

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.

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: 9.769847728806635e-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:
difference: 1.0000000000000001e-10
```

```
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:

1, 2

1. When the input value is negative, it will get zero gradient flow.
2. When the input value is really small or large, it will get close to zero gradient flow

"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)[0], x, dout)
dw_num = eval_numerical_gradient_array(lambda w: affine_relu_forward(x, w, b)[0], w, dout)
db_num = eval_numerical_gradient_array(lambda b: affine_relu_forward(x, w, b)[0], 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:  6.750562121603446e-11
dw error:  8.162015570444288e-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 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
```

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, 15.49994135, 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)
```

```

Testing training loss (no regularization)
Running numeric gradient check with reg = 0.0
W1 relative error: 1.22e-08
W2 relative error: 3.48e-10
b1 relative error: 6.55e-09
b2 relative error: 4.33e-10
Running numeric gradient check with reg = 0.7
W1 relative error: 8.18e-07
W2 relative error: 2.85e-08
b1 relative error: 1.09e-09
b2 relative error: 7.76e-10

```

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

model = TwoLayerNet(reg = 0.5)

solver = Solver(model, data,
                update_rule='sgd',
                optim_config={
                    'learning_rate': 1e-3,
                },
                lr_decay=0.95,
                num_epochs=10, batch_size=100,
                print_every=100)

solver.train()

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

(Iteration 1 / 4900) loss: 2.378087
(Epoch 0 / 10) train acc: 0.164000; val_acc: 0.134000
(Iteration 101 / 4900) loss: 1.905230
(Iteration 201 / 4900) loss: 2.049261
(Iteration 301 / 4900) loss: 1.713070
(Iteration 401 / 4900) loss: 1.582693
(Epoch 1 / 10) train acc: 0.445000; val_acc: 0.451000
(Iteration 501 / 4900) loss: 1.673854
(Iteration 601 / 4900) loss: 1.542241
(Iteration 701 / 4900) loss: 1.661488
(Iteration 801 / 4900) loss: 1.698760
(Iteration 901 / 4900) loss: 1.523104
(Epoch 2 / 10) train acc: 0.486000; val_acc: 0.470000
(Iteration 1001 / 4900) loss: 1.575808
(Iteration 1101 / 4900) loss: 1.559150
(Iteration 1201 / 4900) loss: 1.512748
(Iteration 1301 / 4900) loss: 1.403581
(Iteration 1401 / 4900) loss: 1.582180
(Epoch 3 / 10) train acc: 0.509000; val_acc: 0.478000
(Iteration 1501 / 4900) loss: 1.520756
(Iteration 1601 / 4900) loss: 1.502950
(Iteration 1701 / 4900) loss: 1.504738
(Iteration 1801 / 4900) loss: 1.639809
(Iteration 1901 / 4900) loss: 1.588290
(Epoch 4 / 10) train acc: 0.494000; val_acc: 0.479000
(Iteration 2001 / 4900) loss: 1.617043

```

```

(Iteration 2101 / 4900) loss: 1.571903
(Iteration 2201 / 4900) loss: 1.580801
(Iteration 2301 / 4900) loss: 1.319308
(Iteration 2401 / 4900) loss: 1.405106
(Epoch 5 / 10) train acc: 0.523000; val_acc: 0.487000
(Iteration 2501 / 4900) loss: 1.462077
(Iteration 2601 / 4900) loss: 1.478338
(Iteration 2701 / 4900) loss: 1.452321
(Iteration 2801 / 4900) loss: 1.501195
(Iteration 2901 / 4900) loss: 1.444775
(Epoch 6 / 10) train acc: 0.553000; val_acc: 0.507000
(Iteration 3001 / 4900) loss: 1.377007
(Iteration 3101 / 4900) loss: 1.252317
(Iteration 3201 / 4900) loss: 1.703810
(Iteration 3301 / 4900) loss: 1.449870
(Iteration 3401 / 4900) loss: 1.579887
(Epoch 7 / 10) train acc: 0.532000; val_acc: 0.490000
(Iteration 3501 / 4900) loss: 1.427413
(Iteration 3601 / 4900) loss: 1.284990
(Iteration 3701 / 4900) loss: 1.474521
(Iteration 3801 / 4900) loss: 1.398789
(Iteration 3901 / 4900) loss: 1.239221
(Epoch 8 / 10) train acc: 0.517000; val_acc: 0.492000
(Iteration 4001 / 4900) loss: 1.355188
(Iteration 4101 / 4900) loss: 1.393246
(Iteration 4201 / 4900) loss: 1.277518
(Iteration 4301 / 4900) loss: 1.229724
(Iteration 4401 / 4900) loss: 1.707493
(Epoch 9 / 10) train acc: 0.570000; val_acc: 0.514000
(Iteration 4501 / 4900) loss: 1.374712
(Iteration 4601 / 4900) loss: 1.592834
(Iteration 4701 / 4900) loss: 1.579716
(Iteration 4801 / 4900) loss: 1.336092
(Epoch 10 / 10) train acc: 0.566000; val_acc: 0.477000

```

In [11]:

```
# Run this cell to visualize training loss and train / val accuracy
```

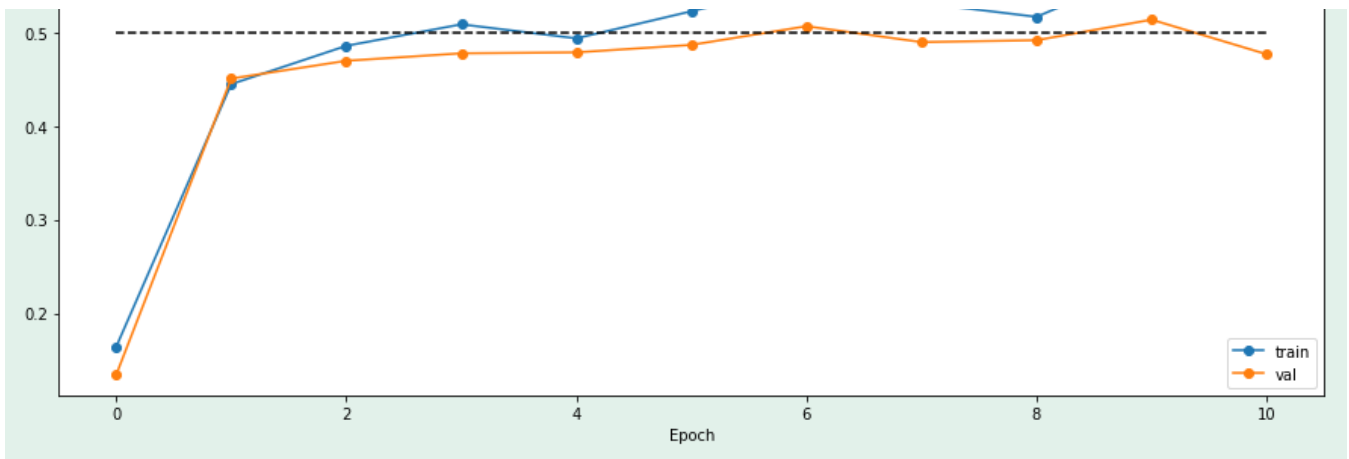
```

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, 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: 1.14e-08
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
```

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 **weight initialization scale** to overfit and achieve 100% training accuracy within 20 epochs.

In [13]:

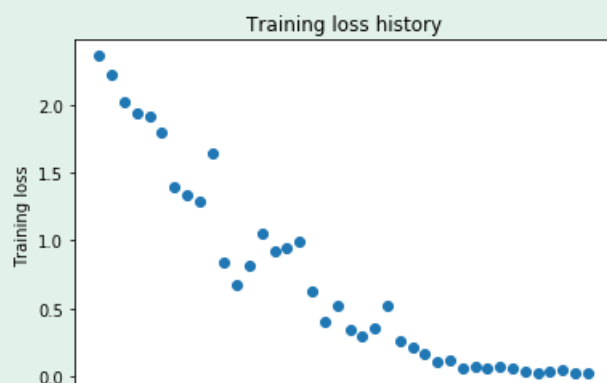
```
# TODO: Use a three-layer Net to overfit 50 training examples by
# tweaking just the learning rate and initialization scale.

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-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
```



0 5 10 15 20 25 30 35 40
Iteration

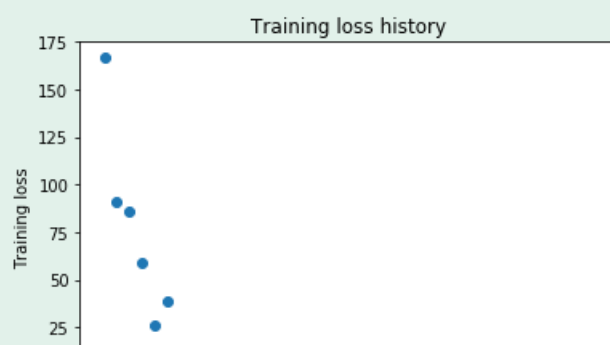
Now try to use a five-layer network with 100 units on each layer to overfit 50 training examples. Again, you will have to adjust the learning rate and weight initialization scale, but you should be able to achieve 100% training accuracy within 20 epochs.

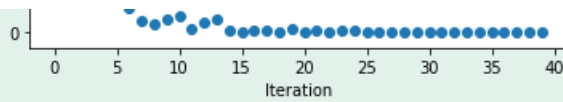
In [14]:

```
# 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 = 3e-4 # Experiment with this!  
weight_scale = 1e-1 # 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: 0.160000; val_acc: 0.120000  
(Epoch 1 / 20) train acc: 0.240000; val_acc: 0.095000  
(Epoch 2 / 20) train acc: 0.320000; val_acc: 0.123000  
(Epoch 3 / 20) train acc: 0.420000; val_acc: 0.120000  
(Epoch 4 / 20) train acc: 0.640000; val_acc: 0.128000  
(Epoch 5 / 20) train acc: 0.680000; val_acc: 0.112000  
(Iteration 11 / 40) loss: 8.208668  
(Epoch 6 / 20) train acc: 0.820000; val_acc: 0.128000  
(Epoch 7 / 20) train acc: 0.860000; val_acc: 0.118000  
(Epoch 8 / 20) train acc: 0.920000; val_acc: 0.118000  
(Epoch 9 / 20) train acc: 0.960000; val_acc: 0.122000  
(Epoch 10 / 20) train acc: 0.920000; val_acc: 0.120000  
(Iteration 21 / 40) loss: 0.206859  
(Epoch 11 / 20) train acc: 0.940000; val_acc: 0.111000  
(Epoch 12 / 20) train acc: 0.960000; val_acc: 0.122000  
(Epoch 13 / 20) train acc: 0.980000; val_acc: 0.128000  
(Epoch 14 / 20) train acc: 0.960000; val_acc: 0.118000  
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.122000  
(Iteration 31 / 40) loss: 0.000895  
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.119000  
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.120000  
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.120000  
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.120000  
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.120000
```





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:

The five-layer net is harder to train and to find correct hyperparameters (learning rate and weight scale). And five-layer network is more sensitive to the initialization hyperparameter scale. Because five-layer net converges much faster and produces more precise result. It will be more sensitive to initial set up.

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> 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 [17]:

```
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()

```

```

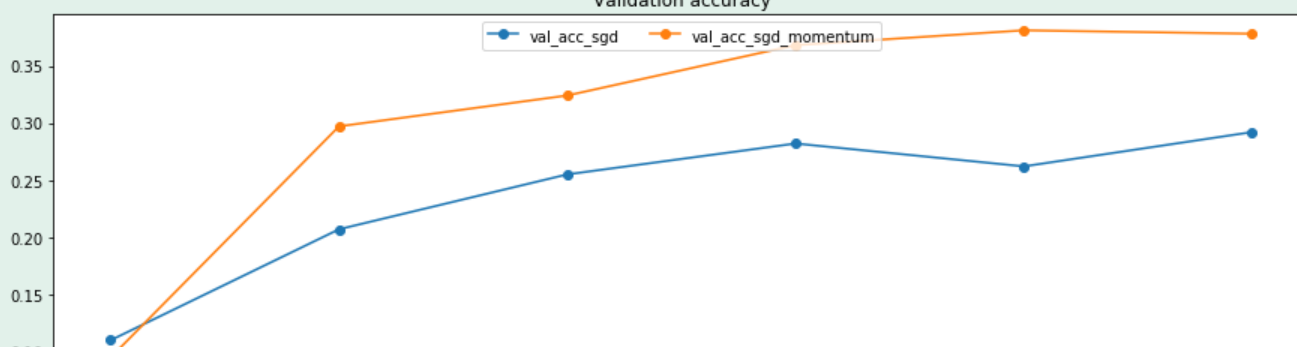
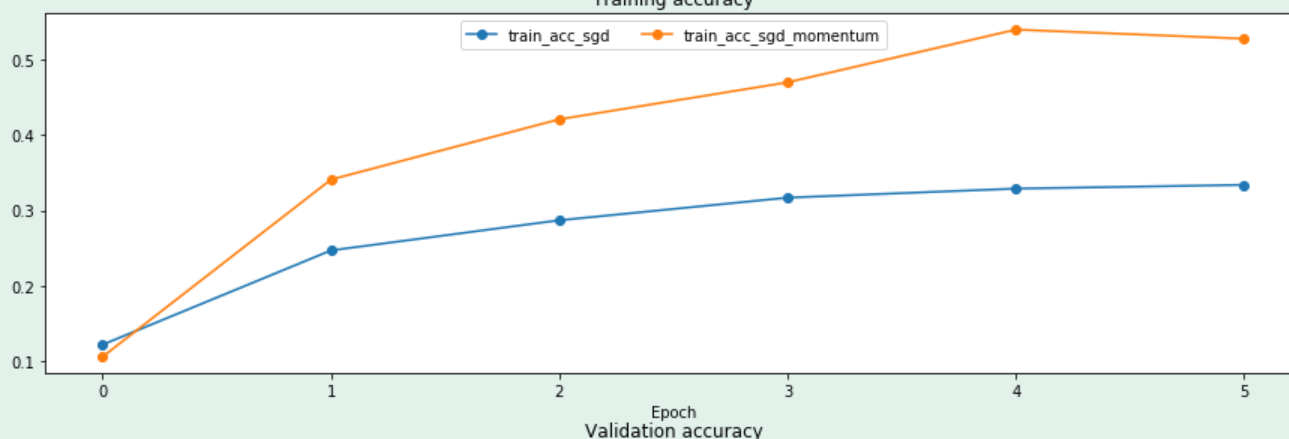
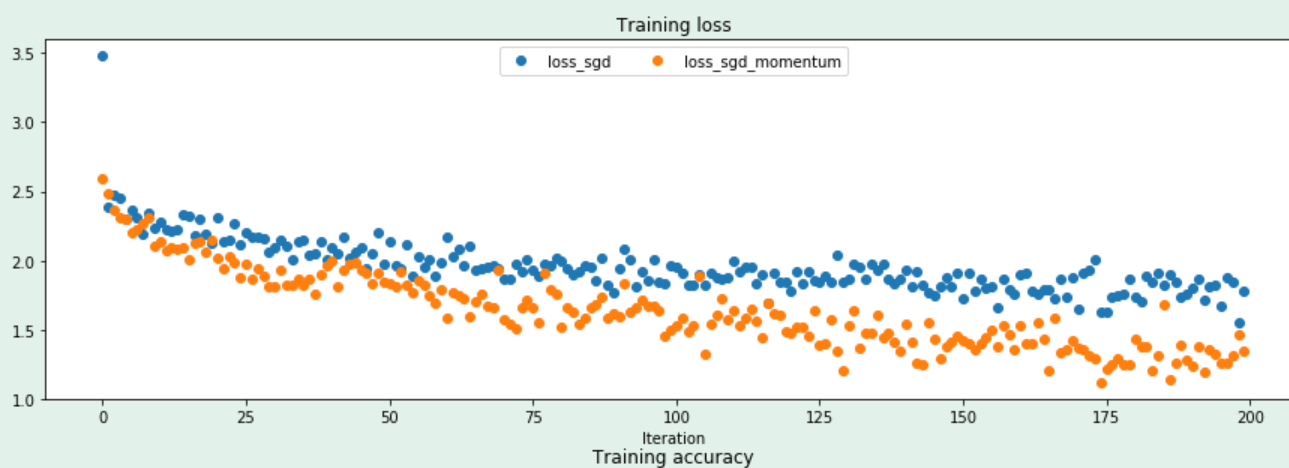
running with sgd
(Iteration 1 / 200) loss: 3.476928
(Epoch 0 / 5) train acc: 0.122000; val_acc: 0.110000
(Iteration 11 / 200) loss: 2.273918
(Iteration 21 / 200) loss: 2.306955
(Iteration 31 / 200) loss: 2.093337
(Epoch 1 / 5) train acc: 0.247000; val_acc: 0.207000
(Iteration 41 / 200) loss: 2.089881
(Iteration 51 / 200) loss: 2.142988
(Iteration 61 / 200) loss: 2.169558
(Iteration 71 / 200) loss: 1.870478
(Epoch 2 / 5) train acc: 0.287000; val_acc: 0.255000
(Iteration 81 / 200) loss: 1.994859
(Iteration 91 / 200) loss: 1.948483
(Iteration 101 / 200) loss: 1.956554
(Iteration 111 / 200) loss: 1.997367
(Epoch 3 / 5) train acc: 0.317000; val_acc: 0.282000
(Iteration 121 / 200) loss: 1.786599
(Iteration 131 / 200) loss: 1.863651
(Iteration 141 / 200) loss: 1.936477
(Iteration 151 / 200) loss: 1.732275
(Epoch 4 / 5) train acc: 0.329000; val_acc: 0.262000
(Iteration 161 / 200) loss: 1.903142
(Iteration 171 / 200) loss: 1.649538
(Iteration 181 / 200) loss: 1.742841

```

```
(Iteration 191 / 200) loss: 1.806653
(Epoch 5 / 5) train acc: 0.334000; val_acc: 0.292000
```

running with `sgd_momentum`

```
(Iteration 1 / 200) loss: 2.589166
(Epoch 0 / 5) train acc: 0.106000; val_acc: 0.095000
(Iteration 11 / 200) loss: 2.139318
(Iteration 21 / 200) loss: 2.017108
(Iteration 31 / 200) loss: 1.813453
(Epoch 1 / 5) train acc: 0.341000; val_acc: 0.297000
(Iteration 41 / 200) loss: 2.002500
(Iteration 51 / 200) loss: 1.830173
(Iteration 61 / 200) loss: 1.584339
(Iteration 71 / 200) loss: 1.579533
(Epoch 2 / 5) train acc: 0.421000; val_acc: 0.324000
(Iteration 81 / 200) loss: 1.523871
(Iteration 91 / 200) loss: 1.597867
(Iteration 101 / 200) loss: 1.527771
(Iteration 111 / 200) loss: 1.638278
(Epoch 3 / 5) train acc: 0.470000; val_acc: 0.368000
(Iteration 121 / 200) loss: 1.482900
(Iteration 131 / 200) loss: 1.533665
(Iteration 141 / 200) loss: 1.547938
(Iteration 151 / 200) loss: 1.419212
(Epoch 4 / 5) train acc: 0.540000; val_acc: 0.381000
(Iteration 161 / 200) loss: 1.531513
(Iteration 171 / 200) loss: 1.371091
(Iteration 181 / 200) loss: 1.438618
(Iteration 191 / 200) loss: 1.241694
(Epoch 5 / 5) train acc: 0.528000; val_acc: 0.378000
```





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." COURSERA: Neural Networks for Machine Learning 4 (2012).

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

In [18]:

```
# 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: 9.524687511038133e-08
cache error: 2.6477955807156126e-09
```

In [19]:

```
# 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],
    [0.54036842, 0.52993714, 0.51950586, 0.50907458, 0.4986433]
```

```
[ 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: 1.1395691798535431e-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 [20]:

```
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: 2.760209
(Epoch 0 / 5) train acc: 0.150000; val_acc: 0.124000
(Iteration 11 / 200) loss: 2.110602
(Iteration 21 / 200) loss: 1.883176
(Iteration 31 / 200) loss: 1.842409
(Epoch 1 / 5) train acc: 0.357000; val_acc: 0.312000
(Iteration 41 / 200) loss: 1.912024
(Iteration 51 / 200) loss: 1.622724
(Iteration 61 / 200) loss: 1.598922
(Iteration 71 / 200) loss: 1.647785
(Epoch 2 / 5) train acc: 0.441000; val_acc: 0.360000
```

```

(Iteration 81 / 200) loss: 1.517582
(Iteration 91 / 200) loss: 1.621971
(Iteration 101 / 200) loss: 1.456345
(Iteration 111 / 200) loss: 1.491637
(Epoch 3 / 5) train acc: 0.489000; val_acc: 0.367000
(Iteration 121 / 200) loss: 1.607893
(Iteration 131 / 200) loss: 1.394200
(Iteration 141 / 200) loss: 1.321148
(Iteration 151 / 200) loss: 1.349208
(Epoch 4 / 5) train acc: 0.545000; val_acc: 0.381000
(Iteration 161 / 200) loss: 1.248493
(Iteration 171 / 200) loss: 1.153209
(Iteration 181 / 200) loss: 1.493787
(Iteration 191 / 200) loss: 1.214109
(Epoch 5 / 5) train acc: 0.593000; val_acc: 0.393000

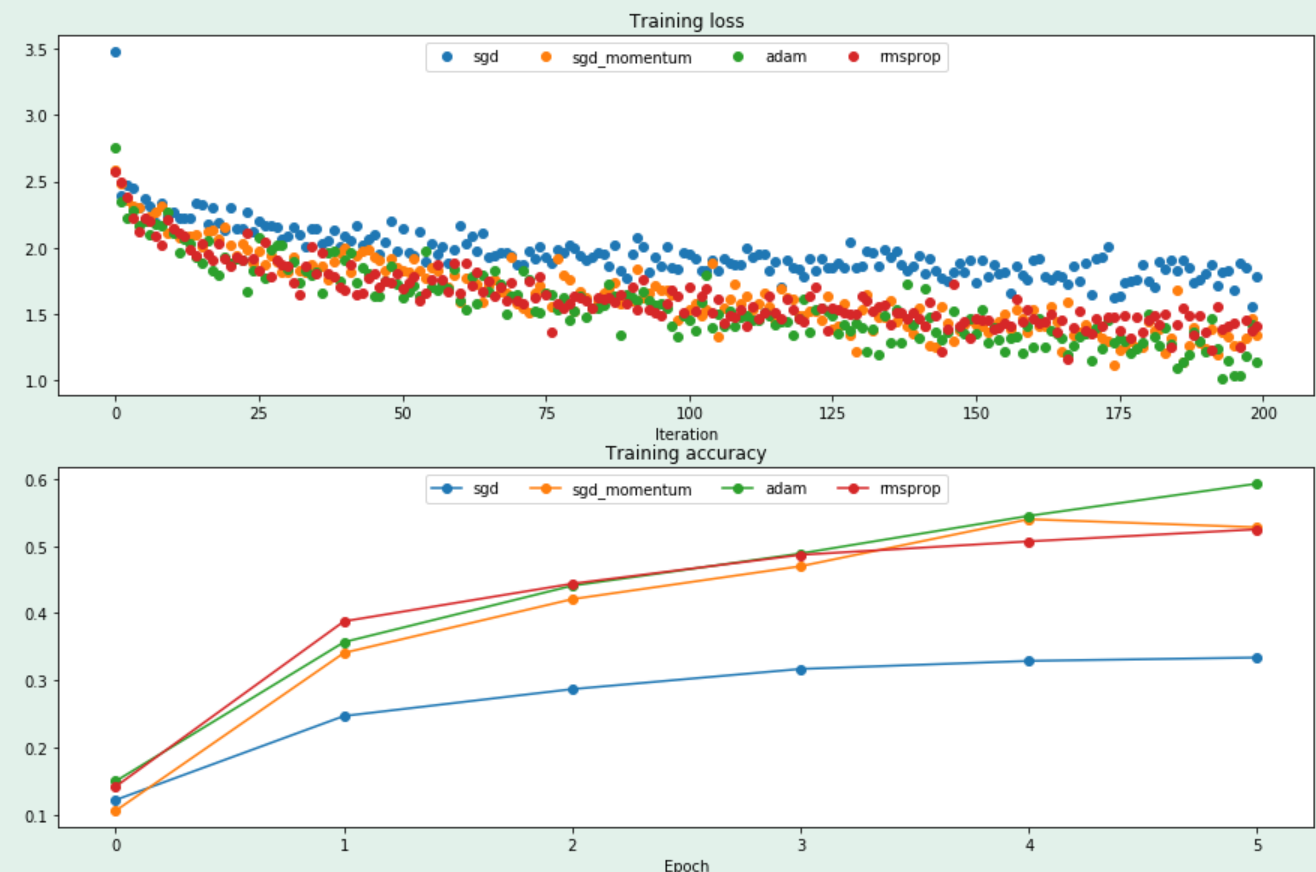
```

