464191
(Iteration 2831 / 4900) loss: 1.460037
(Iteration 2841 / 4900) loss: 1.417209
(Iteration 2851 / 4900) loss: 1.341453
(Iteration 2861 / 4900) loss: 1.360330
(Iteration 2871 / 4900) loss: 1.576963
(Iteration 2881 / 4900) loss: 1.599018
(Iteration 2891 / 4900) loss: 1.671984
(Iteration 2901 / 4900) loss: 1.700702
(Iteration 2911 / 4900) loss: 1.326301
(Iteration 2921 / 4900) loss: 1.524570
(Iteration 2931 / 4900) loss: 1.528632
(Epoch 6 / 10) train acc: 0.514000; val_acc: 0.522000
(Iteration 2941 / 4900) loss: 1.490302
(Iteration 2951 / 4900) loss: 1.530544
(Iteration 2961 / 4900) loss: 1.582536
(Iteration 2971 / 4900) loss: 1.618277
(Iteration 2981 / 4900) loss: 1.316296
(Iteration 2991 / 4900) loss: 1.453802
(Iteration 3001 / 4900) loss: 1.244796
(Iteration 3011 / 4900) loss: 1.360696
(Iteration 3021 / 4900) loss: 1.207763
(Iteration 3031 / 4900) loss: 1.392070
(Iteration 3041 / 4900) loss: 1.463814
(Iteration 3051 / 4900) loss: 1.327665
(Iteration 3061 / 4900) loss: 1.533577
(Iteration 3071 / 4900) loss: 1.498893
```

```
(Iteration 3081 / 4900) loss: 1.478662
(Iteration 3091 / 4900) loss: 1.337832
(Iteration 3101 / 4900) loss: 1.459955
(Iteration 3111 / 4900) loss: 1.579114
(Iteration 3121 / 4900) loss: 1.464574
(Iteration 3131 / 4900) loss: 1.565014
(Iteration 3141 / 4900) loss: 1.639127
(Iteration 3151 / 4900) loss: 1.416621
(Iteration 3161 / 4900) loss: 1.344334
(Iteration 3171 / 4900) loss: 1.228916
(Iteration 3181 / 4900) loss: 1.404916
(Iteration 3191 / 4900) loss: 1.247255
(Iteration 3201 / 4900) loss: 1.497720
(Iteration 3211 / 4900) loss: 1.523882
(Iteration 3221 / 4900) loss: 1.430716
(Iteration 3231 / 4900) loss: 1.509705
(Iteration 3241 / 4900) loss: 1.424874
(Iteration 3251 / 4900) loss: 1.387541
(Iteration 3261 / 4900) loss: 1.515863
(Iteration 3271 / 4900) loss: 1.464078
(Iteration 3281 / 4900) loss: 1.305884
(Iteration 3291 / 4900) loss: 1.523274
(Iteration 3301 / 4900) loss: 1.491001
(Iteration 3311 / 4900) loss: 1.409142
(Iteration 3321 / 4900) loss: 1.396137
(Iteration 3331 / 4900) loss: 1.333454
(Iteration 3341 / 4900) loss: 1.680996
(Iteration 3351 / 4900) loss: 1.651269
(Iteration 3361 / 4900) loss: 1.428943
(Iteration 3371 / 4900) loss: 1.569815
(Iteration 3381 / 4900) loss: 1.392307
(Iteration 3391 / 4900) loss: 1.559233
(Iteration 3401 / 4900) loss: 1.580548
(Iteration 3411 / 4900) loss: 1.379782
(Iteration 3421 / 4900) loss: 1.387185
(Epoch 7 / 10) train acc: 0.538000; val_acc: 0.519000
(Iteration 3431 / 4900) loss: 1.514417
(Iteration 3441 / 4900) loss: 1.512202
(Iteration 3451 / 4900) loss: 1.589685
(Iteration 3461 / 4900) loss: 1.474607
(Iteration 3471 / 4900) loss: 1.455556
(Iteration 3481 / 4900) loss: 1.573458
(Iteration 3491 / 4900) loss: 1.435063
(Iteration 3501 / 4900) loss: 1.409900
(Iteration 3511 / 4900) loss: 1.370527
(Iteration 3521 / 4900) loss: 1.306618
(Iteration 3531 / 4900) loss: 1.560820
(Iteration 3541 / 4900) loss: 1.303177
(Iteration 3551 / 4900) loss: 1.381669
(Iteration 3561 / 4900) loss: 1.415892
(Iteration 3571 / 4900) loss: 1.542831
(Iteration 3581 / 4900) loss: 1.498623
(Iteration 3591 / 4900) loss: 1.491420
(Iteration 3601 / 4900) loss: 1.314307
(Iteration 3611 / 4900) loss: 1.588824
(Iteration 3621 / 4900) loss: 1.441917
(Iteration 3631 / 4900) loss: 1.464109
```