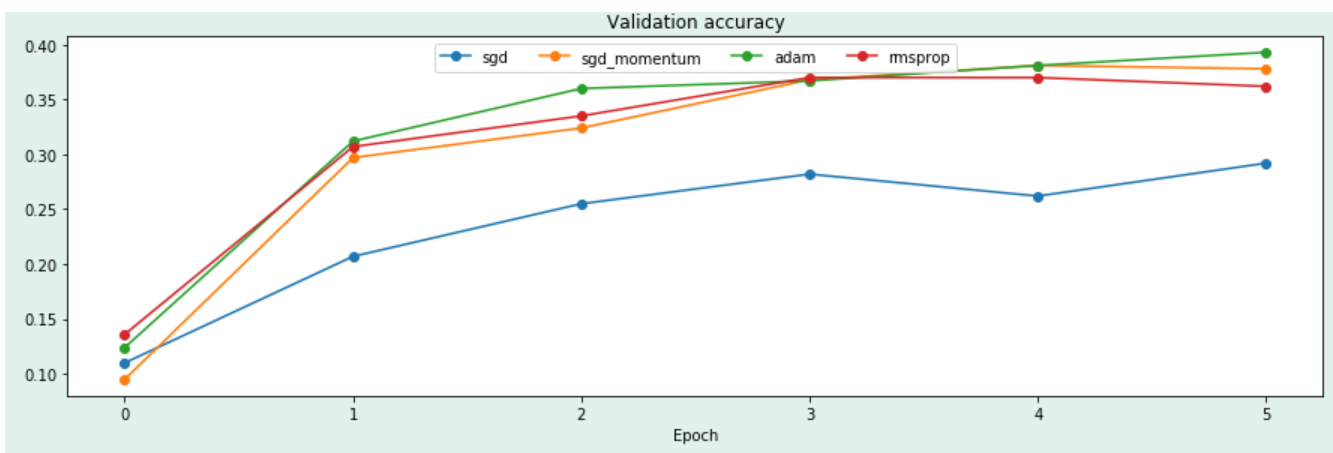
running with rmsprop

```

(Iteration 1 / 200) loss: 2.571494
(Epoch 0 / 5) train acc: 0.142000; val_acc: 0.136000
(Iteration 11 / 200) loss: 2.148984
(Iteration 21 / 200) loss: 1.858964
(Iteration 31 / 200) loss: 1.865406
(Epoch 1 / 5) train acc: 0.388000; val_acc: 0.307000
(Iteration 41 / 200) loss: 1.682452
(Iteration 51 / 200) loss: 1.690672
(Iteration 61 / 200) loss: 1.653560
(Iteration 71 / 200) loss: 1.602213
(Epoch 2 / 5) train acc: 0.444000; val_acc: 0.335000
(Iteration 81 / 200) loss: 1.633064
(Iteration 91 / 200) loss: 1.700801
(Iteration 101 / 200) loss: 1.698334
(Iteration 111 / 200) loss: 1.405926
(Epoch 3 / 5) train acc: 0.487000; val_acc: 0.370000
(Iteration 121 / 200) loss: 1.445123
(Iteration 131 / 200) loss: 1.634678
(Iteration 141 / 200) loss: 1.460120
(Iteration 151 / 200) loss: 1.462433
(Epoch 4 / 5) train acc: 0.507000; val_acc: 0.370000
(Iteration 161 / 200) loss: 1.462897
(Iteration 171 / 200) loss: 1.357323
(Iteration 181 / 200) loss: 1.451439
(Iteration 191 / 200) loss: 1.361932
(Epoch 5 / 5) train acc: 0.525000; val_acc: 0.362000

```





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:

Because every time the square of dw is added to cache. As the learning progresses, the value of cache will become larger and larger, and in the process of w update, it needs to be divided by cache, which will cause the actual learning rate to be smaller and smaller.

Adam would not have this kind of issue. Because in Adam, m and v are similar to cache but they are moving averages of squared gradients. Hyperparameters β_1 and β_2 will respectively make m and v to be leaky so the learning rate updates do not get monotonically smaller.

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 [23]:

```
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_scales = [1e-2, 1e-3]
learning_rates = [1e-5, 1e-4, 1e-3]
regulation_strengths = [1e-4, 1e-3]

best_val_accuracy = 0

for weight_scale in weight_scales:
    for lr in learning_rates:
```



```

for reg in regulation_strengths:
    model = FullyConnectedNet([200, 100],
                               weight_scale=weight_scale,
                               reg=reg)

    solver = Solver(model,
                    data,
                    num_epochs=10,
                    batch_size=200,
                    update_rule='adam',
                    optim_config={
                        'learning_rate': lr
                    },
                    verbose=False)

    solver.train()

    val_accuracy = solver.best_val_acc

    if best_val_accuracy < val_accuracy:
        best_val_accuracy = val_accuracy
        best_model = model

    print('learning rate %e weight scale %e regulation strength %e val accuracy: %f' % (lr,
weight_scale, reg, val_accuracy))

print('Best validation accuracy achieved is: %f' % best_val_accuracy)

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

```

learning rate 1.000000e-05 weight scale 1.000000e-02 regulation strength 1.000000e-04 val accuracy: 0.483000
learning rate 1.000000e-05 weight scale 1.000000e-02 regulation strength 1.000000e-03 val accuracy: 0.468000
learning rate 1.000000e-04 weight scale 1.000000e-02 regulation strength 1.000000e-04 val accuracy: 0.542000
learning rate 1.000000e-04 weight scale 1.000000e-02 regulation strength 1.000000e-03 val accuracy: 0.531000
learning rate 1.000000e-03 weight scale 1.000000e-02 regulation strength 1.000000e-04 val accuracy: 0.525000
learning rate 1.000000e-03 weight scale 1.000000e-02 regulation strength 1.000000e-03 val accuracy: 0.509000
learning rate 1.000000e-05 weight scale 1.000000e-03 regulation strength 1.000000e-04 val accuracy: 0.447000
learning rate 1.000000e-05 weight scale 1.000000e-03 regulation strength 1.000000e-03 val accuracy: 0.445000
learning rate 1.000000e-04 weight scale 1.000000e-03 regulation strength 1.000000e-04 val accuracy: 0.531000
learning rate 1.000000e-04 weight scale 1.000000e-03 regulation strength 1.000000e-03 val accuracy: 0.522000
learning rate 1.000000e-03 weight scale 1.000000e-03 regulation strength 1.000000e-04 val accuracy: 0.509000
learning rate 1.000000e-03 weight scale 1.000000e-03 regulation strength 1.000000e-03 val accuracy: 0.505000
Best validation accuracy achieved is: 0.542000

Test your model!

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

In [25]:

```

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: 0.542
Test set accuracy: 0.534

Batch Normalization

One way to make deep networks easier to train is to use more sophisticated optimization procedures such as SGD+momentum, RMSProp, or Adam. Another strategy is to change the architecture of the network to make it easier to train. One idea along these lines is batch normalization which was proposed by [1] in 2015.

The idea is relatively straightforward. Machine learning methods tend to work better when their input data consists of uncorrelated features with zero mean and unit variance. When training a neural network, we can preprocess the data before feeding it to the network to explicitly decorrelate its features; this will ensure that the first layer of the network sees data that follows a nice distribution. However, even if we preprocess the input data, the activations at deeper layers of the network will likely no longer be decorrelated and will no longer have zero mean or unit variance since they are output from earlier layers in the network. Even worse, during the training process the distribution of features at each layer of the network will shift as the weights of each layer are updated.

The authors of [1] hypothesize that the shifting distribution of features inside deep neural networks may make training deep networks more difficult. To overcome this problem, [1] proposes to insert batch normalization layers into the network. At training time, a batch normalization layer uses a minibatch of data to estimate the mean and standard deviation of each feature. These estimated means and standard deviations are then used to center and normalize the features of the minibatch. A running average of these means and standard deviations is kept during training, and at test time these running averages are used to center and normalize features.

It is possible that this normalization strategy could reduce the representational power of the network, since it may sometimes be optimal for certain layers to have features that are not zero-mean or unit variance. To this end, the batch normalization layer includes learnable shift and scale parameters for each feature dimension.

[1] [Sergey Ioffe and Christian Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift", ICML 2015.](#)

Batch normalization: forward

In the file `cs231n/layers.py`, implement the batch normalization forward pass in the function `batchnorm_forward`. Once you have done so, run the following to test your implementation.

Referencing the paper linked to above in [1] may be helpful!

In [3]:

```
# Check the training-time forward pass by checking means and variances
# of features both before and after batch normalization

# Simulate the forward pass for a two-layer network
np.random.seed(231)
N, D1, D2, D3 = 200, 50, 60, 3
X = np.random.randn(N, D1)
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)
a = np.maximum(0, X.dot(W1)).dot(W2)

print('Before batch normalization:')
print_mean_std(a,axis=0)

gamma = np.ones((D3,))
beta = np.zeros((D3,))
# Means should be close to zero and stds close to one
print('After batch normalization (gamma=1, beta=0)')
a_norm, _ = batchnorm_forward(a, gamma, beta, {'mode': 'train'})
print_mean_std(a_norm,axis=0)

gamma = np.asarray([1.0, 2.0, 3.0])
beta = np.asarray([11.0, 12.0, 13.0])
# Now means should be close to beta and stds close to gamma
print('After batch normalization (gamma=', gamma, ', beta=', beta, ')')
a_norm, _ = batchnorm_forward(a, gamma, beta, {'mode': 'train'})
print_mean_std(a_norm,axis=0)
```

```
Before batch normalization:
means:  [ -2.3814598 -13.18038246  1.91780462]
```

```
stds: [27.18502186 34.21455511 37.68611762]

After batch normalization (gamma=1, beta=0)
means: [3.99680289e-17 6.93889390e-17 4.41313652e-17]
stds: [0.99999999 1. 1. ]

After batch normalization (gamma= [1. 2. 3.] , beta= [11. 12. 13.] )
means: [11. 12. 13.]
stds: [0.99999999 1.99999999 2.99999999]
```

In [4]:

```
# Check the test-time forward pass by running the training-time
# forward pass many times to warm up the running averages, and then
# checking the means and variances of activations after a test-time
# forward pass.

np.random.seed(231)
N, D1, D2, D3 = 200, 50, 60, 3
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)

bn_param = {'mode': 'train'}
gamma = np.ones(D3)
beta = np.zeros(D3)

for t in range(50):
    X = np.random.randn(N, D1)
    a = np.maximum(0, X.dot(W1)).dot(W2)
    batchnorm_forward(a, gamma, beta, bn_param)

bn_param['mode'] = 'test'
X = np.random.randn(N, D1)
a = np.maximum(0, X.dot(W1)).dot(W2)
a_norm, _ = batchnorm_forward(a, gamma, beta, bn_param)

# Means should be close to zero and stds close to one, but will be
# noisier than training-time forward passes.
print('After batch normalization (test-time):')
print_mean_std(a_norm, axis=0)

After batch normalization (test-time):
means: [-0.03927354 -0.04349152 -0.10452688]
stds: [1.01531428 1.01238373 0.97819988]
```

Batch normalization: backward

Now implement the backward pass for batch normalization in the function `batchnorm_backward`.

To derive the backward pass you should write out the computation graph for batch normalization and backprop through each of the intermediate nodes. Some intermediates may have multiple outgoing branches; make sure to sum gradients across these branches in the backward pass.

Once you have finished, run the following to numerically check your backward pass.

In [5]:

```
# Gradient check batchnorm backward pass
np.random.seed(231)
N, D = 4, 5
x = 5 * np.random.randn(N, D) + 12
gamma = np.random.randn(D)
beta = np.random.randn(D)
dout = np.random.randn(N, D)

bn_param = {'mode': 'train'}
fx = lambda x: batchnorm_forward(x, gamma, beta, bn_param)[0]
fg = lambda a: batchnorm_forward(x, a, beta, bn_param)[0]
fb = lambda b: batchnorm_forward(x, gamma, b, bn_param)[0]

dx_num = eval_numerical_gradient_array(fx, x, dout)
da_num = eval_numerical_gradient_array(fg, gamma.copy(), dout)
db_num = eval_numerical_gradient_array(fb, beta.copy(), dout)
```

```
_, cache = batchnorm_forward(x, gamma, beta, bn_param)
dx, dgamma, dbeta = batchnorm_backward(dout, cache)
#You should expect to see relative errors between 1e-13 and 1e-8
print('dx error: ', rel_error(dx_num, dx))
print('dgamma error: ', rel_error(da_num, dgamma))
print('dbeta error: ', rel_error(db_num, dbeta))

dx error:  1.7029261167605239e-09
dgamma error:  7.420414216247087e-13
dbeta error:  2.8795057655839487e-12
```

Batch normalization: alternative backward

In class we talked about two different implementations for the sigmoid backward pass. One strategy is to write out a computation graph composed of simple operations and backprop through all intermediate values. Another strategy is to work out the derivatives on paper. For example, you can derive a very simple formula for the sigmoid function's backward pass by simplifying gradients on paper.

Surprisingly, it turns out that you can do a similar simplification for the batch normalization backward pass too!

In the forward pass, given a set of inputs $X = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_N \end{bmatrix}$,

we first calculate the mean μ and variance v . With μ and v calculated, we can calculate the standard deviation σ and normalized data Y . The equations and graph illustration below describe the computation (y_i is the i -th element of the vector Y).

$$\mu = \frac{1}{N} \sum_{k=1}^N x_k \quad v = \frac{1}{N} \sum_{k=1}^N (x_k - \mu)^2$$

$$\sigma = \sqrt{v + \epsilon} \quad y_i = \frac{x_i - \mu}{\sigma}$$

In [6]:

```
np.random.seed(231)
N, D = 100, 500
x = 5 * np.random.randn(N, D) + 12
gamma = np.random.randn(D)
beta = np.random.randn(D)
dout = np.random.randn(N, D)

bn_param = {'mode': 'train'}
out, cache = batchnorm_forward(x, gamma, beta, bn_param)

t1 = time.time()
dx1, dgamma1, dbeta1 = batchnorm_backward(dout, cache)
t2 = time.time()
dx2, dgamma2, dbeta2 = batchnorm_backward_alt(dout, cache)
t3 = time.time()

print('dx difference: ', rel_error(dx1, dx2))
print('dgamma difference: ', rel_error(dgamma1, dgamma2))
print('dbeta difference: ', rel_error(dbeta1, dbeta2))
print('speedup: %.2fx' % ((t2 - t1) / (t3 - t2)))

dx difference:  1.917873695737547e-12
dgamma difference:  0.0
dbeta difference:  0.0
speedup: 1.98x
```

Fully Connected Nets with Batch Normalization

Now that you have a working implementation for batch normalization, go back to your `FullyConnectedNet` in the file

`cs231n/classifiers/fully_net.py`. Modify your implementation to add batch normalization

`cs231n/classifiers/fc_net.py` . Modify your implementation to add batch normalization.

Concretely, when the `normalization` flag is set to `"batchnorm"` in the constructor, you should insert a batch normalization layer before each ReLU nonlinearity. The outputs from the last layer of the network should not be normalized. Once you are done, run the following to gradient-check your implementation.

HINT: You might find it useful to define an additional helper layer similar to those in the file `cs231n/layer_utils.py` . If you decide to do so, do it in the file `cs231n/classifiers/fc_net.py` .

In [7]:

```
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,))

# You should expect losses between 1e-4~1e-10 for W,
# losses between 1e-08~1e-10 for b,
# and losses between 1e-08~1e-09 for beta and gammas.
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,
                              normalization='batchnorm')

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

    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])))
    if reg == 0: print()
```

```
Running check with reg = 0
Initial loss: 2.2611955101340957
W1 relative error: 1.10e-04
W2 relative error: 2.85e-06
W3 relative error: 3.92e-10
b1 relative error: 4.44e-08
b2 relative error: 2.22e-08
b3 relative error: 4.78e-11
beta1 relative error: 7.33e-09
beta2 relative error: 1.89e-09
gamma1 relative error: 7.57e-09
gamma2 relative error: 1.96e-09
```

```
Running check with reg = 3.14
Initial loss: 6.996533220108303
W1 relative error: 1.98e-06
W2 relative error: 2.29e-06
W3 relative error: 1.11e-08
b1 relative error: 5.55e-09
b2 relative error: 5.55e-09
b3 relative error: 2.23e-10
beta1 relative error: 6.65e-09
beta2 relative error: 3.48e-09
gamma1 relative error: 5.94e-09
gamma2 relative error: 4.14e-09
```

Batchnorm for deep networks

Run the following to train a six-layer network on a subset of 1000 training examples both with and without batch normalization.

In [8]:

```
np.random.seed(231)
# Try training a very deep net with batchnorm
hidden_dims = [100, 100, 100, 100, 100]

num_train = 1000
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']
}
```

```

    'x_val': data['x_val'],
    'y_val': data['y_val'],
}

weight_scale = 2e-2
bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, normalization='batchnorm')
model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, normalization=None)

print('Solver with batch norm:')
bn_solver = Solver(bn_model, small_data,
                   num_epochs=10, batch_size=50,
                   update_rule='adam',
                   optim_config={
                       'learning_rate': 1e-3,
                   },
                   verbose=True, print_every=20)
bn_solver.train()

print('\nSolver without batch norm:')
solver = Solver(model, small_data,
                num_epochs=10, batch_size=50,
                update_rule='adam',
                optim_config={
                    'learning_rate': 1e-3,
                },
                verbose=True, print_every=20)
solver.train()

```

```

Solver with batch norm:
(Iteration 1 / 200) loss: 2.340974
(Epoch 0 / 10) train acc: 0.107000; val_acc: 0.115000
(Epoch 1 / 10) train acc: 0.313000; val_acc: 0.265000
(Iteration 21 / 200) loss: 2.039345
(Epoch 2 / 10) train acc: 0.395000; val_acc: 0.278000
(Iteration 41 / 200) loss: 2.047471
(Epoch 3 / 10) train acc: 0.483000; val_acc: 0.317000
(Iteration 61 / 200) loss: 1.739554
(Epoch 4 / 10) train acc: 0.524000; val_acc: 0.318000
(Iteration 81 / 200) loss: 1.246973
(Epoch 5 / 10) train acc: 0.590000; val_acc: 0.336000
(Iteration 101 / 200) loss: 1.352696
(Epoch 6 / 10) train acc: 0.640000; val_acc: 0.322000
(Iteration 121 / 200) loss: 1.012431
(Epoch 7 / 10) train acc: 0.667000; val_acc: 0.336000
(Iteration 141 / 200) loss: 1.178837
(Epoch 8 / 10) train acc: 0.693000; val_acc: 0.321000
(Iteration 161 / 200) loss: 0.762896
(Epoch 9 / 10) train acc: 0.776000; val_acc: 0.348000
(Iteration 181 / 200) loss: 0.864004
(Epoch 10 / 10) train acc: 0.774000; val_acc: 0.308000

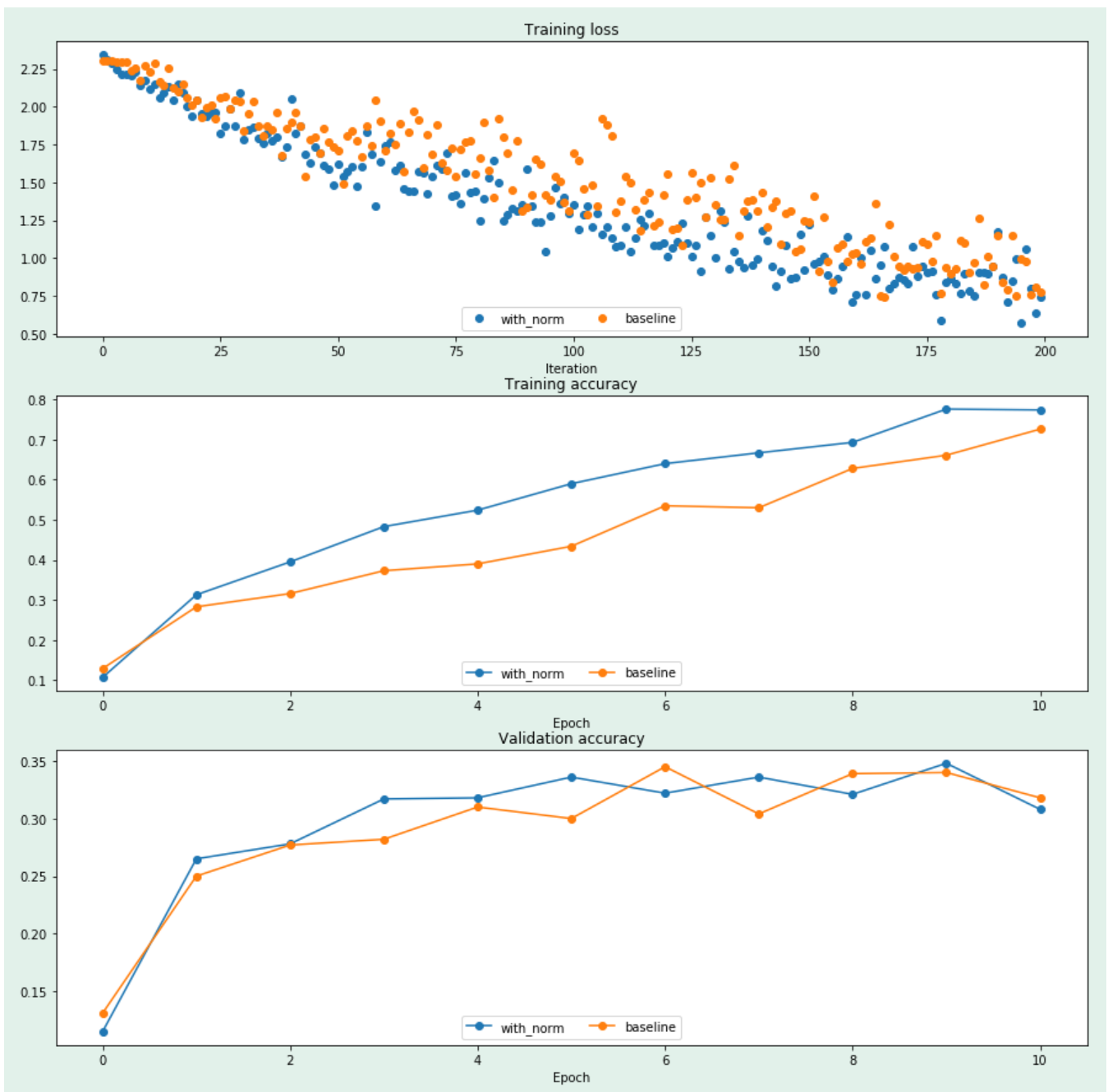
```

```

Solver without batch norm:
(Iteration 1 / 200) loss: 2.302332
(Epoch 0 / 10) train acc: 0.129000; val_acc: 0.131000
(Epoch 1 / 10) train acc: 0.283000; val_acc: 0.250000
(Iteration 21 / 200) loss: 2.041970
(Epoch 2 / 10) train acc: 0.316000; val_acc: 0.277000
(Iteration 41 / 200) loss: 1.900473
(Epoch 3 / 10) train acc: 0.373000; val_acc: 0.282000
(Iteration 61 / 200) loss: 1.713156
(Epoch 4 / 10) train acc: 0.390000; val_acc: 0.310000
(Iteration 81 / 200) loss: 1.662209
(Epoch 5 / 10) train acc: 0.434000; val_acc: 0.300000
(Iteration 101 / 200) loss: 1.696059
(Epoch 6 / 10) train acc: 0.535000; val_acc: 0.345000
(Iteration 121 / 200) loss: 1.557986
(Epoch 7 / 10) train acc: 0.530000; val_acc: 0.304000
(Iteration 141 / 200) loss: 1.432189
(Epoch 8 / 10) train acc: 0.628000; val_acc: 0.339000
(Iteration 161 / 200) loss: 1.033932
(Epoch 9 / 10) train acc: 0.661000; val_acc: 0.340000
(Iteration 181 / 200) loss: 0.901034
(Epoch 10 / 10) train acc: 0.726000; val_acc: 0.318000

```

Run the following to visualize the results from two networks trained above. You should find that using batch normalization helps the network to converge much faster.



Batch normalization and initialization

We will now run a small experiment to study the interaction of batch normalization and weight initialization.

The first cell will train 8-layer networks both with and without batch normalization using different scales for weight initialization. The second layer will plot training accuracy, validation set accuracy, and training loss as a function of the weight initialization scale.

```
Running weight scale 1 / 20
Running weight scale 2 / 20
Running weight scale 3 / 20
Running weight scale 4 / 20
Running weight scale 5 / 20
Running weight scale 6 / 20
Running weight scale 7 / 20
Running weight scale 8 / 20
Running weight scale 9 / 20
Running weight scale 10 / 20
Running weight scale 11 / 20
Running weight scale 12 / 20
Running weight scale 13 / 20
Running weight scale 14 / 20
Running weight scale 15 / 20
Running weight scale 16 / 20
Running weight scale 17 / 20
Running weight scale 18 / 20
```


Running weight scale 19 / 20
Running weight scale 20 / 20



Inline Question 1:

Describe the results of this experiment. How does the scale of weight initialization affect models with/without batch normalization differently, and why?

Answer:

According to the plot, the batchnorm case usually has both higher validation accuracy and training accuracy than baseline case. And the final training loss for batchnorm case is lower.

As the scale of weight initialization increases, the training accuracy and validation accuracy both increase until the scale reaches about 10^{-1} and start to decrease after that. The training loss behaves almost in the opposite way.

With batch normalization, the accuracy is higher when the value of scale of weight is very small and the change of accuracy is more smoothly and regularly compared to the situation without batch normalization. Because with normalization, the result will be adjusted if the scale of input features is extremely different. In this way, the gradient descent can reduce the oscillations when approaching the minimum point and converge faster. It reduces the impact from earlier layers on later layers.

Batch normalization and batch size

We will now run a small experiment to study the interaction of batch normalization and batch size.

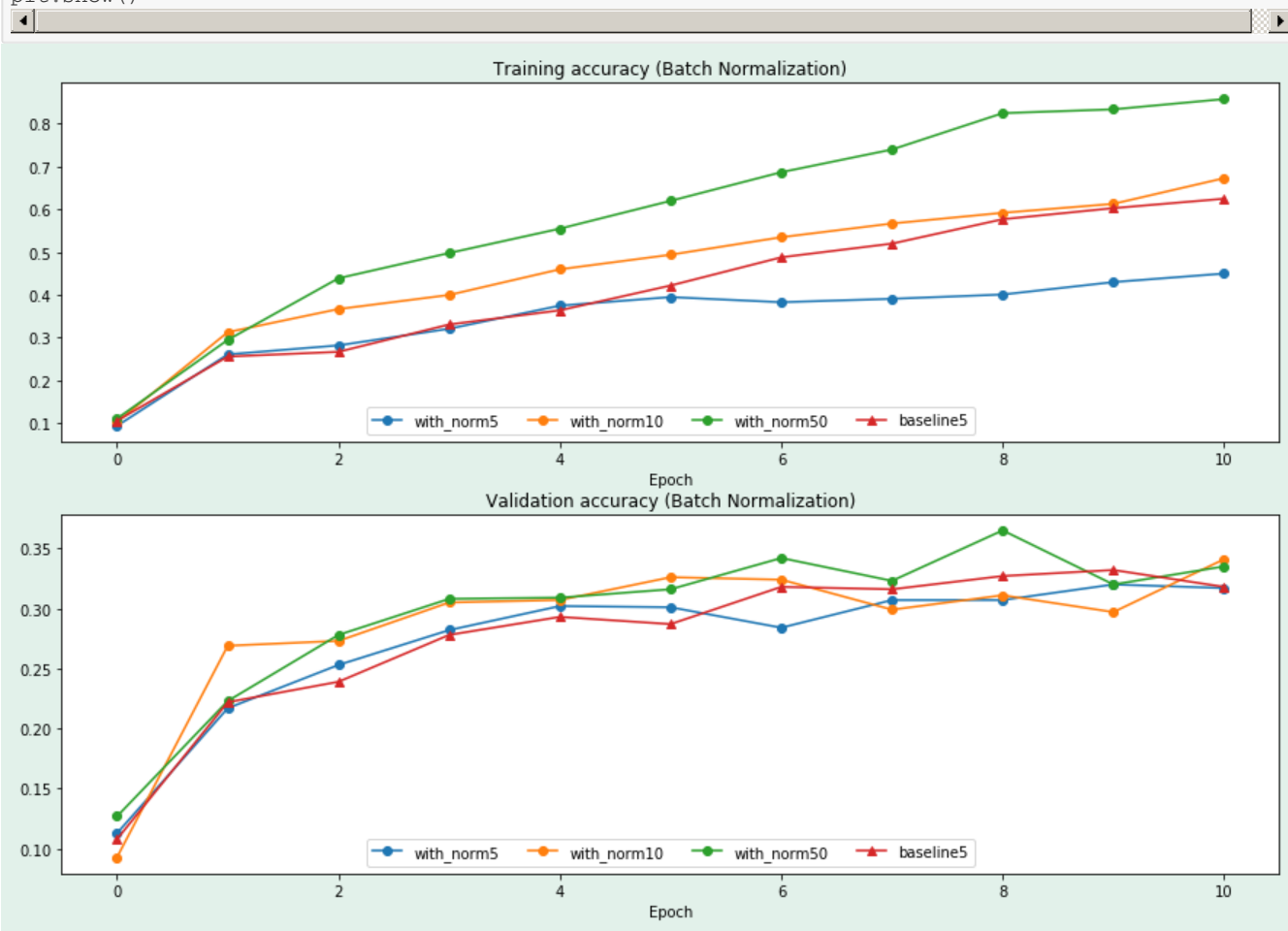
The first cell will train 6-layer networks both with and without batch normalization using different batch sizes. The second layer will plot training accuracy and validation set accuracy over time.

```
No normalization: batch size = 5
Normalization: batch size = 5
Normalization: batch size = 10
Normalization: batch size = 50
```

In [13]:

```
plt.subplot(2, 1, 1)
plot_training_history('Training accuracy (Batch Normalization)', 'Epoch', solver_bsize,
bn_solvers_bsize, \
                    lambda x: x.train_acc_history, bl_marker='-^', bn_marker='-o', labels=batch_s
zes)
plt.subplot(2, 1, 2)
plot_training_history('Validation accuracy (Batch Normalization)', 'Epoch', solver_bsize,
bn_solvers_bsize, \
                    lambda x: x.val_acc_history, bl_marker='-^', bn_marker='-o', labels=batch_siz
s)

plt.gcf().set_size_inches(15, 10)
plt.show()
```



Inline Question 2:

Describe the results of this experiment. What does this imply about the relationship between batch normalization and batch size? Why is this relationship observed?

Answer:

When the batch size is 5, baseline case has slightly better training than batchnorm case as epoch gets larger. But as the batch size gets larger, BN cases tend to converge faster and perform better on training accuracy, but the validation accuracy is not significantly influenced.

This result shows that batch normalization tends to improve the training speed and accuracy but it doesn't have much influence on validation. And when the batch size is small, this normalization may even have an adverse effect. Because increasing batch size can make the steps more accurate since the sampling will be closer to the real population.

Layer Normalization

Batch normalization has proved to be effective in making networks easier to train, but the dependency on batch size makes it less useful in complex networks which have a cap on the input batch size due to hardware limitations.

Several alternatives to batch normalization have been proposed to mitigate this problem; one such technique is Layer Normalization [2]. Instead of normalizing over the batch, we normalize over the features. In other words, when using Layer Normalization, each feature vector corresponding to a single datapoint is normalized based on the sum of all terms within that feature vector.

[2] [Ba, Jimmy Lei, Jamie Ryan Kiros, and Geoffrey E. Hinton. "Layer Normalization." stat 1050 \(2016\): 21.](#)

Inline Question 3:

Which of these data preprocessing steps is analogous to batch normalization, and which is analogous to layer normalization?

1. Scaling each image in the dataset, so that the RGB channels for each row of pixels within an image sums up to 1.
2. Scaling each image in the dataset, so that the RGB channels for all pixels within an image sums up to 1.
3. Subtracting the mean image of the dataset from each image in the dataset.
4. Setting all RGB values to either 0 or 1 depending on a given threshold.

Answer:

1, 2 are like layer normalization. 3 is like batch normalization

Layer Normalization: Implementation

Now you'll implement layer normalization. This step should be relatively straightforward, as conceptually the implementation is almost identical to that of batch normalization. One significant difference though is that for layer normalization, we do not keep track of the moving moments, and the testing phase is identical to the training phase, where the mean and variance are directly calculated per datapoint.

Here's what you need to do:

- In `cs231n/layers.py`, implement the forward pass for layer normalization in the function `layernorm_backward`.

Run the cell below to check your results.

- In `cs231n/layers.py`, implement the backward pass for layer normalization in the function `layernorm_backward`.

Run the second cell below to check your results.

- Modify `cs231n/classifiers/fc_net.py` to add layer normalization to the `FullyConnectedNet`. When the `normalization` flag is set to `"layernorm"` in the constructor, you should insert a layer normalization layer before each ReLU nonlinearity.

Run the third cell below to run the batch size experiment on layer normalization.

In [14]:

```
# Check the training-time forward pass by checking means and variances
# of features both before and after layer normalization

# Simulate the forward pass for a two-layer network
np.random.seed(231)
N, D1, D2, D3 = 4, 50, 60, 3
X = np.random.randn(N, D1)
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)
a = np.maximum(0, X.dot(W1)).dot(W2)
```

```

print('Before layer normalization:')
print_mean_std(a,axis=1)

gamma = np.ones(D3)
beta = np.zeros(D3)
# Means should be close to zero and stds close to one
print('After layer normalization (gamma=1, beta=0)')
a_norm, _ = layernorm_forward(a, gamma, beta, {'mode': 'train'})
print_mean_std(a_norm,axis=1)

gamma = np.asarray([3.0,3.0,3.0])
beta = np.asarray([5.0,5.0,5.0])
# Now means should be close to beta and stds close to gamma
print('After layer normalization (gamma=', gamma, ', beta=', beta, ')')
a_norm, _ = layernorm_forward(a, gamma, beta, {'mode': 'train'})
print_mean_std(a_norm,axis=1)

```

```

Before layer normalization:
means:  [-59.06673243 -47.60782686 -43.31137368 -26.40991744]
stds:   [10.07429373 28.39478981 35.28360729  4.01831507]

After layer normalization (gamma=1, beta=0)
means:  [-4.81096644e-16  0.00000000e+00  7.40148683e-17 -5.55111512e-16]
stds:   [0.99999995 0.99999999 1.          0.99999969]

After layer normalization (gamma= [3. 3. 3.] , beta= [5. 5. 5.] )
means:  [5. 5. 5. 5.]
stds:   [2.99999985 2.99999998 2.99999999 2.99999907]

```

In [15]:

```

# Gradient check batchnorm backward pass
np.random.seed(231)
N, D = 4, 5
x = 5 * np.random.randn(N, D) + 12
gamma = np.random.randn(D)
beta = np.random.randn(D)
dout = np.random.randn(N, D)

ln_param = {}
fx = lambda x: layernorm_forward(x, gamma, beta, ln_param)[0]
fg = lambda a: layernorm_forward(x, a, beta, ln_param)[0]
fb = lambda b: layernorm_forward(x, gamma, b, ln_param)[0]

dx_num = eval_numerical_gradient_array(fx, x, dout)
da_num = eval_numerical_gradient_array(fg, gamma.copy(), dout)
db_num = eval_numerical_gradient_array(fb, beta.copy(), dout)

_, cache = layernorm_forward(x, gamma, beta, ln_param)
dx, dgamma, dbeta = layernorm_backward(dout, cache)

#You should expect to see relative errors between 1e-12 and 1e-8
print('dx error: ', rel_error(dx_num, dx))
print('dgamma error: ', rel_error(da_num, dgamma))
print('dbeta error: ', rel_error(db_num, dbeta))

dx error:  1.433615657860454e-09
dgamma error:  4.519489546032799e-12
dbeta error:  2.276445013433725e-12

```

Layer Normalization and batch size

We will now run the previous batch size experiment with layer normalization instead of batch normalization. Compared to the previous experiment, you should see a markedly smaller influence of batch size on the training history!

In [16]:

```

ln_solvers_bsize, solver_bsize, batch_sizes = run_batchsize_experiments('layernorm')

plt.subplot(2, 1, 1)
plot_training_history('Training accuracy (Layer Normalization)', 'Epoch', solver_bsize,
ln_solvers_bsize, \
lambda x: x.train_acc_history, bl_marker='-^', bn_marker='-o', labels=batch_s

```

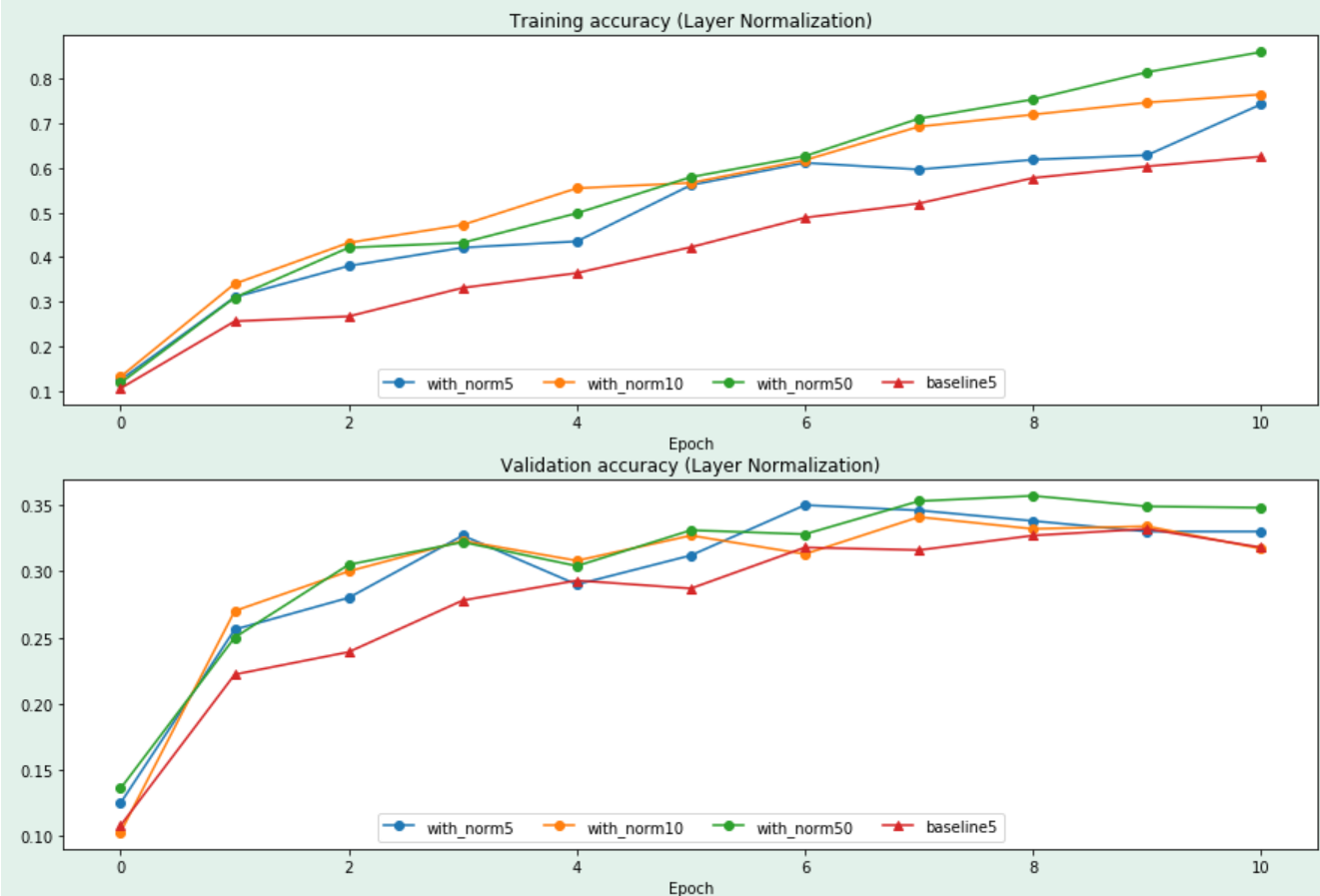
```

zes)
plt.subplot(2, 1, 2)
plot_training_history('Validation accuracy (Layer Normalization)', 'Epoch', solver_bsize,
ln_solvers_bsize, \
                    lambda x: x.val_acc_history, bl_marker='^-^', bn_marker='-o', labels=batch_siz
s)

plt.gcf().set_size_inches(15, 10)
plt.show()

```

No normalization: batch size = 5
 Normalization: batch size = 5
 Normalization: batch size = 10
 Normalization: batch size = 50



Inline Question 4:

When is layer normalization likely to not work well, and why?

1. Using it in a very deep network
2. Having a very small dimension of features
3. Having a high regularization term

Answer:

2,3

1. When the dimension of features is very small, layer normalization may not perform well due to the lack of data about features.
2. When the regularization term is very high, the weights of affine layers will be greatly influenced and the output from affine layer will be really small. In this way, the effect from normalization layer will be reduced.

Dropout

Dropout [1] is a technique for regularizing neural networks by randomly setting some output activations to zero during the forward pass. In this exercise you will implement a dropout layer and modify your fully-connected network to optionally use dropout.

[1] [Geoffrey E. Hinton et al. "Improving neural networks by preventing co-adaptation of feature detectors". arXiv 2012](#)

Dropout forward pass

In the file `cs231n/layers.py`, implement the forward pass for dropout. Since dropout behaves differently during training and testing, make sure to implement the operation for both modes.

Once you have done so, run the cell below to test your implementation.

In [3]:

```
np.random.seed(231)
x = np.random.randn(500, 500) + 10

for p in [0.25, 0.4, 0.7]:
    out, _ = dropout_forward(x, {'mode': 'train', 'p': p})
    out_test, _ = dropout_forward(x, {'mode': 'test', 'p': p})

    print('Running tests with p = ', p)
    print('Mean of input: ', x.mean())
    print('Mean of train-time output: ', out.mean())
    print('Mean of test-time output: ', out_test.mean())
    print('Fraction of train-time output set to zero: ', (out == 0).mean())
    print('Fraction of test-time output set to zero: ', (out_test == 0).mean())
    print()
```

```
Running tests with p = 0.25
Mean of input: 10.000207878477502
Mean of train-time output: 10.014059116977283
Mean of test-time output: 10.000207878477502
Fraction of train-time output set to zero: 0.749784
Fraction of test-time output set to zero: 0.0
```