```
(Iteration 3641 / 4900) loss: 1.554027
(Iteration 3651 / 4900) loss: 1.672120
(Iteration 3661 / 4900) loss: 1.448953
(Iteration 3671 / 4900) loss: 1.468854
(Iteration 3681 / 4900) loss: 1.409235
(Iteration 3691 / 4900) loss: 1.374531
(Iteration 3701 / 4900) loss: 1.483481
(Iteration 3711 / 4900) loss: 1.214689
(Iteration 3721 / 4900) loss: 1.440445
(Iteration 3731 / 4900) loss: 1.316762
(Iteration 3741 / 4900) loss: 1.577804
(Iteration 3751 / 4900) loss: 1.532934
(Iteration 3761 / 4900) loss: 1.677935
(Iteration 3771 / 4900) loss: 1.532323
(Iteration 3781 / 4900) loss: 1.382981
(Iteration 3791 / 4900) loss: 1.375729
(Iteration 3801 / 4900) loss: 1.357316
(Iteration 3811 / 4900) loss: 1.394309
(Iteration 3821 / 4900) loss: 1.357784
(Iteration 3831 / 4900) loss: 1.279952
(Iteration 3841 / 4900) loss: 1.475929
(Iteration 3851 / 4900) loss: 1.366373
(Iteration 3861 / 4900) loss: 1.639102
(Iteration 3871 / 4900) loss: 1.510731
(Iteration 3881 / 4900) loss: 1.590582
(Iteration 3891 / 4900) loss: 1.473270
(Iteration 3901 / 4900) loss: 1.373501
(Iteration 3911 / 4900) loss: 1.303971
(Epoch 8 / 10) train acc: 0.562000; val_acc: 0.508000
(Iteration 3921 / 4900) loss: 1.423376
(Iteration 3931 / 4900) loss: 1.468155
(Iteration 3941 / 4900) loss: 1.418921
(Iteration 3951 / 4900) loss: 1.369891
(Iteration 3961 / 4900) loss: 1.163782
(Iteration 3971 / 4900) loss: 1.555019
(Iteration 3981 / 4900) loss: 1.334379
(Iteration 3991 / 4900) loss: 1.359831
(Iteration 4001 / 4900) loss: 1.220367
(Iteration 4011 / 4900) loss: 1.378576
(Iteration 4021 / 4900) loss: 1.479768
(Iteration 4031 / 4900) loss: 1.417902
(Iteration 4041 / 4900) loss: 1.335871
(Iteration 4051 / 4900) loss: 1.300520
(Iteration 4061 / 4900) loss: 1.373932
(Iteration 4071 / 4900) loss: 1.293765
(Iteration 4081 / 4900) loss: 1.441203
(Iteration 4091 / 4900) loss: 1.357552
(Iteration 4101 / 4900) loss: 1.363083

(Iteration 4111 / 4900) loss: 1.312332
(Iteration 4121 / 4900) loss: 1.431205
(Iteration 4131 / 4900) loss: 1.395403
(Iteration 4141 / 4900) loss: 1.195140
(Iteration 4151 / 4900) loss: 1.390286
(Iteration 4161 / 4900) loss: 1.354413
(Iteration 4171 / 4900) loss: 1.515228
(Iteration 4181 / 4900) loss: 1.360029
```