```
Running tests with p = 0.4
Mean of input: 10.000207878477502
Mean of train-time output: 9.977917658761159
Mean of test-time output: 10.000207878477502
Fraction of train-time output set to zero: 0.600796
Fraction of test-time output set to zero: 0.0
```

```
Running tests with p = 0.7
Mean of input: 10.000207878477502
Mean of train-time output: 9.987811912159426
Mean of test-time output: 10.000207878477502
Fraction of train-time output set to zero: 0.30074
Fraction of test-time output set to zero: 0.0
```

Dropout backward pass

In the file `cs231n/layers.py`, implement the backward pass for dropout. After doing so, run the following cell to numerically gradient-check your implementation.

In [4]:

```
np.random.seed(231)
x = np.random.randn(10, 10) + 10
dout = np.random.randn(*x.shape)

dropout_param = {'mode': 'train', 'p': 0.2, 'seed': 123}
out, cache = dropout_forward(x, dropout_param)
dx = dropout_backward(dout, cache)
```

```
dx_num = eval_numerical_gradient_array(lambda xx: dropout_forward(xx, dropout_param)[0], x, dout)

# Error should be around e-10 or less
print('dx relative error: ', rel_error(dx, dx_num))

dx relative error: 5.44560814873387e-11
```

Inline Question 1:

What happens if we do not divide the values being passed through inverse dropout by p in the dropout layer? Why does that happen?

Answer:

The test output will not be identical to the training output.

Because if we use dropout without the division, the expected training output from input x will become $px + (1-p)0 = px$ while the test output is still x .

Fully-connected nets with Dropout

In the file `cs231n/classifiers/fc_net.py`, modify your implementation to use dropout. Specifically, if the constructor of the network receives a value that is not 1 for the `dropout` parameter, then the net should add a dropout layer immediately after every ReLU nonlinearity. After doing so, run the following to numerically gradient-check your implementation.

In [5]:

```
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 dropout in [1, 0.75, 0.5]:
    print('Running check with dropout = ', dropout)
    model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
                              weight_scale=5e-2, dtype=np.float64,
                              dropout=dropout, seed=123)

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

    # Relative errors should be around e-6 or less; Note that it's fine
    # if for dropout=1 you have W2 error be on the order of e-5.
    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])))
    print()
```

```
Running check with dropout = 1
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 dropout = 0.75
Initial loss: 2.302371489704412
W1 relative error: 1.90e-07
W2 relative error: 4.76e-06
W3 relative error: 2.60e-08
b1 relative error: 4.73e-09
b2 relative error: 1.82e-09
b3 relative error: 1.70e-10

Running check with dropout = 0.5
Initial loss: 2.3042759220785896
W1 relative error: 3.11e-07
```

```
W2 relative error: 1.84e-08
W3 relative error: 5.35e-08
b1 relative error: 2.58e-08
b2 relative error: 2.99e-09
b3 relative error: 1.13e-10
```

Regularization experiment

As an experiment, we will train a pair of two-layer networks on 500 training examples: one will use no dropout, and one will use a keep probability of 0.25. We will then visualize the training and validation accuracies of the two networks over time.

In [6]:

```
# Train two identical nets, one with dropout and one without
np.random.seed(231)
num_train = 500
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 = {}
dropout_choices = [1, 0.25]
for dropout in dropout_choices:
    model = FullyConnectedNet([500], dropout=dropout)
    print(dropout)

    solver = Solver(model, small_data,
                    num_epochs=25, batch_size=100,
                    update_rule='adam',
                    optim_config={
                        'learning_rate': 5e-4,
                    },
                    verbose=True, print_every=100)
    solver.train()
    solvers[dropout] = solver
    print()
```

```
1
(Iteration 1 / 125) loss: 7.856643
(Epoch 0 / 25) train acc: 0.260000; val_acc: 0.184000
(Epoch 1 / 25) train acc: 0.416000; val_acc: 0.258000
(Epoch 2 / 25) train acc: 0.482000; val_acc: 0.276000
(Epoch 3 / 25) train acc: 0.532000; val_acc: 0.277000
(Epoch 4 / 25) train acc: 0.600000; val_acc: 0.271000
(Epoch 5 / 25) train acc: 0.708000; val_acc: 0.299000
(Epoch 6 / 25) train acc: 0.722000; val_acc: 0.282000
(Epoch 7 / 25) train acc: 0.832000; val_acc: 0.255000
(Epoch 8 / 25) train acc: 0.878000; val_acc: 0.269000
(Epoch 9 / 25) train acc: 0.902000; val_acc: 0.275000
(Epoch 10 / 25) train acc: 0.890000; val_acc: 0.261000
(Epoch 11 / 25) train acc: 0.930000; val_acc: 0.282000
(Epoch 12 / 25) train acc: 0.958000; val_acc: 0.300000
(Epoch 13 / 25) train acc: 0.964000; val_acc: 0.305000
(Epoch 14 / 25) train acc: 0.962000; val_acc: 0.318000
(Epoch 15 / 25) train acc: 0.966000; val_acc: 0.304000
(Epoch 16 / 25) train acc: 0.982000; val_acc: 0.307000
(Epoch 17 / 25) train acc: 0.968000; val_acc: 0.322000
(Epoch 18 / 25) train acc: 0.990000; val_acc: 0.319000
(Epoch 19 / 25) train acc: 0.984000; val_acc: 0.300000
(Epoch 20 / 25) train acc: 0.970000; val_acc: 0.306000
(Iteration 101 / 125) loss: 0.121149
(Epoch 21 / 25) train acc: 0.978000; val_acc: 0.307000
(Epoch 22 / 25) train acc: 0.954000; val_acc: 0.315000
(Epoch 23 / 25) train acc: 0.964000; val_acc: 0.320000
(Epoch 24 / 25) train acc: 0.994000; val_acc: 0.304000
(Epoch 25 / 25) train acc: 0.978000; val_acc: 0.304000
```

```
0.25
(Iteration 1 / 125) loss: 17.318479
(Epoch 0 / 25) train acc: 0.230000; val_acc: 0.177000
(Epoch 1 / 25) train acc: 0.378000; val_acc: 0.243000
```



```

(Epoch 2 / 25) train acc: 0.402000; val_acc: 0.254000
(Epoch 3 / 25) train acc: 0.502000; val_acc: 0.276000
(Epoch 4 / 25) train acc: 0.528000; val_acc: 0.298000
(Epoch 5 / 25) train acc: 0.562000; val_acc: 0.297000
(Epoch 6 / 25) train acc: 0.628000; val_acc: 0.291000
(Epoch 7 / 25) train acc: 0.622000; val_acc: 0.299000
(Epoch 8 / 25) train acc: 0.684000; val_acc: 0.312000
(Epoch 9 / 25) train acc: 0.716000; val_acc: 0.289000
(Epoch 10 / 25) train acc: 0.724000; val_acc: 0.297000
(Epoch 11 / 25) train acc: 0.760000; val_acc: 0.309000
(Epoch 12 / 25) train acc: 0.788000; val_acc: 0.284000
(Epoch 13 / 25) train acc: 0.822000; val_acc: 0.314000
(Epoch 14 / 25) train acc: 0.828000; val_acc: 0.349000
(Epoch 15 / 25) train acc: 0.852000; val_acc: 0.340000
(Epoch 16 / 25) train acc: 0.856000; val_acc: 0.301000
(Epoch 17 / 25) train acc: 0.850000; val_acc: 0.299000
(Epoch 18 / 25) train acc: 0.862000; val_acc: 0.331000
(Epoch 19 / 25) train acc: 0.874000; val_acc: 0.337000
(Epoch 20 / 25) train acc: 0.872000; val_acc: 0.309000
(Iteration 101 / 125) loss: 4.035628
(Epoch 21 / 25) train acc: 0.900000; val_acc: 0.333000
(Epoch 22 / 25) train acc: 0.904000; val_acc: 0.300000
(Epoch 23 / 25) train acc: 0.886000; val_acc: 0.292000
(Epoch 24 / 25) train acc: 0.886000; val_acc: 0.309000
(Epoch 25 / 25) train acc: 0.898000; val_acc: 0.321000

```

In [7]:

```
# Plot train and validation accuracies of the two models
```

```

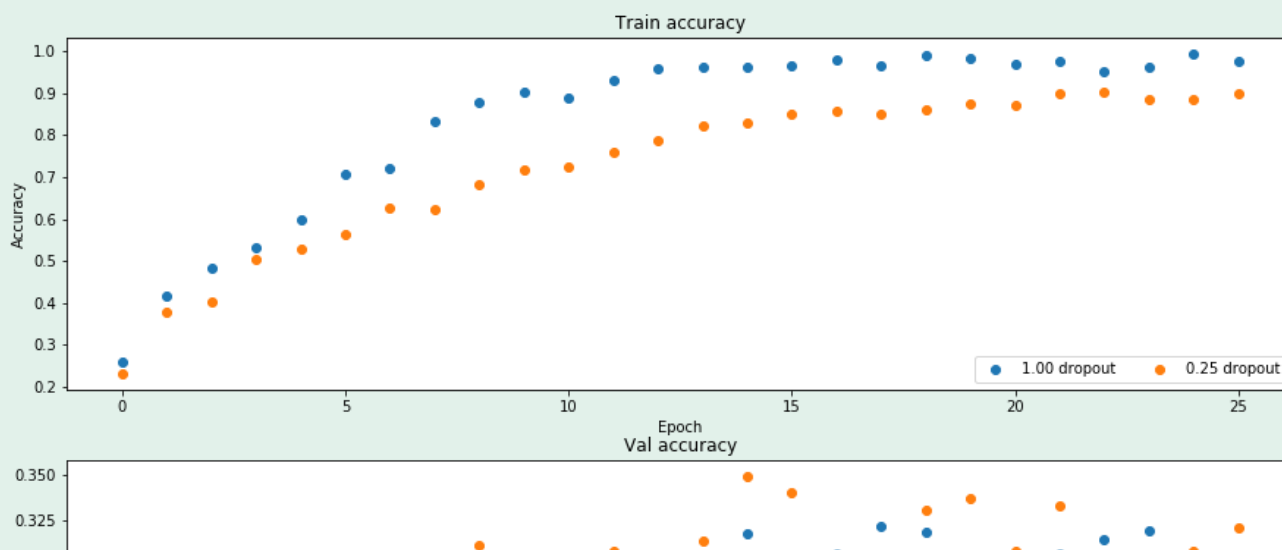
train_accs = []
val_accs = []
for dropout in dropout_choices:
    solver = solvers[dropout]
    train_accs.append(solver.train_acc_history[-1])
    val_accs.append(solver.val_acc_history[-1])

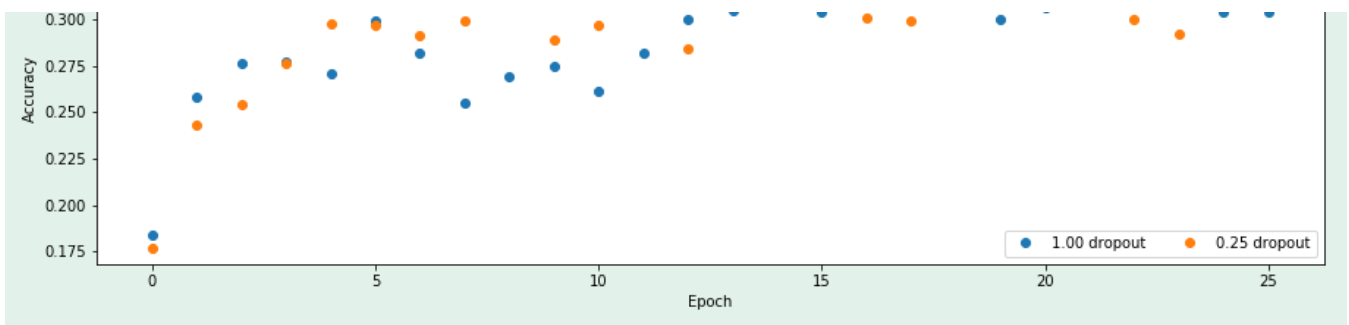
plt.subplot(3, 1, 1)
for dropout in dropout_choices:
    plt.plot(solvers[dropout].train_acc_history, 'o', label='%.2f dropout' % dropout)
plt.title('Train accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend(ncol=2, loc='lower right')

plt.subplot(3, 1, 2)
for dropout in dropout_choices:
    plt.plot(solvers[dropout].val_acc_history, 'o', label='%.2f dropout' % dropout)
plt.title('Val accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend(ncol=2, loc='lower right')

plt.gcf().set_size_inches(15, 15)
plt.show()

```





Inline Question 2:

Compare the validation and training accuracies with and without dropout -- what do your results suggest about dropout as a regularizer?

Answer:

The training accuracy with dropout is slightly worse. But at test time, the validation accuracy with dropout is slightly better.

Based on these results we can see that dropout plays a role as a regularizer. It prevents models from overfitting to training data so that at test time they generalize better.

Inline Question 3:

Suppose we are training a deep fully-connected network for image classification, with dropout after hidden layers (parameterized by keep probability p). If we are concerned about overfitting, how should we modify p (if at all) when we decide to decrease the size of the hidden layers (that is, the number of nodes in each layer)?

Answer:

We should slight decrease the value of p .

Convolutional Networks

So far we have worked with deep fully-connected networks, using them to explore different optimization strategies and network architectures. Fully-connected networks are a good testbed for experimentation because they are very computationally efficient, but in practice all state-of-the-art results use convolutional networks instead.

First you will implement several layer types that are used in convolutional networks. You will then use these layers to train a convolutional network on the CIFAR-10 dataset.

Convolution: Naive forward pass

The core of a convolutional network is the convolution operation. In the file `cs231n/layers.py`, implement the forward pass for the convolution layer in the function `conv_forward_naive`.

You don't have to worry too much about efficiency at this point; just write the code in whatever way you find most clear.

You can test your implementation by running the following:

In [3]:

```
x_shape = (2, 3, 4, 4)
w_shape = (3, 3, 4, 4)
x = np.linspace(-0.1, 0.5, num=np.prod(x_shape)).reshape(x_shape)
w = np.linspace(-0.2, 0.3, num=np.prod(w_shape)).reshape(w_shape)
b = np.linspace(-0.1, 0.2, num=3)

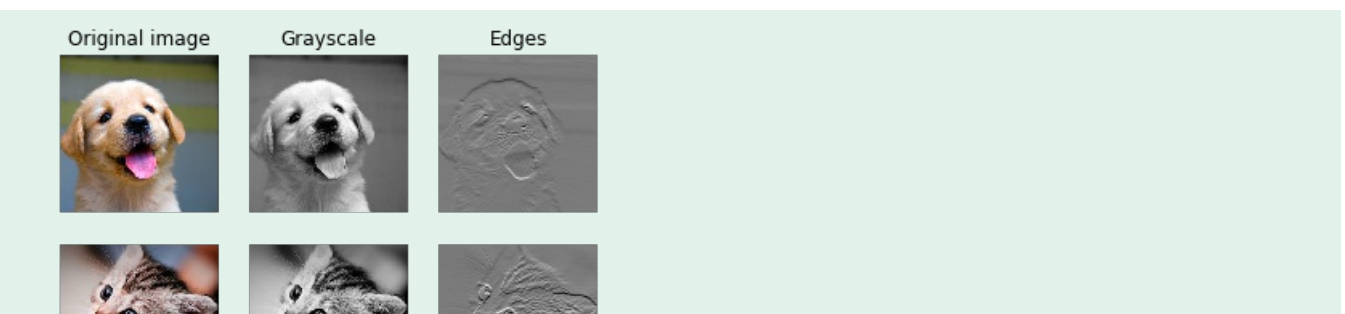
conv_param = {'stride': 2, 'pad': 1}
out, _ = conv_forward_naive(x, w, b, conv_param)
correct_out = np.array([[[[-0.08759809, -0.10987781],
                          [-0.18387192, -0.2109216 ]],
                        [[ 0.21027089,  0.21661097],
                          [ 0.22847626,  0.23004637]],
                        [[ 0.50813986,  0.54309974],
                          [ 0.64082444,  0.67101435]]],
                       [[[-0.98053589, -1.03143541],
                          [-1.19128892, -1.24695841]],
                        [[ 0.69108355,  0.66880383],
                          [ 0.59480972,  0.56776003]],
                        [[ 2.36270298,  2.36904306],
                          [ 2.38090835,  2.38247847]]]])

# Compare your output to ours; difference should be around e-8
print('Testing conv_forward_naive')
print('difference: ', rel_error(out, correct_out))
```

```
Testing conv_forward_naive
difference:  2.2121476417505994e-08
```

Aside: Image processing via convolutions

As fun way to both check your implementation and gain a better understanding of the type of operation that convolutional layers can perform, we will set up an input containing two images and manually set up filters that perform common image processing operations (grayscale conversion and edge detection). The convolution forward pass will apply these operations to each of the input images. We can then visualize the results as a sanity check.





Convolution: Naive backward pass

Implement the backward pass for the convolution operation in the function `conv_backward_naive` in the file `cs231n/layers.py`. Again, you don't need to worry too much about computational efficiency.

When you are done, run the following to check your backward pass with a numeric gradient check.

In [5]:

```
np.random.seed(231)
x = np.random.randn(4, 3, 5, 5)
w = np.random.randn(2, 3, 3, 3)
b = np.random.randn(2,)
dout = np.random.randn(4, 2, 5, 5)
conv_param = {'stride': 1, 'pad': 1}

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

out, cache = conv_forward_naive(x, w, b, conv_param)
dx, dw, db = conv_backward_naive(dout, cache)

# Your errors should be around e-8 or less.
print('Testing conv_backward_naive function')
print('dx error: ', rel_error(dx, dx_num))
print('dw error: ', rel_error(dw, dw_num))
print('db error: ', rel_error(db, db_num))
```

```
Testing conv_backward_naive function
dx error:  1.159803161159293e-08
dw error:  2.2471264748452487e-10
db error:  3.3726153958780465e-11
```

Max-Pooling: Naive forward

Implement the forward pass for the max-pooling operation in the function `max_pool_forward_naive` in the file `cs231n/layers.py`. Again, don't worry too much about computational efficiency.

Check your implementation by running the following:

In [6]:

```
x_shape = (2, 3, 4, 4)
x = np.linspace(-0.3, 0.4, num=np.prod(x_shape)).reshape(x_shape)
pool_param = {'pool_width': 2, 'pool_height': 2, 'stride': 2}

out, _ = max_pool_forward_naive(x, pool_param)

correct_out = np.array([[[[-0.26315789, -0.24842105],
                           [-0.20421053, -0.18947368]],
                          [[-0.14526316, -0.13052632],
                           [-0.08631579, -0.07157895]],
                          [[-0.02736842, -0.01263158],
                           [ 0.03157895,  0.04631579]]],
                        [[[ 0.09052632,  0.10526316],
                           [ 0.14947368,  0.16421053]],
                          [[ 0.20842105,  0.22315789],
                           [ 0.26736842,  0.28210526]],
                          [[ 0.32631579,  0.34105263],
                           [ 0.38526316,  0.4          ]]]])

# Compare your output with ours. Difference should be on the order of e-8.
print('Testing max_pool_forward_naive function:')
```

```
print('difference: ', rel_error(out, correct_out))
```

```
Testing max_pool_forward_naive function:  
difference: 4.1666665157267834e-08
```

Max-Pooling: Naive backward

Implement the backward pass for the max-pooling operation in the function `max_pool_backward_naive` in the file `cs231n/layers.py`. You don't need to worry about computational efficiency.

Check your implementation with numeric gradient checking by running the following:

In [7]:

```
np.random.seed(231)
x = np.random.randn(3, 2, 8, 8)
dout = np.random.randn(3, 2, 4, 4)
pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}

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

out, cache = max_pool_forward_naive(x, pool_param)
dx = max_pool_backward_naive(dout, cache)

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

Testing max_pool_backward_naive function:
dx error: 3.27562514223145e-12
```

Fast layers

Making convolution and pooling layers fast can be challenging. To spare you the pain, we've provided fast implementations of the forward and backward passes for convolution and pooling layers in the file `cs231n/fast_layers.py`.

The fast convolution implementation depends on a Cython extension; to compile it you need to run the following from the `cs231n` directory:

```
python setup.py build_ext --inplace
```

The API for the fast versions of the convolution and pooling layers is exactly the same as the naive versions that you implemented above: the forward pass receives data, weights, and parameters and produces outputs and a cache object; the backward pass receives upstream derivatives and the cache object and produces gradients with respect to the data and weights.

NOTE: The fast implementation for pooling will only perform optimally if the pooling regions are non-overlapping and tile the input. If these conditions are not met then the fast pooling implementation will not be much faster than the naive implementation.

You can compare the performance of the naive and fast versions of these layers by running the following:

In [8]:

```
# Rel errors should be around e-9 or less
from cs231n.fast_layers import conv_forward_fast, conv_backward_fast
from time import time
np.random.seed(231)
x = np.random.randn(100, 3, 31, 31)
w = np.random.randn(25, 3, 3, 3)
b = np.random.randn(25,)
dout = np.random.randn(100, 25, 16, 16)
conv_param = {'stride': 2, 'pad': 1}

t0 = time()
out_naive, cache_naive = conv_forward_naive(x, w, b, conv_param)
t1 = time()
out_fast, cache_fast = conv_forward_fast(x, w, b, conv_param)
t2 = time()

print('Testing conv_forward_fast:')
print('Naive: %fs' % (t1 - t0))
print('Fast: %fs' % (t2 - t1))
print('Speedup: %fx' % ((t1 - t0) / (t2 - t1)))
```

```

print('Speedup: %fx' % ((t1 - t0) / (t2 - t1)))
print('Difference: ', rel_error(out_naive, out_fast))

t0 = time()
dx_naive, dw_naive, db_naive = conv_backward_naive(dout, cache_naive)
t1 = time()
dx_fast, dw_fast, db_fast = conv_backward_fast(dout, cache_fast)
t2 = time()

print('\nTesting conv_backward_fast:')
print('Naive: %fs' % (t1 - t0))
print('Fast: %fs' % (t2 - t1))
print('Speedup: %fx' % ((t1 - t0) / (t2 - t1)))
print('dx difference: ', rel_error(dx_naive, dx_fast))
print('dw difference: ', rel_error(dw_naive, dw_fast))
print('db difference: ', rel_error(db_naive, db_fast))

```

```

Testing conv_forward_fast:
Naive: 5.172025s
Fast: 0.016499s
Speedup: 313.474278x
Difference: 4.926407851494105e-11

Testing conv_backward_fast:
Naive: 8.520645s
Fast: 0.008236s
Speedup: 1034.539731x
dx difference: 1.949764775345631e-11
dw difference: 3.7012612707710095e-13
db difference: 3.1393858025571252e-15

```

In [9]:

```

# Relative errors should be close to 0.0
from cs231n.fast_layers import max_pool_forward_fast, max_pool_backward_fast
np.random.seed(231)
x = np.random.randn(100, 3, 32, 32)
dout = np.random.randn(100, 3, 16, 16)
pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}

t0 = time()
out_naive, cache_naive = max_pool_forward_naive(x, pool_param)
t1 = time()
out_fast, cache_fast = max_pool_forward_fast(x, pool_param)
t2 = time()

print('Testing pool_forward_fast:')
print('Naive: %fs' % (t1 - t0))
print('fast: %fs' % (t2 - t1))
print('speedup: %fx' % ((t1 - t0) / (t2 - t1)))
print('difference: ', rel_error(out_naive, out_fast))

t0 = time()
dx_naive = max_pool_backward_naive(dout, cache_naive)
t1 = time()
dx_fast = max_pool_backward_fast(dout, cache_fast)
t2 = time()

print('\nTesting pool_backward_fast:')
print('Naive: %fs' % (t1 - t0))
print('fast: %fs' % (t2 - t1))
print('speedup: %fx' % ((t1 - t0) / (t2 - t1)))
print('dx difference: ', rel_error(dx_naive, dx_fast))

```

```

Testing pool_forward_fast:
Naive: 0.383643s
fast: 0.001968s
speedup: 194.926348x
difference: 0.0

Testing pool_backward_fast:
Naive: 0.996497s
fast: 0.012093s
speedup: 82.405580x
dx difference: 0.0

```

Convolutional "sandwich" layers

Convolutional sandwich layers

Previously we introduced the concept of "sandwich" layers that combine multiple operations into commonly used patterns. In the file `cs231n/layer_utils.py` you will find sandwich layers that implement a few commonly used patterns for convolutional networks.

Run the cells below to sanity check they're working.

In [10]:

```
from cs231n.layer_utils import conv_relu_pool_forward, conv_relu_pool_backward
np.random.seed(231)
x = np.random.randn(2, 3, 16, 16)
w = np.random.randn(3, 3, 3, 3)
b = np.random.randn(3,)
dout = np.random.randn(2, 3, 8, 8)
conv_param = {'stride': 1, 'pad': 1}
pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}

out, cache = conv_relu_pool_forward(x, w, b, conv_param, pool_param)
dx, dw, db = conv_relu_pool_backward(dout, cache)

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

# Relative errors should be around e-8 or less
print('Testing conv_relu_pool')
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 conv_relu_pool
dx error: 6.514336569263308e-09
dw error: 1.490843753539445e-08
db error: 2.037390356217257e-09
```

In [11]:

```
from cs231n.layer_utils import conv_relu_forward, conv_relu_backward
np.random.seed(231)
x = np.random.randn(2, 3, 8, 8)
w = np.random.randn(3, 3, 3, 3)
b = np.random.randn(3,)
dout = np.random.randn(2, 3, 8, 8)
conv_param = {'stride': 1, 'pad': 1}

out, cache = conv_relu_forward(x, w, b, conv_param)
dx, dw, db = conv_relu_backward(dout, cache)

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

# Relative errors should be around e-8 or less
print('Testing conv_relu:')
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 conv_relu:
dx error: 3.5600610115232832e-09
dw error: 2.2497700915729298e-10
db error: 1.3087619975802167e-10
```

Three-layer ConvNet

Now that you have implemented all the necessary layers, we can put them together into a simple convolutional network.

Open the file `cs231n/classifiers/cnn.py` and complete the implementation of the `ThreeLayerConvNet` class. Remember you can use the fast/sandwich layers (already imported for you) in your implementation. Run the following cells to help you debug:

Sanity check loss

After you build a new network, one of the first things you should do is sanity check the loss. When we use the softmax loss, we expect the loss for random weights (and no regularization) to be about $\log(C)$ for C classes. When we add regularization the loss should go up slightly.

In [12]:

```
model = ThreeLayerConvNet()

N = 50
X = np.random.randn(N, 3, 32, 32)
y = np.random.randint(10, size=N)

loss, grads = model.loss(X, y)
print('Initial loss (no regularization): ', loss)

model.reg = 0.5
loss, grads = model.loss(X, y)
print('Initial loss (with regularization): ', loss)

Initial loss (no regularization):  2.302586071243987
Initial loss (with regularization):  2.508255638232932
```

Gradient check

After the loss looks reasonable, use numeric gradient checking to make sure that your backward pass is correct. When you use numeric gradient checking you should use a small amount of artificial data and a small number of neurons at each layer. Note: correct implementations may still have relative errors up to the order of $e-2$.

In [13]:

```
num_inputs = 2
input_dim = (3, 16, 16)
reg = 0.0
num_classes = 10
np.random.seed(231)
X = np.random.randn(num_inputs, *input_dim)
y = np.random.randint(num_classes, size=num_inputs)

model = ThreeLayerConvNet(num_filters=3, filter_size=3,
                           input_dim=input_dim, hidden_dim=7,
                           dtype=np.float64)

loss, grads = model.loss(X, y)
# Errors should be small, but correct implementations may have
# relative errors up to the order of e-2
for param_name in sorted(grads):
    f = lambda _: model.loss(X, y)[0]
    param_grad_num = eval_numerical_gradient(f, model.params[param_name], verbose=False, h=1e-6)
    e = rel_error(param_grad_num, grads[param_name])
    print('%s max relative error: %e' % (param_name, rel_error(param_grad_num, grads[param_name])))

W1 max relative error: 1.380104e-04
W2 max relative error: 1.822723e-02
W3 max relative error: 3.064049e-04
b1 max relative error: 3.477652e-05
b2 max relative error: 2.516375e-03
b3 max relative error: 7.945660e-10
```

Overfit small data

A nice trick is to train your model with just a few training samples. You should be able to overfit small datasets, which will result in very high training accuracy and comparatively low validation accuracy.

In [16]:

```
np.random.seed(231)

num_train = 100
```



```

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'],
}

model = ThreeLayerConvNet(weight_scale=1e-2)

solver = Solver(model, small_data,
                num_epochs=15, batch_size=50,
                update_rule='adam',
                optim_config={
                    'learning_rate': 1e-3,
                },
                verbose=True, print_every=1)

solver.train()

(Iteration 1 / 30) loss: 2.414060
(Epoch 0 / 15) train acc: 0.200000; val_acc: 0.137000
(Iteration 2 / 30) loss: 3.102925
(Epoch 1 / 15) train acc: 0.140000; val_acc: 0.087000
(Iteration 3 / 30) loss: 2.270331
(Iteration 4 / 30) loss: 2.096705
(Epoch 2 / 15) train acc: 0.240000; val_acc: 0.094000
(Iteration 5 / 30) loss: 1.838880
(Iteration 6 / 30) loss: 1.934188
(Epoch 3 / 15) train acc: 0.510000; val_acc: 0.173000
(Iteration 7 / 30) loss: 1.827912
(Iteration 8 / 30) loss: 1.639574
(Epoch 4 / 15) train acc: 0.520000; val_acc: 0.188000
(Iteration 9 / 30) loss: 1.330082
(Iteration 10 / 30) loss: 1.756115
(Epoch 5 / 15) train acc: 0.630000; val_acc: 0.167000
(Iteration 11 / 30) loss: 1.024162
(Iteration 12 / 30) loss: 1.041826
(Epoch 6 / 15) train acc: 0.750000; val_acc: 0.229000
(Iteration 13 / 30) loss: 1.142777
(Iteration 14 / 30) loss: 0.835706
(Epoch 7 / 15) train acc: 0.790000; val_acc: 0.247000
(Iteration 15 / 30) loss: 0.587786
(Iteration 16 / 30) loss: 0.645509
(Epoch 8 / 15) train acc: 0.820000; val_acc: 0.252000
(Iteration 17 / 30) loss: 0.786844
(Iteration 18 / 30) loss: 0.467054
(Epoch 9 / 15) train acc: 0.820000; val_acc: 0.178000
(Iteration 19 / 30) loss: 0.429880
(Iteration 20 / 30) loss: 0.635498
(Epoch 10 / 15) train acc: 0.900000; val_acc: 0.206000
(Iteration 21 / 30) loss: 0.365807
(Iteration 22 / 30) loss: 0.284220
(Epoch 11 / 15) train acc: 0.820000; val_acc: 0.201000
(Iteration 23 / 30) loss: 0.469343
(Iteration 24 / 30) loss: 0.509369
(Epoch 12 / 15) train acc: 0.920000; val_acc: 0.211000
(Iteration 25 / 30) loss: 0.111638
(Iteration 26 / 30) loss: 0.145389
(Epoch 13 / 15) train acc: 0.930000; val_acc: 0.213000
(Iteration 27 / 30) loss: 0.155576
(Iteration 28 / 30) loss: 0.143400
(Epoch 14 / 15) train acc: 0.960000; val_acc: 0.212000
(Iteration 29 / 30) loss: 0.158156
(Iteration 30 / 30) loss: 0.118937
(Epoch 15 / 15) train acc: 0.990000; val_acc: 0.220000

```

Plotting the loss, training accuracy, and validation accuracy should show clear overfitting:

In [17]:

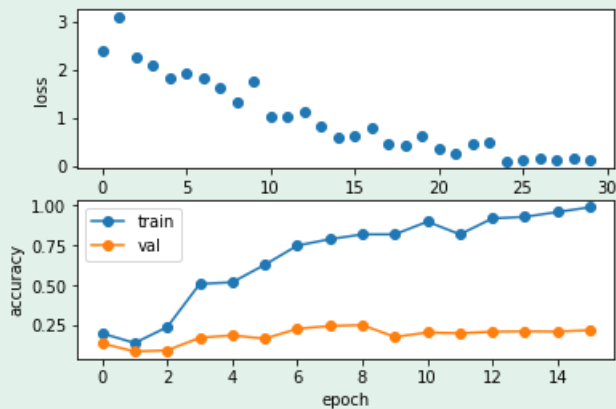
```

plt.subplot(2, 1, 1)
plt.plot(solver.loss_history, 'o')
plt.xlabel('iteration')
plt.ylabel('loss')

plt.subplot(2, 1, 2)
plt.plot(solver.train_acc_history, '-o')
plt.plot(solver.val_acc_history, '-o')

```

```
plt.legend(['train', 'val'], loc='upper left')
plt.xlabel('epoch')
plt.ylabel('accuracy')
plt.show()
```



Train the net

By training the three-layer convolutional network for one epoch, you should achieve greater than 40% accuracy on the training set:

In [16]:

```
model = ThreeLayerConvNet(weight_scale=0.001, hidden_dim=500, reg=0.001)

solver = Solver(model, data,
                 num_epochs=1, batch_size=50,
                 update_rule='adam',
                 optim_config={
                     'learning_rate': 1e-3,
                 },
                 verbose=True, print_every=20)

solver.train()
```

```
(Iteration 1 / 980) loss: 2.304740
(Epoch 0 / 1) train acc: 0.103000; val_acc: 0.107000
(Iteration 21 / 980) loss: 2.098229
(Iteration 41 / 980) loss: 1.949740
(Iteration 61 / 980) loss: 1.824802
(Iteration 81 / 980) loss: 1.879293
(Iteration 101 / 980) loss: 1.923165
(Iteration 121 / 980) loss: 1.725399
(Iteration 141 / 980) loss: 1.884197
(Iteration 161 / 980) loss: 1.935079
(Iteration 181 / 980) loss: 1.784737
(Iteration 201 / 980) loss: 1.908147
(Iteration 221 / 980) loss: 1.885975
(Iteration 241 / 980) loss: 1.573188
(Iteration 261 / 980) loss: 1.732478
(Iteration 281 / 980) loss: 1.817697
(Iteration 301 / 980) loss: 1.752375
(Iteration 321 / 980) loss: 1.832898
(Iteration 341 / 980) loss: 1.564610
(Iteration 361 / 980) loss: 1.866280
(Iteration 381 / 980) loss: 1.356685
(Iteration 401 / 980) loss: 1.876740
(Iteration 421 / 980) loss: 1.553664
(Iteration 441 / 980) loss: 1.646373
(Iteration 461 / 980) loss: 1.794048
(Iteration 481 / 980) loss: 1.652758
(Iteration 501 / 980) loss: 1.687621
(Iteration 521 / 980) loss: 1.722508
(Iteration 541 / 980) loss: 1.745398
(Iteration 561 / 980) loss: 1.624082
(Iteration 581 / 980) loss: 1.203774
(Iteration 601 / 980) loss: 1.654945
(Iteration 621 / 980) loss: 1.525178
(Iteration 641 / 980) loss: 1.579597
(Iteration 661 / 980) loss: 1.760286
(Iteration 681 / 980) loss: 1.653154
(Iteration 701 / 980) loss: 1.520100
```

```
(Iteration 721 / 980) loss: 1.524231
(Iteration 741 / 980) loss: 1.609275
(Iteration 761 / 980) loss: 1.685576
(Iteration 781 / 980) loss: 1.866236
(Iteration 801 / 980) loss: 1.682262
(Iteration 821 / 980) loss: 1.857055
(Iteration 841 / 980) loss: 1.556042
(Iteration 861 / 980) loss: 1.646650
(Iteration 881 / 980) loss: 1.657959
(Iteration 901 / 980) loss: 1.423653
(Iteration 921 / 980) loss: 1.588974
(Iteration 941 / 980) loss: 1.613119
(Iteration 961 / 980) loss: 1.616299
(Epoch 1 / 1) train acc: 0.496000; val_acc: 0.489000
```

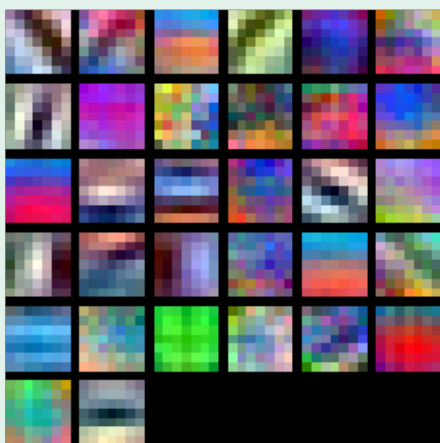
Visualize Filters

You can visualize the first-layer convolutional filters from the trained network by running the following:

In [17]:

```
from cs231n.vis_utils import visualize_grid

grid = visualize_grid(model.params['W1'].transpose(0, 2, 3, 1))
plt.imshow(grid.astype('uint8'))
plt.axis('off')
plt.gcf().set_size_inches(5, 5)
plt.show()
```



Spatial Batch Normalization

We already saw that batch normalization is a very useful technique for training deep fully-connected networks. As proposed in the original paper (link in `BatchNormalization.ipynb`), batch normalization can also be used for convolutional networks, but we need to tweak it a bit; the modification will be called "spatial batch normalization."

Normally batch-normalization accepts inputs of shape (N, D) and produces outputs of shape (N, D) , where we normalize across the minibatch dimension N . For data coming from convolutional layers, batch normalization needs to accept inputs of shape (N, C, H, W) and produce outputs of shape (N, C, H, W) where the N dimension gives the minibatch size and the (H, W) dimensions give the spatial size of the feature map.

If the feature map was produced using convolutions, then we expect every feature channel's statistics e.g. mean, variance to be relatively consistent both between different images, and different locations within the same image -- after all, every feature channel is produced by the same convolutional filter! Therefore spatial batch normalization computes a mean and variance for each of the C feature channels by computing statistics over the minibatch dimension N as well the spatial dimensions H and W .

[1] [Sergey Ioffe and Christian Szegedy. "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift". ICML 2015.](#)

Spatial batch normalization: forward

In the file `cs231n/layers.py`, implement the forward pass for spatial batch normalization in the function `spatial_batchnorm_forward`. Check your implementation by running the following:

In [12]:

```
np.random.seed(231)
# Check the training-time forward pass by checking means and variances
# of features both before and after spatial batch normalization

N, C, H, W = 2, 3, 4, 5
x = 4 * np.random.randn(N, C, H, W) + 10

print('Before spatial batch normalization:')
print('  Shape: ', x.shape)
print('  Means: ', x.mean(axis=(0, 2, 3)))
print('  Stds: ', x.std(axis=(0, 2, 3)))

# Means should be close to zero and stds close to one
gamma, beta = np.ones(C), np.zeros(C)
bn_param = {'mode': 'train'}
out, _ = spatial_batchnorm_forward(x, gamma, beta, bn_param)
print('After spatial batch normalization:')
print('  Shape: ', out.shape)
print('  Means: ', out.mean(axis=(0, 2, 3)))
print('  Stds: ', out.std(axis=(0, 2, 3)))

# Means should be close to beta and stds close to gamma
gamma, beta = np.asarray([3, 4, 5]), np.asarray([6, 7, 8])
out, _ = spatial_batchnorm_forward(x, gamma, beta, bn_param)
print('After spatial batch normalization (nontrivial gamma, beta):')
print('  Shape: ', out.shape)
print('  Means: ', out.mean(axis=(0, 2, 3)))
print('  Stds: ', out.std(axis=(0, 2, 3)))
```

```
Before spatial batch normalization:
  Shape: (2, 3, 4, 5)
  Means: [9.33463814 8.90909116 9.11056338]
  Stds:  [3.61447857 3.19347686 3.5168142 ]
After spatial batch normalization:
  Shape: (2, 3, 4, 5)
  Means: [ 6.18949336e-16  5.99520433e-16 -1.22124533e-16]
  Stds:  [0.99999962 0.99999951 0.9999996 ]
After spatial batch normalization (nontrivial gamma, beta):
  Shape: (2, 3, 4, 5)
  Means: [6. 7. 8.]
  Stds:  [2.99999885 3.99999804 4.99999798]
```

In [13]:

```
np.random.seed(231)
# Check the test-time forward pass by running the training-time
# forward pass many times to warm up the running averages, and then
# checking the means and variances of activations after a test-time
# forward pass.
N, C, H, W = 10, 4, 11, 12

bn_param = {'mode': 'train'}
gamma = np.ones(C)
beta = np.zeros(C)
for t in range(50):
    x = 2.3 * np.random.randn(N, C, H, W) + 13
    spatial_batchnorm_forward(x, gamma, beta, bn_param)
bn_param['mode'] = 'test'
x = 2.3 * np.random.randn(N, C, H, W) + 13
a_norm, _ = spatial_batchnorm_forward(x, gamma, beta, bn_param)

# Means should be close to zero and stds close to one, but will be
# noisier than training-time forward passes.
print('After spatial batch normalization (test-time):')
print('  means: ', a_norm.mean(axis=(0, 2, 3)))
print('  stds: ', a_norm.std(axis=(0, 2, 3)))
```

```
After spatial batch normalization (test-time):
  means: [-0.08034406 0.07562881 0.05716371 0.04378383]
  stds:  [0.96718744 1.0299714 1.02887624 1.00585577]
```

Spatial batch normalization: backward

In the file `cs231n/layers.py`, implement the backward pass for spatial batch normalization in the function `spatial_batchnorm_backward`. Run the following to check your implementation using a numeric gradient check:

In [14]:

```
np.random.seed(231)
N, C, H, W = 2, 3, 4, 5
x = 5 * np.random.randn(N, C, H, W) + 12
gamma = np.random.randn(C)
beta = np.random.randn(C)
dout = np.random.randn(N, C, H, W)

bn_param = {'mode': 'train'}
fx = lambda x: spatial_batchnorm_forward(x, gamma, beta, bn_param)[0]
fg = lambda a: spatial_batchnorm_forward(x, gamma, beta, bn_param)[0]
fb = lambda b: spatial_batchnorm_forward(x, gamma, beta, bn_param)[0]

dx_num = eval_numerical_gradient_array(fx, x, dout)
da_num = eval_numerical_gradient_array(fg, gamma, dout)
db_num = eval_numerical_gradient_array(fb, beta, dout)

#You should expect errors of magnitudes between 1e-12~1e-06
_, cache = spatial_batchnorm_forward(x, gamma, beta, bn_param)
dx, dgamma, dbeta = spatial_batchnorm_backward(dout, cache)
print('dx error: ', rel_error(dx_num, dx))
print('dgamma error: ', rel_error(da_num, dgamma))
print('dbeta error: ', rel_error(db_num, dbeta))

dx error:  2.786648197756335e-07
dgamma error:  7.0974817113608705e-12
dbeta error:  3.275608725278405e-12
```

Group Normalization

In the previous notebook, we mentioned that Layer Normalization is an alternative normalization technique that mitigates the batch size limitations of Batch Normalization. However, as the authors of [2] observed, Layer Normalization does not perform as well as Batch Normalization when used with Convolutional Layers:

With fully connected layers, all the hidden units in a layer tend to make similar contributions to the final prediction, and re-centering and rescaling the summed inputs to a layer works well. However, the assumption of similar contributions is no longer true for convolutional neural networks. The large number of the hidden units whose receptive fields lie near the boundary of the image are rarely turned on and thus have very different statistics from the rest of the hidden units within the same layer.

The authors of [3] propose an intermediary technique. In contrast to Layer Normalization, where you normalize over the entire feature per-datapoint, they suggest a consistent splitting of each per-datapoint feature into G groups, and a per-group per-datapoint normalization instead.

****Visual comparison of the normalization techniques discussed so far (image edited from [3])****

Even though an assumption of equal contribution is still being made within each group, the authors hypothesize that this is not as problematic, as innate grouping arises within features for visual recognition. One example they use to illustrate this is that many high-performance handcrafted features in traditional Computer Vision have terms that are explicitly grouped together. Take for example Histogram of Oriented Gradients [4]-- after computing histograms per spatially local block, each per-block histogram is normalized before being concatenated together to form the final feature vector.

You will now implement Group Normalization. Note that this normalization technique that you are to implement in the following cells was introduced and published to ECCV just in 2018 -- this truly is still an ongoing and excitingly active field of research!

[2] [Ba, Jimmy Lei, Jamie Ryan Kiros, and Geoffrey E. Hinton. "Layer Normalization." stat 1050 \(2016\): 21.](#)

[3] [Wu, Yuxin, and Kaiming He. "Group Normalization." arXiv preprint arXiv:1803.08494 \(2018\).](#)

[4] [N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In Computer Vision and Pattern Recognition \(CVPR\), 2005.](#)

Group normalization: forward

In the file `cs231n/layers.py`, implement the forward pass for group normalization in the function

`spatial_groupnorm_forward`. Check your implementation by running the following:

In [15]:

```
np.random.seed(231)
# Check the training-time forward pass by checking means and variances
# of features both before and after spatial batch normalization

N, C, H, W = 2, 6, 4, 5
G = 2
x = 4 * np.random.randn(N, C, H, W) + 10
x_g = x.reshape((N*G,-1))
print('Before spatial group normalization:')
print('  Shape: ', x.shape)
print('  Means: ', x_g.mean(axis=1))
print('  Stds: ', x_g.std(axis=1))

# Means should be close to zero and stds close to one
gamma, beta = np.ones((1,C,1,1)), np.zeros((1,C,1,1))
bn_param = {'mode': 'train'}

out, _ = spatial_groupnorm_forward(x, gamma, beta, G, bn_param)
out_g = out.reshape((N*G,-1))
print('After spatial group normalization:')
print('  Shape: ', out.shape)
print('  Means: ', out_g.mean(axis=1))
print('  Stds: ', out_g.std(axis=1))

Before spatial group normalization:
  Shape: (2, 6, 4, 5)
  Means: [9.72505327 8.51114185 8.9147544  9.43448077]
  Stds:  [3.67070958 3.09892597 4.27043622 3.97521327]
After spatial group normalization:
  Shape: (2, 6, 4, 5)
  Means: [-2.14643118e-16  5.25505565e-16  2.65528340e-16 -3.38618023e-16]
  Stds:  [0.99999963 0.99999948 0.99999973 0.99999968]
```

Spatial group normalization: backward

In the file `cs231n/layers.py`, implement the backward pass for spatial batch normalization in the function

`spatial_groupnorm_backward`. Run the following to check your implementation using a numeric gradient check:

In [23]:

```
np.random.seed(231)
N, C, H, W = 2, 6, 4, 5
G = 2
x = 5 * np.random.randn(N, C, H, W) + 12
gamma = np.random.randn(1,C,1,1)
beta = np.random.randn(1,C,1,1)
dout = np.random.randn(N, C, H, W)

gn_param = {}
fx = lambda x: spatial_groupnorm_forward(x, gamma, beta, G, gn_param)[0]
fg = lambda a: spatial_groupnorm_forward(x, gamma, beta, G, gn_param)[0]
fb = lambda b: spatial_groupnorm_forward(x, gamma, beta, G, gn_param)[0]

dx_num = eval_numerical_gradient_array(fx, x, dout)
da_num = eval_numerical_gradient_array(fg, gamma, dout)
db_num = eval_numerical_gradient_array(fb, beta, dout)

_, cache = spatial_groupnorm_forward(x, gamma, beta, G, gn_param)
dx, dgamma, dbeta = spatial_groupnorm_backward(dout, cache)
#You should expect errors of magnitudes between 1e-12~1e-07
print('dx error: ', rel_error(dx_num, dx))
print('dgamma error: ', rel_error(da_num, dgamma))
print('dbeta error: ', rel_error(db_num, dbeta))

dx error:  7.413109648400194e-08
dgamma error:  9.468195772749234e-12
dbeta error:  3.354494437653335e-12
```


What's this PyTorch business?

You've written a lot of code in this assignment to provide a whole host of neural network functionality. Dropout, Batch Norm, and 2D convolutions are some of the workhorses of deep learning in computer vision. You've also worked hard to make your code efficient and vectorized.

For the last part of this assignment, though, we're going to leave behind your beautiful codebase and instead migrate to one of two popular deep learning frameworks: in this instance, PyTorch (or TensorFlow, if you choose to use that notebook).

Part I. Preparation

First, we load the CIFAR-10 dataset. This might take a couple minutes the first time you do it, but the files should stay cached after that.

In previous parts of the assignment we had to write our own code to download the CIFAR-10 dataset, preprocess it, and iterate through it in minibatches; PyTorch provides convenient tools to automate this process for us.

```
using device: cuda
```

Part II. Barebones PyTorch

PyTorch ships with high-level APIs to help us define model architectures conveniently, which we will cover in Part II of this tutorial. In this section, we will start with the barebone PyTorch elements to understand the autograd engine better. After this exercise, you will come to appreciate the high-level model API more.

We will start with a simple fully-connected ReLU network with two hidden layers and no biases for CIFAR classification. This implementation computes the forward pass using operations on PyTorch Tensors, and uses PyTorch autograd to compute gradients. It is important that you understand every line, because you will write a harder version after the example.

When we create a PyTorch Tensor with `requires_grad=True`, then operations involving that Tensor will not just compute values; they will also build up a computational graph in the background, allowing us to easily backpropagate through the graph to compute gradients of some Tensors with respect to a downstream loss. Concretely if `x` is a Tensor with `x.requires_grad == True` then after backpropagation `x.grad` will be another Tensor holding the gradient of `x` with respect to the scalar loss at the end.

```
Before flattening: tensor([[[[ 0,  1],
      [ 2,  3],
      [ 4,  5]]],

      [[[ 6,  7],
      [ 8,  9],
      [10, 11]]]])
After flattening: tensor([ 0,  1,  2,  3,  4,  5],
      [ 6,  7,  8,  9, 10, 11])
```

```
torch.Size([64, 10])
```

Barebones PyTorch: Three-Layer ConvNet

Here you will complete the implementation of the function `three_layer_convnet`, which will perform the forward pass of a three-layer convolutional network. Like above, we can immediately test our implementation by passing zeros through the network. The network should have the following architecture:

1. A convolutional layer (with bias) with `channel_1` filters, each with shape `KW1 x KH1`, and zero-padding of two
2. ReLU nonlinearity
3. A convolutional layer (with bias) with `channel_2` filters, each with shape `KW2 x KH2`, and zero-padding of one
4. ReLU nonlinearity
5. Fully-connected layer with bias, producing scores for `C` classes.

Note that we have **no softmax activation** here after our fully-connected layer: this is because PyTorch's cross entropy loss performs a softmax activation for you, and by bundling that step in makes computation more efficient

a softmax activation for you, and by bundling that step in makes computation more efficient.

HINT: For convolutions: <http://pytorch.org/docs/stable/nn.html#torch.nn.functional.conv2d>; pay attention to the shapes of convolutional filters!

In [6]:

```
def three_layer_convnet(x, params):
    """
    Performs the forward pass of a three-layer convolutional network with the
    architecture defined above.

    Inputs:
    - x: A PyTorch Tensor of shape (N, 3, H, W) giving a minibatch of images
    - params: A list of PyTorch Tensors giving the weights and biases for the
      network; should contain the following:
      - conv_w1: PyTorch Tensor of shape (channel_1, 3, KH1, KW1) giving weights
        for the first convolutional layer
      - conv_b1: PyTorch Tensor of shape (channel_1,) giving biases for the first
        convolutional layer
      - conv_w2: PyTorch Tensor of shape (channel_2, channel_1, KH2, KW2) giving
        weights for the second convolutional layer
      - conv_b2: PyTorch Tensor of shape (channel_2,) giving biases for the second
        convolutional layer
      - fc_w: PyTorch Tensor giving weights for the fully-connected layer. Can you
        figure out what the shape should be?
      - fc_b: PyTorch Tensor giving biases for the fully-connected layer. Can you
        figure out what the shape should be?

    Returns:
    - scores: PyTorch Tensor of shape (N, C) giving classification scores for x
    """
    conv_w1, conv_b1, conv_w2, conv_b2, fc_w, fc_b = params
    scores = None
    #####
    # TODO: Implement the forward pass for the three-layer ConvNet.          #
    #####
    # ****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)****

    # (32 + 2 * 2 - 5)/1 + 1 = 32
    # (32 + 2 * 1 - 3)/1 + 1 = 32

    # zero-padding of two
    conv1 = F.conv2d(x, weight = conv_w1, bias = conv_b1, padding = 2)
    relu1 = F.relu(conv1)

    # zero-padding of one
    conv2 = F.conv2d(relu1, weight = conv_w2, bias = conv_b2, padding = 1)
    relu2 = F.relu(conv2)

    relu2_flat = flatten(relu2)
    scores = relu2_flat.mm(fc_w) + fc_b

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

After defining the forward pass of the ConvNet above, run the following cell to test your implementation.

When you run this function, scores should have shape (64, 10).

```
torch.Size([64, 10])
```

Barebones PyTorch: Initialization

Let's write a couple utility methods to initialize the weight matrices for our models.

- `random_weight(shape)` initializes a weight tensor with the Kaiming normalization method.
- `zero_weight(shape)` initializes a weight tensor with all zeros. Useful for instantiating bias parameters.

The `random_weight` function uses the Kaiming normal initialization method, described in:

Out [8]:

```
tensor([[ 0.8353,  0.0717,  0.8552, -0.1880,  0.3146],
        [ 0.1567, -0.3085, -0.2543, -0.6643, -0.9352],
        [-0.0133,  0.1022, -0.4872,  0.5114, -0.8374]], device='cuda:0',
        requires_grad=True)
```

Barebones PyTorch: Check Accuracy

When training the model we will use the following function to check the accuracy of our model on the training or validation sets.

When checking accuracy we don't need to compute any gradients; as a result we don't need PyTorch to build a computational graph for us when we compute scores. To prevent a graph from being built we scope our computation under a `torch.no_grad()` context manager.

BareBones PyTorch: Training Loop

We can now set up a basic training loop to train our network. We will train the model using stochastic gradient descent without momentum. We will use `torch.functional.cross_entropy` to compute the loss; you can [read about it here](#).

The training loop takes as input the neural network function, a list of initialized parameters (`[w1, w2]` in our example), and learning rate.

BareBones PyTorch: Train a Two-Layer Network

Now we are ready to run the training loop. We need to explicitly allocate tensors for the fully connected weights, `w1` and `w2`.

Each minibatch of CIFAR has 64 examples, so the tensor shape is `[64, 3, 32, 32]`.

After flattening, `x` shape should be `[64, 3 * 32 * 32]`. This will be the size of the first dimension of `w1`. The second dimension of `w1` is the hidden layer size, which will also be the first dimension of `w2`.

Finally, the output of the network is a 10-dimensional vector that represents the probability distribution over 10 classes.

You don't need to tune any hyperparameters but you should see accuracies above 40% after training for one epoch.

In [11]:

```
hidden_layer_size = 4000
learning_rate = 1e-2

w1 = random_weight((3 * 32 * 32, hidden_layer_size))
w2 = random_weight((hidden_layer_size, 10))

train_part2(two_layer_fc, [w1, w2], learning_rate)
```

```
Iteration 0, loss = 3.4906
Checking accuracy on the val set
Got 157 / 1000 correct (15.70%)
```

```
Iteration 100, loss = 2.6464
Checking accuracy on the val set
Got 326 / 1000 correct (32.60%)
```

```
Iteration 200, loss = 1.9548
Checking accuracy on the val set
Got 388 / 1000 correct (38.80%)
```

```
Iteration 300, loss = 1.9776
Checking accuracy on the val set
Got 380 / 1000 correct (38.00%)
```

```
Iteration 400, loss = 2.3733
Checking accuracy on the val set
Got 409 / 1000 correct (40.90%)
```

```
Iteration 500, loss = 1.8391
Checking accuracy on the val set
Got 441 / 1000 correct (44.10%)
```

```
Iteration 600, loss = 1.9783
Checking accuracy on the val set
Got 375 / 1000 correct (37.50%)
```

```
Iteration 700, loss = 1.4147
Checking accuracy on the val set
Got 438 / 1000 correct (43.80%)
```

BareBones PyTorch: Training a ConvNet

In the below you should use the functions defined above to train a three-layer convolutional network on CIFAR. The network should have the following architecture:

1. Convolutional layer (with bias) with 32 5x5 filters, with zero-padding of 2
2. ReLU
3. Convolutional layer (with bias) with 16 3x3 filters, with zero-padding of 1
4. ReLU
5. Fully-connected layer (with bias) to compute scores for 10 classes

You should initialize your weight matrices using the `random_weight` function defined above, and you should initialize your bias vectors using the `zero_weight` function above.

You don't need to tune any hyperparameters, but if everything works correctly you should achieve an accuracy above 42% after one epoch.

In [12]:

```
learning_rate = 3e-3

channel_1 = 32
channel_2 = 16

conv_w1 = None
conv_b1 = None
conv_w2 = None
conv_b2 = None
fc_w = None
fc_b = None

#####
# TODO: Initialize the parameters of a three-layer ConvNet. #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

# Basically just change tensor.zeros to random_weights and adopt the same implementation
# in the function three_layer_convnet_test

conv_w1 = random_weight((channel_1, 3, 5, 5))
conv_b1 = zero_weight((channel_1,))
conv_w2 = random_weight((channel_2, channel_1, 3, 3))
conv_b2 = zero_weight((channel_2,))
fc_w = random_weight((channel_2 * 32 * 32, 10))
fc_b = zero_weight((10,))

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

params = [conv_w1, conv_b1, conv_w2, conv_b2, fc_w, fc_b]
train_part2(three_layer_convnet, params, learning_rate)
```

```
Iteration 0, loss = 2.4838
Checking accuracy on the val set
Got 110 / 1000 correct (11.00%)
```

```
Iteration 100, loss = 2.1007
Checking accuracy on the val set
Got 321 / 1000 correct (32.10%)
```

```
Iteration 200, loss = 1.7121
Checking accuracy on the val set
Got 387 / 1000 correct (38.70%)
```

```
Iteration 300, loss = 1.7689
Checking accuracy on the val set
Got 419 / 1000 correct (41.90%)
```

```
Iteration 400, loss = 1.5850
Checking accuracy on the val set
Got 454 / 1000 correct (45.40%)
```

```
Iteration 500, loss = 1.5140
Checking accuracy on the val set
Got 453 / 1000 correct (45.30%)
```

```
Iteration 600, loss = 1.5323
Checking accuracy on the val set
Got 456 / 1000 correct (45.60%)
```

```
Iteration 700, loss = 1.4435
Checking accuracy on the val set
Got 472 / 1000 correct (47.20%)
```

Part III. PyTorch Module API

Barebone PyTorch requires that we track all the parameter tensors by hand. This is fine for small networks with a few tensors, but it would be extremely inconvenient and error-prone to track tens or hundreds of tensors in larger networks.

PyTorch provides the `nn.Module` API for you to define arbitrary network architectures, while tracking every learnable parameters for you. In Part II, we implemented SGD ourselves. PyTorch also provides the `torch.optim` package that implements all the common optimizers, such as RMSProp, Adagrad, and Adam. It even supports approximate second-order methods like L-BFGS! You can refer to the [doc](#) for the exact specifications of each optimizer.

To use the Module API, follow the steps below:

1. Subclass `nn.Module`. Give your network class an intuitive name like `TwoLayerFC`.
2. In the constructor `__init__()`, define all the layers you need as class attributes. Layer objects like `nn.Linear` and `nn.Conv2d` are themselves `nn.Module` subclasses and contain learnable parameters, so that you don't have to instantiate the raw tensors yourself. `nn.Module` will track these internal parameters for you. Refer to the [doc](#) to learn more about the dozens of builtin layers. **Warning:** don't forget to call the `super().__init__()` first!
3. In the `forward()` method, define the *connectivity* of your network. You should use the attributes defined in `__init__` as function calls that take tensor as input and output the "transformed" tensor. Do *not* create any new layers with learnable parameters in `forward()`! All of them must be declared upfront in `__init__`.

After you define your Module subclass, you can instantiate it as an object and call it just like the NN forward function in part II.

Module API: Two-Layer Network

Here is a concrete example of a 2-layer fully connected network:

In [13]:

```
class TwoLayerFC(nn.Module):
    def __init__(self, input_size, hidden_size, num_classes):
        super().__init__()
        # assign layer objects to class attributes
        self.fc1 = nn.Linear(input_size, hidden_size)
        # nn.init package contains convenient initialization methods
        # http://pytorch.org/docs/master/nn.html#torch-nn-init
        nn.init.kaiming_normal_(self.fc1.weight)
        self.fc2 = nn.Linear(hidden_size, num_classes)
        nn.init.kaiming_normal_(self.fc2.weight)

    def forward(self, x):
        # forward always defines connectivity
        x = flatten(x)
        scores = self.fc2(F.relu(self.fc1(x)))
```

```

        return scores

def test_TwoLayerFC():
    input_size = 50
    x = torch.zeros((64, input_size), dtype=dtype) # minibatch size 64, feature dimension 50
    model = TwoLayerFC(input_size, 42, 10)
    scores = model(x)
    print(scores.size()) # you should see [64, 10]
test_TwoLayerFC()

torch.Size([64, 10])

```

Module API: Three-Layer ConvNet

It's your turn to implement a 3-layer ConvNet followed by a fully connected layer. The network architecture should be the same as in Part II:

1. Convolutional layer with `channel_1` 5x5 filters with zero-padding of 2
2. ReLU
3. Convolutional layer with `channel_2` 3x3 filters with zero-padding of 1
4. ReLU
5. Fully-connected layer to `num_classes` classes

You should initialize the weight matrices of the model using the Kaiming normal initialization method.

HINT: <http://pytorch.org/docs/stable/nn.html#conv2d>

After you implement the three-layer ConvNet, the `test_ThreeLayerConvNet` function will run your implementation; it should print `(64, 10)` for the shape of the output scores.

In [14]:

```

class ThreeLayerConvNet(nn.Module):
    def __init__(self, in_channel, channel_1, channel_2, num_classes):
        super().__init__()
        #####
        # TODO: Set up the layers you need for a three-layer ConvNet with the #
        # architecture defined above. #
        #####
        # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

        self.conv1 = nn.Conv2d(in_channel, channel_1, 5, padding = 2)
        nn.init.kaiming_normal_(self.conv1.weight)
        nn.init.constant_(self.conv1.bias, 0)

        self.conv2 = nn.Conv2d(channel_1, channel_2, 3, padding = 1)
        nn.init.kaiming_normal_(self.conv2.weight)
        nn.init.constant_(self.conv2.bias, 0)

        self.fc = nn.Linear(channel_2 * 32 * 32, num_classes)
        nn.init.kaiming_normal_(self.fc.weight)
        nn.init.constant_(self.fc.bias, 0)

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

    def forward(self, x):
        scores = None
        #####
        # TODO: Implement the forward function for a 3-layer ConvNet. you #
        # should use the layers you defined in __init__ and specify the #
        # connectivity of those layers in forward() #
        #####
        # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

        # Before moving to Fully-connected layer, flatten the outcome from the second Relu layer
        scores = self.fc(flatten(F.relu(self.conv2(F.relu(self.conv1(x))))))

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

```

```

        return scores

def test_ThreeLayerConvNet():
    x = torch.zeros((64, 3, 32, 32), dtype=dtype) # minibatch size 64, image size [3, 32, 32]
    model = ThreeLayerConvNet(in_channel=3, channel_1=12, channel_2=8, num_classes=10)
    scores = model(x)
    print(scores.size()) # you should see [64, 10]
test_ThreeLayerConvNet()

torch.Size([64, 10])

```

Module API: Check Accuracy

Given the validation or test set, we can check the classification accuracy of a neural network.

This version is slightly different from the one in part II. You don't manually pass in the parameters anymore.

In [15]:

```

def check_accuracy_part34(loader, model):
    if loader.dataset.train:
        print('Checking accuracy on validation set')
    else:
        print('Checking accuracy on test set')
    num_correct = 0
    num_samples = 0
    model.eval() # set model to evaluation mode
    with torch.no_grad():
        for x, y in loader:
            x = x.to(device=device, dtype=dtype) # move to device, e.g. GPU
            y = y.to(device=device, dtype=torch.long)
            scores = model(x)
            _, preds = scores.max(1)
            num_correct += (preds == y).sum()
            num_samples += preds.size(0)
    acc = float(num_correct) / num_samples
    print('Got %d / %d correct (%.2f)' % (num_correct, num_samples, 100 * acc))

```

Module API: Training Loop

We also use a slightly different training loop. Rather than updating the values of the weights ourselves, we use an Optimizer object from the `torch.optim` package, which abstract the notion of an optimization algorithm and provides implementations of most of the algorithms commonly used to optimize neural networks.

In [16]:

```

def train_part34(model, optimizer, epochs=1):
    """
    Train a model on CIFAR-10 using the PyTorch Module API.

    Inputs:
    - model: A PyTorch Module giving the model to train.
    - optimizer: An Optimizer object we will use to train the model
    - epochs: (Optional) A Python integer giving the number of epochs to train for

    Returns: Nothing, but prints model accuracies during training.
    """
    model = model.to(device=device) # move the model parameters to CPU/GPU
    for e in range(epochs):
        for t, (x, y) in enumerate(loader_train):
            model.train() # put model to training mode
            x = x.to(device=device, dtype=dtype) # move to device, e.g. GPU
            y = y.to(device=device, dtype=torch.long)

            scores = model(x)
            loss = F.cross_entropy(scores, y)

            # Zero out all of the gradients for the variables which the optimizer
            # will update.
            optimizer.zero_grad()

            # This is the backwards pass: compute the gradient of the loss with
            # respect to each parameter of the model

```

```

# respect to each parameter of the model.
loss.backward()

# Actually update the parameters of the model using the gradients
# computed by the backwards pass.
optimizer.step()

if t % print_every == 0:
    print('Iteration %d, loss = %.4f' % (t, loss.item()))
    check_accuracy_part34(loader_val, model)
    print()

```

Module API: Train a Two-Layer Network

Now we are ready to run the training loop. In contrast to part II, we don't explicitly allocate parameter tensors anymore.

Simply pass the input size, hidden layer size, and number of classes (i.e. output size) to the constructor of `TwoLayerFC`.

You also need to define an optimizer that tracks all the learnable parameters inside `TwoLayerFC`.

You don't need to tune any hyperparameters, but you should see model accuracies above 40% after training for one epoch.

In [17]:

```

hidden_layer_size = 4000
learning_rate = 1e-2
model = TwoLayerFC(3 * 32 * 32, hidden_layer_size, 10)
optimizer = optim.SGD(model.parameters(), lr=learning_rate)

train_part34(model, optimizer)

```

```

Iteration 0, loss = 2.9949
Checking accuracy on validation set
Got 121 / 1000 correct (12.10)

```

```

Iteration 100, loss = 2.3486
Checking accuracy on validation set
Got 324 / 1000 correct (32.40)

```

```

Iteration 200, loss = 2.1595
Checking accuracy on validation set
Got 373 / 1000 correct (37.30)

```

```

Iteration 300, loss = 1.5870
Checking accuracy on validation set
Got 384 / 1000 correct (38.40)

```

```

Iteration 400, loss = 1.7430
Checking accuracy on validation set
Got 425 / 1000 correct (42.50)

```

```

Iteration 500, loss = 2.2017
Checking accuracy on validation set
Got 405 / 1000 correct (40.50)

```

```

Iteration 600, loss = 1.6490
Checking accuracy on validation set
Got 431 / 1000 correct (43.10)

```

```

Iteration 700, loss = 1.7141
Checking accuracy on validation set
Got 445 / 1000 correct (44.50)

```

Module API: Train a Three-Layer ConvNet

You should now use the Module API to train a three-layer ConvNet on CIFAR. This should look very similar to training the two-layer network! You don't need to tune any hyperparameters, but you should achieve above 45% after training for one epoch.