```
(Iteration 4191 / 4900) loss: 1.578167
(Iteration 4201 / 4900) loss: 1.493550
(Iteration 4211 / 4900) loss: 1.312836
(Iteration 4221 / 4900) loss: 1.517592
(Iteration 4231 / 4900) loss: 1.495008
(Iteration 4241 / 4900) loss: 1.373694
(Iteration 4251 / 4900) loss: 1.450862
(Iteration 4261 / 4900) loss: 1.504509
(Iteration 4271 / 4900) loss: 1.545513
(Iteration 4281 / 4900) loss: 1.320170
(Iteration 4291 / 4900) loss: 1.499207
(Iteration 4301 / 4900) loss: 1.424187
(Iteration 4311 / 4900) loss: 1.547060
(Iteration 4321 / 4900) loss: 1.459724
(Iteration 4331 / 4900) loss: 1.397727
(Iteration 4341 / 4900) loss: 1.564625
(Iteration 4351 / 4900) loss: 1.555447
(Iteration 4361 / 4900) loss: 1.354314
(Iteration 4371 / 4900) loss: 1.440833
(Iteration 4381 / 4900) loss: 1.288322
(Iteration 4391 / 4900) loss: 1.588200
(Iteration 4401 / 4900) loss: 1.437425
(Epoch 9 / 10) train acc: 0.584000; val_acc: 0.526000
(Iteration 4411 / 4900) loss: 1.553572
(Iteration 4421 / 4900) loss: 1.221204
(Iteration 4431 / 4900) loss: 1.495060
(Iteration 4441 / 4900) loss: 1.268140
(Iteration 4451 / 4900) loss: 1.335368
(Iteration 4461 / 4900) loss: 1.550898
(Iteration 4471 / 4900) loss: 1.563515
(Iteration 4481 / 4900) loss: 1.304514
(Iteration 4491 / 4900) loss: 1.385563
(Iteration 4501 / 4900) loss: 1.272674
(Iteration 4511 / 4900) loss: 1.359469
(Iteration 4521 / 4900) loss: 1.281859
(Iteration 4531 / 4900) loss: 1.371623
(Iteration 4541 / 4900) loss: 1.285450
(Iteration 4551 / 4900) loss: 1.486921
(Iteration 4561 / 4900) loss: 1.489606
(Iteration 4571 / 4900) loss: 1.264356
(Iteration 4581 / 4900) loss: 1.426886
(Iteration 4591 / 4900) loss: 1.469634
(Iteration 4601 / 4900) loss: 1.353229
(Iteration 4611 / 4900) loss: 1.483000
(Iteration 4621 / 4900) loss: 1.529013
(Iteration 4631 / 4900) loss: 1.393742
(Iteration 4641 / 4900) loss: 1.421267
(Iteration 4651 / 4900) loss: 1.269265
(Iteration 4661 / 4900) loss: 1.152884
(Iteration 4671 / 4900) loss: 1.358464
(Iteration 4681 / 4900) loss: 1.389666
(Iteration 4691 / 4900) loss: 1.225836
(Iteration 4701 / 4900) loss: 1.546010
(Iteration 4711 / 4900) loss: 1.394259
(Iteration 4721 / 4900) loss: 1.377669
(Iteration 4731 / 4900) loss: 1.405093
(Iteration 4741 / 4900) loss: 1.507539
```



```
(Iteration 4751 / 4900) loss: 1.230656
(Iteration 4761 / 4900) loss: 1.423131
(Iteration 4771 / 4900) loss: 1.323932
(Iteration 4781 / 4900) loss: 1.358263
(Iteration 4791 / 4900) loss: 1.662666
(Iteration 4801 / 4900) loss: 1.505661
(Iteration 4811 / 4900) loss: 1.469058
(Iteration 4821 / 4900) loss: 1.506454
(Iteration 4831 / 4900) loss: 1.325260
(Iteration 4841 / 4900) loss: 1.489108
(Iteration 4851 / 4900) loss: 1.230440
(Iteration 4861 / 4900) loss: 1.300720
(Iteration 4871 / 4900) loss: 1.334097
(Iteration 4881 / 4900) loss: 1.138237
(Iteration 4891 / 4900) loss: 1.214555
(Epoch 10 / 10) train acc: 0.548000; val_acc: 0.518000
```

## Test your model!

Run your best model on the validation and test sets. You should achieve above 50% accuracy on the validation set.

```
In [21]: 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: 0.539
```

```
Test set accuracy: 0.514
```