You should train the model using stochastic gradient descent without momentum.

In [18]:

```

learning_rate = 3e-3

```

```

channel_1 = 32
channel_2 = 16

model = None
optimizer = None
#####
# TODO: Instantiate your ThreeLayerConvNet model and a corresponding optimizer #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

model = ThreeLayerConvNet(3, channel_1, channel_2, 10)
optimizer = optim.SGD(model.parameters(), lr=learning_rate)

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

train_part34(model, optimizer)

Iteration 0, loss = 3.2324
Checking accuracy on validation set
Got 132 / 1000 correct (13.20)

Iteration 100, loss = 1.6502
Checking accuracy on validation set
Got 357 / 1000 correct (35.70)

Iteration 200, loss = 1.7441
Checking accuracy on validation set
Got 403 / 1000 correct (40.30)

Iteration 300, loss = 1.5787
Checking accuracy on validation set
Got 438 / 1000 correct (43.80)

Iteration 400, loss = 1.8331
Checking accuracy on validation set
Got 462 / 1000 correct (46.20)

Iteration 500, loss = 1.8362
Checking accuracy on validation set
Got 468 / 1000 correct (46.80)

Iteration 600, loss = 1.6090
Checking accuracy on validation set
Got 459 / 1000 correct (45.90)

Iteration 700, loss = 1.4880
Checking accuracy on validation set
Got 472 / 1000 correct (47.20)

```

Part IV. PyTorch Sequential API

Part III introduced the PyTorch Module API, which allows you to define arbitrary learnable layers and their connectivity.

For simple models like a stack of feed forward layers, you still need to go through 3 steps: subclass `nn.Module`, assign layers to class attributes in `__init__`, and call each layer one by one in `forward()`. Is there a more convenient way?

Fortunately, PyTorch provides a container Module called `nn.Sequential`, which merges the above steps into one. It is not as flexible as `nn.Module`, because you cannot specify more complex topology than a feed-forward stack, but it's good enough for many use cases.

Sequential API: Two-Layer Network

Let's see how to rewrite our two-layer fully connected network example with `nn.Sequential`, and train it using the training loop defined above.

Again, you don't need to tune any hyperparameters here, but you should achieve above 40% accuracy after one epoch of training.

In [19]:


```

# We need to wrap `flatten` function in a module in order to stack it
# in nn.Sequential
class Flatten(nn.Module):
    def forward(self, x):
        return flatten(x)

hidden_layer_size = 4000
learning_rate = 1e-2

model = nn.Sequential(
    Flatten(),
    nn.Linear(3 * 32 * 32, hidden_layer_size),
    nn.ReLU(),
    nn.Linear(hidden_layer_size, 10),
)

# you can use Nesterov momentum in optim.SGD
optimizer = optim.SGD(model.parameters(), lr=learning_rate,
                       momentum=0.9, nesterov=True)

train_part34(model, optimizer)

```

```

Iteration 0, loss = 2.3438
Checking accuracy on validation set
Got 157 / 1000 correct (15.70)

Iteration 100, loss = 1.8731
Checking accuracy on validation set
Got 382 / 1000 correct (38.20)

Iteration 200, loss = 1.8784
Checking accuracy on validation set
Got 426 / 1000 correct (42.60)

Iteration 300, loss = 1.6531
Checking accuracy on validation set
Got 410 / 1000 correct (41.00)

Iteration 400, loss = 1.6222
Checking accuracy on validation set
Got 433 / 1000 correct (43.30)

Iteration 500, loss = 1.4940
Checking accuracy on validation set
Got 459 / 1000 correct (45.90)

Iteration 600, loss = 1.5447
Checking accuracy on validation set
Got 431 / 1000 correct (43.10)

Iteration 700, loss = 2.0133
Checking accuracy on validation set
Got 464 / 1000 correct (46.40)

```

Sequential API: Three-Layer ConvNet

Here you should use `nn.Sequential` to define and train a three-layer ConvNet with the same architecture we used in Part III:

1. Convolutional layer (with bias) with 32 5x5 filters, with zero-padding of 2
2. ReLU
3. Convolutional layer (with bias) with 16 3x3 filters, with zero-padding of 1
4. ReLU
5. Fully-connected layer (with bias) to compute scores for 10 classes

You should initialize your weight matrices using the `random_weight` function defined above, and you should initialize your bias vectors using the `zero_weight` function above.

You should optimize your model using stochastic gradient descent with Nesterov momentum 0.9.

Again, you don't need to tune any hyperparameters but you should see accuracy above 55% after one epoch of training.

In [20]:

```

channel_1 = 32

```

```

channel_1 = 32
channel_2 = 16
learning_rate = 1e-2

model = None
optimizer = None

#####
# TODO: Rewrite the 2-layer ConvNet with bias from Part III with the #
# Sequential API. #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

# From Piazza, it is OK to use default initialization. No need to do explicit initialization

model = nn.Sequential(
    nn.Conv2d(3, channel_1, 5, padding = 2),
    nn.ReLU(),
    nn.Conv2d(channel_1, channel_2, 3, padding = 1),
    nn.ReLU(),
    Flatten(),
    nn.Linear(channel_2 * 32 * 32, 10)
)

optimizer = optim.SGD(model.parameters(), lr=learning_rate,
                       momentum=0.9, nesterov=True)

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

train_part34(model, optimizer)

```

```

Iteration 0, loss = 2.2971
Checking accuracy on validation set
Got 100 / 1000 correct (10.00)

Iteration 100, loss = 1.7188
Checking accuracy on validation set
Got 422 / 1000 correct (42.20)

Iteration 200, loss = 1.6543
Checking accuracy on validation set
Got 475 / 1000 correct (47.50)

Iteration 300, loss = 1.2017
Checking accuracy on validation set
Got 519 / 1000 correct (51.90)

Iteration 400, loss = 1.4025
Checking accuracy on validation set
Got 535 / 1000 correct (53.50)

Iteration 500, loss = 1.1097
Checking accuracy on validation set
Got 540 / 1000 correct (54.00)

Iteration 600, loss = 1.0994
Checking accuracy on validation set
Got 574 / 1000 correct (57.40)

Iteration 700, loss = 1.0963
Checking accuracy on validation set
Got 578 / 1000 correct (57.80)

```

Part V. CIFAR-10 open-ended challenge

In this section, you can experiment with whatever ConvNet architecture you'd like on CIFAR-10.

Now it's your job to experiment with architectures, hyperparameters, loss functions, and optimizers to train a model that achieves **at least 70% accuracy** on the CIFAR-10 **validation** set within 10 epochs. You can use the `check_accuracy` and `train` functions from above. You can use either `nn.Module` or `nn.Sequential` API.

Describe what you did at the end of this notebook.

Here are the official API documentation for each component. One note: what we call in the class "spatial batch norm" is called "BatchNorm2D" in PyTorch.

- Layers in torch.nn package: <http://pytorch.org/docs/stable/nn.html>
- Activations: <http://pytorch.org/docs/stable/nn.html#non-linear-activations>
- Loss functions: <http://pytorch.org/docs/stable/nn.html#loss-functions>
- Optimizers: <http://pytorch.org/docs/stable/optim.html>

Things you might try:

- **Filter size:** Above we used 5x5; would smaller filters be more efficient?
- **Number of filters:** Above we used 32 filters. Do more or fewer do better?
- **Pooling vs Strided Convolution:** Do you use max pooling or just stride convolutions?
- **Batch normalization:** Try adding spatial batch normalization after convolution layers and vanilla batch normalization after affine layers. Do your networks train faster?
- **Network architecture:** The network above has two layers of trainable parameters. Can you do better with a deep network?
Good architectures to try include:
 - [conv-relu-pool]xN -> [affine]xM -> [softmax or SVM]
 - [conv-relu-conv-relu-pool]xN -> [affine]xM -> [softmax or SVM]
 - [batchnorm-relu-conv]xN -> [affine]xM -> [softmax or SVM]
- **Global Average Pooling:** Instead of flattening and then having multiple affine layers, perform convolutions until your image gets small (7x7 or so) and then perform an average pooling operation to get to a 1x1 image picture (1, 1, Filter#), which is then reshaped into a (Filter#) vector. This is used in [Google's Inception Network](#) (See Table 1 for their architecture).
- **Regularization:** Add l2 weight regularization, or perhaps use Dropout.

Tips for training

For each network architecture that you try, you should tune the learning rate and other hyperparameters. When doing this there are a couple important things to keep in mind:

- If the parameters are working well, you should see improvement within a few hundred iterations
- Remember the coarse-to-fine approach for hyperparameter tuning: start by testing a large range of hyperparameters for just a few training iterations to find the combinations of parameters that are working at all.
- Once you have found some sets of parameters that seem to work, search more finely around these parameters. You may need to train for more epochs.
- You should use the validation set for hyperparameter search, and save your test set for evaluating your architecture on the best parameters as selected by the validation set.

Going above and beyond

If you are feeling adventurous there are many other features you can implement to try and improve your performance. You are **not required** to implement any of these, but don't miss the fun if you have time!

- Alternative optimizers: you can try Adam, Adagrad, RMSprop, etc.
- Alternative activation functions such as leaky ReLU, parametric ReLU, ELU, or MaxOut.
- Model ensembles
- Data augmentation
- New Architectures
 - [ResNets](#) where the input from the previous layer is added to the output.
 - [DenseNets](#) where inputs into previous layers are concatenated together.
 - [This blog has an in-depth overview](#)

Have fun and happy training!

In [24]:

```
#####  
# TODO: #  
# Experiment with any architectures, optimizers, and hyperparameters. #  
# Achieve AT LEAST 70% accuracy on the *validation set* within 10 epochs. #  
# #  
# Note that you can use the check_accuracy function to evaluate on either #  
# the test set or the validation set, by passing either loader_test or #  
# loader_val as the second argument to check_accuracy. You should not touch #  
# the test set until you have finished your architecture and hyperparameter #  
# tuning, and only run the test set once at the end to report a final value. #
```

```
#####
model = None
optimizer = None

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

#  $(32 + 2 * 2 - 5) / 1 + 1 = 32$ 
layer1 = nn.Sequential(
    nn.Conv2d(3, 16, kernel_size=5, padding=2),
    nn.ReLU(),
)

#  $(32 + 2 * 2 - 5) / 1 + 1 = 32$ 
#  $(32 - 2) / 2 + 1 = 16$ 
layer2 = nn.Sequential(
    nn.Conv2d(16, 32, kernel_size=3, padding=1),
    nn.ReLU(),
    nn.MaxPool2d(2)
)

#  $(16 + 2 * 1 - 3) / 1 + 1 = 16$ 
layer3 = nn.Sequential(
    nn.Conv2d(32, 48, kernel_size=3, padding=1),
    nn.ReLU(),
)

#  $(16 + 2 * 1 - 3) / 1 + 1 = 16$ 
#  $(16 - 2) / 2 + 1 = 8$ 
layer4 = nn.Sequential(
    nn.Conv2d(48, 64, kernel_size=3, padding=1),
    nn.ReLU(),
    nn.MaxPool2d(2)
)

fc = nn.Linear(64*8*8, 10)

model = nn.Sequential(
    layer1,
    layer2,
    layer3,
    layer4,
    Flatten(),
    fc
)

learning_rate = 1e-3

optimizer = optim.Adam(model.parameters(), lr=learning_rate)

# Print training status every epoch
print_every = 1000

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

# You should get at least 70% accuracy
train_part34(model, optimizer, epochs=10)
```

```
Iteration 0, loss = 2.2974
Checking accuracy on validation set
Got 98 / 1000 correct (9.80)
```

```
Iteration 0, loss = 0.9447
Checking accuracy on validation set
Got 585 / 1000 correct (58.50)
```

```
Iteration 0, loss = 1.0508
Checking accuracy on validation set
Got 684 / 1000 correct (68.40)
```

```
Iteration 0, loss = 0.7837
Checking accuracy on validation set
Got 716 / 1000 correct (71.60)
```

```
Iteration 0, loss = 0.5664
Checking accuracy on validation set
Got 750 / 1000 correct (75.00)
```

```
Iteration 0, loss = 0.5874
Checking accuracy on validation set
Got 737 / 1000 correct (73.70)
```

```
Iteration 0, loss = 0.5890
Checking accuracy on validation set
Got 740 / 1000 correct (74.00)
```

```
Iteration 0, loss = 0.3552
Checking accuracy on validation set
Got 743 / 1000 correct (74.30)
```

```
Iteration 0, loss = 0.3969
Checking accuracy on validation set
Got 753 / 1000 correct (75.30)
```

```
Iteration 0, loss = 0.2919
Checking accuracy on validation set
Got 744 / 1000 correct (74.40)
```

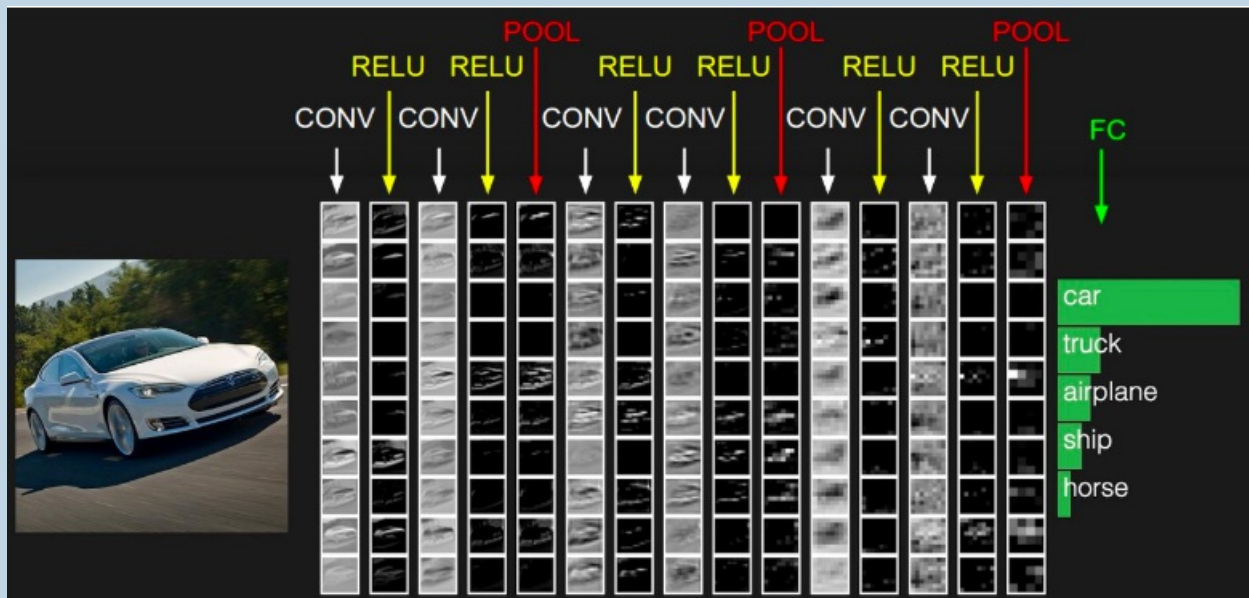
Describe what you did

In the cell below you should write an explanation of what you did, any additional features that you implemented, and/or any graphs that you made in the process of training and evaluating your network.

TODO:

I constructed a 5-layer convolutional network, the first four layers have similar structure as the example in the class slide, and the last layer is a fully-connected layer. They can be represented as below.

[(conv -> relu) -> (conv -> relu -> pool)] * 2 -> fc



Test set -- run this only once

Now that we've gotten a result we're happy with, we test our final model on the test set (which you should store in `best_model`). Think about how this compares to your validation set accuracy.

In [25]:

```
best_model = model
check_accuracy(test_data_loader, best_model)
```

```
check_accuracy_part34(loader_test, best_model)
```

```
Checking accuracy on test set  
Got 7327 / 10000 correct (73.27)
```

What's this TensorFlow business?

You've written a lot of code in this assignment to provide a whole host of neural network functionality. Dropout, Batch Norm, and 2D convolutions are some of the workhorses of deep learning in computer vision. You've also worked hard to make your code efficient and vectorized.

For the last part of this assignment, though, we're going to leave behind your beautiful codebase and instead migrate to one of two popular deep learning frameworks: in this instance, TensorFlow (or PyTorch, if you choose to work with that notebook).

Part I: Preparation

First, we load the CIFAR-10 dataset. This might take a few minutes to download the first time you run it, but after that the files should be cached on disk and loading should be faster.

In previous parts of the assignment we used CS231N-specific code to download and read the CIFAR-10 dataset; however the `tf.keras.datasets` package in TensorFlow provides prebuilt utility functions for loading many common datasets.

For the purposes of this assignment we will still write our own code to preprocess the data and iterate through it in minibatches. The `tf.data` package in TensorFlow provides tools for automating this process, but working with this package adds extra complication and is beyond the scope of this notebook. However using `tf.data` can be much more efficient than the simple approach used in this notebook, so you should consider using it for your project.

In []:

```
# We can iterate through a dataset like this:
for t, (x, y) in enumerate(train_dset):
    print(t, x.shape, y.shape)
    if t > 5: break
```

You can optionally **use GPU by setting the flag to True below**. It's not necessary to use a GPU for this assignment; if you are working on Google Cloud then we recommend that you do not use a GPU, as it will be significantly more expensive.

Barebones TensorFlow: Define a Two-Layer Network

We will now implement our first neural network with TensorFlow: a fully-connected ReLU network with two hidden layers and no biases on the CIFAR10 dataset. For now we will use only low-level TensorFlow operators to define the network; later we will see how to use the higher-level abstractions provided by `tf.keras` to simplify the process.

We will define the forward pass of the network in the function `two_layer_fc`; this will accept TensorFlow Tensors for the inputs and weights of the network, and return a TensorFlow Tensor for the scores.

After defining the network architecture in the `two_layer_fc` function, we will test the implementation by checking the shape of the output.

It's important that you read and understand this implementation.

Barebones TensorFlow: Three-Layer ConvNet

Here you will complete the implementation of the function `three_layer_convnet` which will perform the forward pass of a three-layer convolutional network. The network should have the following architecture:

1. A convolutional layer (with bias) with `channel_1` filters, each with shape `KW1 x KH1`, and zero-padding of two
2. ReLU nonlinearity
3. A convolutional layer (with bias) with `channel_2` filters, each with shape `KW2 x KH2`, and zero-padding of one
4. ReLU nonlinearity
5. Fully-connected layer with bias, producing scores for `C` classes.

HINT: For convolutions: https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/nn/conv2d; be careful with padding!

HINT: For biases: <https://www.tensorflow.org/performance/xla/broadcasting>

In []:

```
def three_layer_convnet(x, params):
    """
    A three-layer convolutional network with the architecture described above.

    Inputs:
    - x: A TensorFlow Tensor of shape (N, H, W, 3) giving a minibatch of images
    - params: A list of TensorFlow Tensors giving the weights and biases for the
      network; should contain the following:
      - conv_w1: TensorFlow Tensor of shape (KH1, KW1, 3, channel_1) giving
        weights for the first convolutional layer.
      - conv_b1: TensorFlow Tensor of shape (channel_1,) giving biases for the
        first convolutional layer.
      - conv_w2: TensorFlow Tensor of shape (KH2, KW2, channel_1, channel_2)
        giving weights for the second convolutional layer
      - conv_b2: TensorFlow Tensor of shape (channel_2,) giving biases for the
        second convolutional layer.
      - fc_w: TensorFlow Tensor giving weights for the fully-connected layer.
        Can you figure out what the shape should be?
      - fc_b: TensorFlow Tensor giving biases for the fully-connected layer.
        Can you figure out what the shape should be?
    """
    conv_w1, conv_b1, conv_w2, conv_b2, fc_w, fc_b = params
    scores = None
    #####
    # TODO: Implement the forward pass for the three-layer ConvNet. #
    #####
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    pass

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

After defining the forward pass of the three-layer ConvNet above, run the following cell to test your implementation. Like the two-layer network, we run the graph on a batch of zeros just to make sure the function doesn't crash, and produces outputs of the correct shape.

When you run this function, `scores_np` should have shape `(64, 10)`.

Barebones TensorFlow: Training Step

We now define the `training_step` function performs a single training step. This will take three basic steps:

1. Compute the loss
2. Compute the gradient of the loss with respect to all network weights
3. Make a weight update step using (stochastic) gradient descent.

We need to use a few new TensorFlow functions to do all of this:

- For computing the cross-entropy loss we'll use `tf.nn.sparse_softmax_cross_entropy_with_logits`:
https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/nn/sparse_softmax_cross_entropy_with_logits
- For averaging the loss across a minibatch of data we'll use `tf.reduce_mean`:
https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/reduce_mean
- For computing gradients of the loss with respect to the weights we'll use `tf.GradientTape` (useful for Eager execution):
https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/GradientTape
- We'll mutate the weight values stored in a TensorFlow Tensor using `tf.assign_sub` ("sub" is for subtraction):
https://www.tensorflow.org/api_docs/python/tf/assign_sub

Barebones TensorFlow: Initialization

We'll use the following utility method to initialize the weight matrices for our models using Kaiming's normalization method.

[1] He et al, *Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification*, ICCV 2015, <https://arxiv.org/abs/1502.01852>

In []:

```
def create_matrix_with_kaiming_normal(shape):
    if len(shape) == 2:
        fan_in, fan_out = shape[0], shape[1]
    elif len(shape) == 4:
        fan_in, fan_out = np.prod(shape[:3]), shape[3]
    return tf.keras.backend.random_normal(shape) * np.sqrt(2.0 / fan_in)
```

Barebones TensorFlow: Train a Two-Layer Network

We are finally ready to use all of the pieces defined above to train a two-layer fully-connected network on CIFAR-10.

We just need to define a function to initialize the weights of the model, and call `train_part2`.

Defining the weights of the network introduces another important piece of TensorFlow API: `tf.Variable`. A TensorFlow Variable is a Tensor whose value is stored in the graph and persists across runs of the computational graph; however unlike constants defined with `tf.zeros` or `tf.random_normal`, the values of a Variable can be mutated as the graph runs; these mutations will persist across graph runs. Learnable parameters of the network are usually stored in Variables.

You don't need to tune any hyperparameters, but you should achieve validation accuracies above 40% after one epoch of training.

In []:

```
def two_layer_fc_init():
    """
    Initialize the weights of a two-layer network, for use with the
    two_layer_network function defined above.
    You can use the `create_matrix_with_kaiming_normal` helper!

    Inputs: None

    Returns: A list of:
    - w1: TensorFlow tf.Variable giving the weights for the first layer
    - w2: TensorFlow tf.Variable giving the weights for the second layer
    """
    hidden_layer_size = 4000
    w1 = tf.Variable(create_matrix_with_kaiming_normal((3 * 32 * 32, 4000)))
    w2 = tf.Variable(create_matrix_with_kaiming_normal((4000, 10)))
    return [w1, w2]

learning_rate = 1e-2
train_part2(two_layer_fc, two_layer_fc_init, learning_rate)
```

Barebones TensorFlow: Train a three-layer ConvNet

We will now use TensorFlow to train a three-layer ConvNet on CIFAR-10.

You need to implement the `three_layer_convnet_init` function. Recall that the architecture of the network is:

1. Convolutional layer (with bias) with 32 5x5 filters, with zero-padding 2
2. ReLU
3. Convolutional layer (with bias) with 16 3x3 filters, with zero-padding 1
4. ReLU
5. Fully-connected layer (with bias) to compute scores for 10 classes

You don't need to do any hyperparameter tuning, but you should see validation accuracies above 43% after one epoch of training.

In []:

```
def three_layer_convnet_init():
    """
    Initialize the weights of a Three-Layer ConvNet, for use with the
    three_layer_convnet function defined above.
    You can use the `create_matrix_with_kaiming_normal` helper!

    Inputs: None
```

```

Returns a list containing:
- conv_w1: TensorFlow tf.Variable giving weights for the first conv layer
- conv_b1: TensorFlow tf.Variable giving biases for the first conv layer
- conv_w2: TensorFlow tf.Variable giving weights for the second conv layer
- conv_b2: TensorFlow tf.Variable giving biases for the second conv layer
- fc_w: TensorFlow tf.Variable giving weights for the fully-connected layer
- fc_b: TensorFlow tf.Variable giving biases for the fully-connected layer
"""

params = None
#####
# TODO: Initialize the parameters of the three-layer network. #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

pass

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

learning_rate = 3e-3
train_part2(three_layer_convnet, three_layer_convnet_init, learning_rate)

```

Keras Model Subclassing API: Three-Layer ConvNet

Now it's your turn to implement a three-layer ConvNet using the `tf.keras.Model` API. Your model should have the same architecture used in Part II:

1. Convolutional layer with 5 x 5 kernels, with zero-padding of 2
2. ReLU nonlinearity
3. Convolutional layer with 3 x 3 kernels, with zero-padding of 1
4. ReLU nonlinearity
5. Fully-connected layer to give class scores
6. Softmax nonlinearity

You should initialize the weights of your network using the same initialization method as was used in the two-layer network above.

Hint: Refer to the documentation for `tf.keras.layers.Conv2D` and `tf.keras.layers.Dense`:

https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/keras/layers/Conv2D

https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/keras/layers/Dense

In []:

```

class ThreeLayerConvNet(tf.keras.Model):
    def __init__(self, channel_1, channel_2, num_classes):
        super(ThreeLayerConvNet, self).__init__()
        #####
        # TODO: Implement the __init__ method for a three-layer ConvNet. You #
        # should instantiate layer objects to be used in the forward pass. #
        #####
        # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

        pass

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

    def call(self, x, training=False):
        scores = None
        #####
        # TODO: Implement the forward pass for a three-layer ConvNet. You #
        # should use the layer objects defined in the __init__ method. #
        #####
        # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

        pass

```

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

Once you complete the implementation of the `ThreeLayerConvNet` above you can run the following to ensure that your implementation does not crash and produces outputs of the expected shape.

In []:

```
def test_ThreeLayerConvNet():
    channel_1, channel_2, num_classes = 12, 8, 10
    model = ThreeLayerConvNet(channel_1, channel_2, num_classes)
    with tf.device(device):
        x = tf.zeros((64, 3, 32, 32))
        scores = model(x)
        print(scores.shape)

test_ThreeLayerConvNet()
```

Keras Model Subclassing API: Eager Training

While keras models have a builtin training loop (using the `model.fit`), sometimes you need more customization. Here's an example, of a training loop implemented with eager execution.

In particular, notice `tf.GradientTape`. Automatic differentiation is used in the backend for implementing backpropagation in frameworks like TensorFlow. During eager execution, `tf.GradientTape` is used to trace operations for computing gradients later. A particular `tf.GradientTape` can only compute one gradient; subsequent calls to `tape` will throw a runtime error.

TensorFlow 2.0 ships with easy-to-use built-in metrics under `tf.keras.metrics` module. Each metric is an object, and we can use `update_state()` to add observations and `reset_state()` to clear all observations. We can get the current result of a metric by calling `result()` on the metric object.

Keras Model Subclassing API: Train a Two-Layer Network

We can now use the tools defined above to train a two-layer network on CIFAR-10. We define the `model_init_fn` and `optimizer_init_fn` that construct the model and optimizer respectively when called. Here we want to train the model using stochastic gradient descent with no momentum, so we construct a `tf.keras.optimizers.SGD` function; you can [read about it here](#).

You don't need to tune any hyperparameters here, but you should achieve validation accuracies above 40% after one epoch of training.

In []:

```
hidden_size, num_classes = 4000, 10
learning_rate = 1e-2

def model_init_fn():
    return TwoLayerFC(hidden_size, num_classes)

def optimizer_init_fn():
    return tf.keras.optimizers.SGD(learning_rate=learning_rate)

train_part34(model_init_fn, optimizer_init_fn)
```

Keras Model Subclassing API: Train a Three-Layer ConvNet

Here you should use the tools we've defined above to train a three-layer ConvNet on CIFAR-10. Your ConvNet should use 32 filters in the first convolutional layer and 16 filters in the second layer.

To train the model you should use gradient descent with Nesterov momentum 0.9.

HINT: https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/optimizers/SGD

You don't need to perform any hyperparameter tuning, but you should achieve validation accuracies above 50% after training for one epoch

epoch.

In []:

```
learning_rate = 3e-3
channel_1, channel_2, num_classes = 32, 16, 10

def model_init_fn():
    model = None
    #####
    # TODO: Complete the implementation of model_fn. #
    #####
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    pass

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

def optimizer_init_fn():
    optimizer = None
    #####
    # TODO: Complete the implementation of model_fn. #
    #####
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    pass

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

train_part34(model_init_fn, optimizer_init_fn)
```

Part IV: Keras Sequential API

In Part III we introduced the `tf.keras.Model` API, which allows you to define models with any number of learnable layers and with arbitrary connectivity between layers.

However for many models you don't need such flexibility - a lot of models can be expressed as a sequential stack of layers, with the output of each layer fed to the next layer as input. If your model fits this pattern, then there is an even easier way to define your model: using `tf.keras.Sequential`. You don't need to write any custom classes; you simply call the `tf.keras.Sequential` constructor with a list containing a sequence of layer objects.

One complication with `tf.keras.Sequential` is that you must define the shape of the input to the model by passing a value to the `input_shape` of the first layer in your model.

Keras Sequential API: Two-Layer Network

In this subsection, we will rewrite the two-layer fully-connected network using `tf.keras.Sequential`, and train it using the training loop defined above.

You don't need to perform any hyperparameter tuning here, but you should see validation accuracies above 40% after training for one epoch.

In []:

```
learning_rate = 1e-2

def model_init_fn():
    input_shape = (32, 32, 3)
    hidden_layer_size, num_classes = 4000, 10
    initializer = tf.initializers.VarianceScaling(scale=2.0)
    layers = [
        tf.keras.layers.Flatten(input_shape=input_shape),
        tf.keras.layers.Dense(hidden_layer_size, activation='relu',
                               kernel_initializer=initializer),
        tf.keras.layers.Dense(num_classes, activation='softmax')
```

```

        tf.keras.layers.Dense(num_classes, activation='softmax',
                                kernel_initializer=initializer),
    ]
    model = tf.keras.Sequential(layers)
    return model

def optimizer_init_fn():
    return tf.keras.optimizers.SGD(learning_rate=learning_rate)

train_part34(model_init_fn, optimizer_init_fn)

```

Abstracting Away the Training Loop

In the previous examples, we used a customised training loop to train models (e.g. `train_part34`). Writing your own training loop is only required if you need more flexibility and control during training your model. Alternately, you can also use built-in APIs like `tf.keras.Model.fit()` and `tf.keras.Model.evaluate` to train and evaluate a model. Also remember to configure your model for training by calling `tf.keras.Model.compile`.

You don't need to perform any hyperparameter tuning here, but you should see validation and test accuracies above 42% after training for one epoch.

In []:

```

model = model_init_fn()
model.compile(optimizer=tf.keras.optimizers.SGD(learning_rate=learning_rate),
              loss='sparse_categorical_crossentropy',
              metrics=[tf.keras.metrics.sparse_categorical_accuracy])
model.fit(X_train, y_train, batch_size=64, epochs=1, validation_data=(X_val, y_val))
model.evaluate(X_test, y_test)

```

Keras Sequential API: Three-Layer ConvNet

Here you should use `tf.keras.Sequential` to reimplement the same three-layer ConvNet architecture used in Part II and Part III. As a reminder, your model should have the following architecture:

1. Convolutional layer with 32 5x5 kernels, using zero padding of 2
2. ReLU nonlinearity
3. Convolutional layer with 16 3x3 kernels, using zero padding of 1
4. ReLU nonlinearity
5. Fully-connected layer giving class scores
6. Softmax nonlinearity

You should initialize the weights of the model using a `tf.initializers.VarianceScaling` as above.

You should train the model using Nesterov momentum 0.9.

You don't need to perform any hyperparameter search, but you should achieve accuracy above 45% after training for one epoch.

In []:

```

def model_init_fn():
    model = None
    #####
    # TODO: Construct a three-layer ConvNet using tf.keras.Sequential. #
    #####
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    pass

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

learning_rate = 5e-4
def optimizer_init_fn():
    optimizer = None
    #####
    # TODO: Complete the implementation of model_fn. #
    #####
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

```

```

pass

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

train_part34(model_init_fn, optimizer_init_fn)

```

We will also train this model with the built-in training loop APIs provided by TensorFlow.

In []:

```

model = model_init_fn()
model.compile(optimizer='sgd',
              loss='sparse_categorical_crossentropy',
              metrics=[tf.keras.metrics.sparse_categorical_accuracy])
model.fit(X_train, y_train, batch_size=64, epochs=1, validation_data=(X_val, y_val))
model.evaluate(X_test, y_test)

```

Part IV: Functional API

Demonstration with a Two-Layer Network

In the previous section, we saw how we can use `tf.keras.Sequential` to stack layers to quickly build simple models. But this comes at the cost of losing flexibility.

Often we will have to write complex models that have non-sequential data flows: a layer can have **multiple inputs and/or outputs**, such as stacking the output of 2 previous layers together to feed as input to a third! (Some examples are residual connections and dense blocks.)

In such cases, we can use Keras functional API to write models with complex topologies such as:

1. Multi-input models
2. Multi-output models
3. Models with shared layers (the same layer called several times)
4. Models with non-sequential data flows (e.g. residual connections)

Writing a model with Functional API requires us to create a `tf.keras.Model` instance and explicitly write input tensors and output tensors for this model.

Keras Functional API: Train a Two-Layer Network

You can now train this two-layer network constructed using the functional API.

You don't need to perform any hyperparameter tuning here, but you should see validation accuracies above 40% after training for one epoch.

In []:

```

input_shape = (32, 32, 3)
hidden_size, num_classes = 4000, 10
learning_rate = 1e-2

def model_init_fn():
    return two_layer_fc_functional(input_shape, hidden_size, num_classes)

def optimizer_init_fn():
    return tf.keras.optimizers.SGD(learning_rate=learning_rate)

train_part34(model_init_fn, optimizer_init_fn)

```

Part V: CIFAR-10 open-ended challenge

In this section you can experiment with whatever ConvNet architecture you'd like on CIFAR-10.

You should experiment with architectures, hyperparameters, loss functions, regularization, or anything else you can think of to train a

You should experiment with architecture, hyperparameters, loss functions, regularization, or anything else you can think of to train a model that achieves **at least 70% accuracy** on the **validation** set within 10 epochs. You can use the built-in train function, the `train_part34` function from above, or implement your own training loop.

Describe what you did at the end of the notebook.

Some things you can try:

- **Filter size:** Above we used 5x5 and 3x3; is this optimal?
- **Number of filters:** Above we used 16 and 32 filters. Would more or fewer do better?
- **Pooling:** We didn't use any pooling above. Would this improve the model?
- **Normalization:** Would your model be improved with batch normalization, layer normalization, group normalization, or some other normalization strategy?
- **Network architecture:** The ConvNet above has only three layers of trainable parameters. Would a deeper model do better?
- **Global average pooling:** Instead of flattening after the final convolutional layer, would global average pooling do better? This strategy is used for example in Google's Inception network and in Residual Networks.
- **Regularization:** Would some kind of regularization improve performance? Maybe weight decay or dropout?

NOTE: Batch Normalization / Dropout

If you are using Batch Normalization and Dropout, remember to pass `is_training=True` if you use the `train_part34()` function. BatchNorm and Dropout layers have different behaviors at training and inference time. `training` is a specific keyword argument reserved for this purpose in any `tf.keras.Model`'s `call()` function. Read more about this here :

https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/keras/layers/BatchNormalization#methods

https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/keras/layers/Dropout#methods

Tips for training

For each network architecture that you try, you should tune the learning rate and other hyperparameters. When doing this there are a couple important things to keep in mind:

- If the parameters are working well, you should see improvement within a few hundred iterations
- Remember the coarse-to-fine approach for hyperparameter tuning: start by testing a large range of hyperparameters for just a few training iterations to find the combinations of parameters that are working at all.
- Once you have found some sets of parameters that seem to work, search more finely around these parameters. You may need to train for more epochs.
- You should use the validation set for hyperparameter search, and save your test set for evaluating your architecture on the best parameters as selected by the validation set.

Going above and beyond

If you are feeling adventurous there are many other features you can implement to try and improve your performance. You are **not required** to implement any of these, but don't miss the fun if you have time!

- Alternative optimizers: you can try Adam, Adagrad, RMSprop, etc.
- Alternative activation functions such as leaky ReLU, parametric ReLU, ELU, or MaxOut.
- Model ensembles
- Data augmentation
- New Architectures
 - [ResNets](#) where the input from the previous layer is added to the output.
 - [DenseNets](#) where inputs into previous layers are concatenated together.
 - [This blog has an in-depth overview](#)

Have fun and happy training!

In []:

```
class CustomConvNet(tf.keras.Model):
    def __init__(self):
        super(CustomConvNet, self).__init__()
        #####
        # TODO: Construct a model that performs well on CIFAR-10 #
        #####
        # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

        pass

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

```

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

def call(self, input_tensor, training=False):
    #####
    # TODO: Construct a model that performs well on CIFAR-10                        #
    #####
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    pass

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

    return x

device = '/device:GPU:0'   # Change this to a CPU/GPU as you wish!
# device = '/cpu:0'       # Change this to a CPU/GPU as you wish!
print_every = 700
num_epochs = 10

model = CustomConvNet()

def model_init_fn():
    return CustomConvNet()

def optimizer_init_fn():
    learning_rate = 1e-3
    return tf.keras.optimizers.Adam(learning_rate)

train_part34(model_init_fn, optimizer_init_fn, num_epochs=num_epochs, is_training=True)

```

Describe what you did

In the cell below you should write an explanation of what you did, any additional features that you implemented, and/or any graphs that you made in the process of training and evaluating your network.

TODO: Tell us what you did

1 layers.py

```
1 from builtins import range
2 import numpy as np
3
4
5 def affine_forward(x, w, b):
6     """
7     Computes the forward pass for an affine (fully-connected) layer.
8
9     The input x has shape (N, d_1, ..., d_k) and contains a minibatch of N
10    examples, where each example x[i] has shape (d_1, ..., d_k). We will
11    reshape each input into a vector of dimension  $\bar{D} = d_1 * \dots * d_k$ , and
12    then transform it to an output vector of dimension M.
13
14    Inputs:
15    - x: A numpy array containing input data, of shape (N, d_1, ..., d_k)
16    - w: A numpy array of weights, of shape (D, M)
17    - b: A numpy array of biases, of shape (M,)
18
19    Returns a tuple of:
20    - out: output, of shape (N, M)
21    - cache: (x, w, b)
22    """
23    out = None
24    #####
25    # TODO: Implement the affine forward pass. Store the result in out. You #
26    # will need to reshape the input into rows. #
27    #####
28    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
29
30    N = x.shape[0]
31    x_row = np.reshape(x, (N, -1)) # Reshape x into (N, D) matrix
32    out = x_row.dot(w) + b
33
34    # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
35    #####
36    #                               END OF YOUR CODE                               #
37    #####
38    cache = (x, w, b)
39    return out, cache
40
41
42 def affine_backward(dout, cache):
43     """
44     Computes the backward pass for an affine layer.
45
46     Inputs:
47     - dout: Upstream derivative, of shape (N, M)
48     - cache: Tuple of:
49       - x: Input data, of shape (N, d_1, ..., d_k)
50       - w: Weights, of shape (D, M)
51       - b: Biases, of shape (M,)
52
53     Returns a tuple of:
54     - dx: Gradient with respect to x, of shape (N, d_1, ..., d_k)
55     - dw: Gradient with respect to w, of shape (D, M)
56     - db: Gradient with respect to b, of shape (M,)
57     """
58    x, w, b = cache
59    dx, dw, db = None, None, None
60    #####
61    # TODO: Implement the affine backward pass. #
62    #####
63    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
64
65    N = x.shape[0]
66
67    x_row = np.reshape(x, (N, -1))
68
69    # Based on the shape of dx, dw and db, can get the calculation formula
70    dx = dout.dot(w.T).reshape(x.shape)
71    dw = x_row.T.dot(dout)
72    db = np.sum(dout, axis = 0)
73
```

```

74 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
75 #####
76 #                                     END OF YOUR CODE                                     #
77 #####
78 return dx, dw, db
79
80
81 def relu_forward(x):
82     """
83     Computes the forward pass for a layer of rectified linear units (ReLU).
84
85     Input:
86     - x: Inputs, of any shape
87
88     Returns a tuple of:
89     - out: Output, of the same shape as x
90     - cache: x
91     """
92     out = None
93     #####
94     # TODO: Implement the ReLU forward pass.                                     #
95     #####
96     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
97
98     # Must truly copy the variables into new variables
99     out = x.copy()
100     out[out < 0] = 0
101
102
103     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
104     #####
105     #                                     END OF YOUR CODE                                     #
106     #####
107     cache = x
108     return out, cache
109
110
111 def relu_backward(dout, cache):
112     """
113     Computes the backward pass for a layer of rectified linear units (ReLU).
114
115     Input:
116     - dout: Upstream derivatives, of any shape
117     - cache: Input x, of same shape as dout
118
119     Returns:
120     - dx: Gradient with respect to x
121     """
122     dx, x = None, cache
123     #####
124     # TODO: Implement the ReLU backward pass.                                     #
125     #####
126     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
127
128     dx = x
129     dx[dx < 0] = 0
130     dx[dx > 0] = 1
131     dx = np.multiply(dx, dout)
132
133
134     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
135     #####
136     #                                     END OF YOUR CODE                                     #
137     #####
138     return dx
139
140
141 def batchnorm_forward(x, gamma, beta, bn_param):
142     """
143     Forward pass for batch normalization.
144
145     During training the sample mean and (uncorrected) sample variance are
146     computed from minibatch statistics and used to normalize the incoming data.
147     During training we also keep an exponentially decaying running mean of the
148     mean and variance of each feature, and these averages are used to normalize
149     data at test-time.

```

```

150
151 At each timestep we update the running averages for mean and variance using
152 an exponential decay based on the momentum parameter:
153
154 running_mean = momentum * running_mean + (1 - momentum) * sample_mean
155 running_var = momentum * running_var + (1 - momentum) * sample_var
156
157 Note that the batch normalization paper suggests a different test-time
158 behavior: they compute sample mean and variance for each feature using a
159 large number of training images rather than using a running average. For
160 this implementation we have chosen to use running averages instead since
161 they do not require an additional estimation step; the torch7
162 implementation of batch normalization also uses running averages.
163
164 Input:
165 - x: Data of shape (N, D)
166 - gamma: Scale parameter of shape (D,)
167 - beta: Shift parameter of shape (D,)
168 - bn_param: Dictionary with the following keys:
169   - mode: 'train' or 'test'; required
170   - eps: Constant for numeric stability
171   - momentum: Constant for running mean / variance.
172   - running_mean: Array of shape (D,) giving running mean of features
173   - running_var: Array of shape (D,) giving running variance of features
174
175 Returns a tuple of:
176 - out: of shape (N, D)
177 - cache: A tuple of values needed in the backward pass
178 """
179 mode = bn_param['mode']
180 eps = bn_param.get('eps', 1e-5)
181 momentum = bn_param.get('momentum', 0.9)
182
183 N, D = x.shape
184 running_mean = bn_param.get('running_mean', np.zeros(D, dtype=x.dtype))
185 running_var = bn_param.get('running_var', np.zeros(D, dtype=x.dtype))
186
187 out, cache = None, None
188 if mode == 'train':
189     #####
190     # TODO: Implement the training-time forward pass for batch norm.      #
191     # Use minibatch statistics to compute the mean and variance, use      #
192     # these statistics to normalize the incoming data, and scale and      #
193     # shift the normalized data using gamma and beta.                      #
194     #                                                                      #
195     # You should store the output in the variable out. Any intermediates  #
196     # that you need for the backward pass should be stored in the cache   #
197     # variable.                                                            #
198     #                                                                      #
199     # You should also use your computed sample mean and variance together #
200     # with the momentum variable to update the running mean and running  #
201     # variance, storing your result in the running_mean and running_var  #
202     # variables.                                                           #
203     #                                                                      #
204     # Note that though you should be keeping track of the running         #
205     # variance, you should normalize the data based on the standard       #
206     # deviation (square root of variance) instead!                       #
207     # Referencing the original paper (https://arxiv.org/abs/1502.03167)  #
208     # might prove to be helpful.                                         #
209     #####
210     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
211
212     sample_mean = np.mean(x, axis = 0)
213     sample_var = np.var(x, axis = 0)
214     running_mean = momentum * running_mean + (1 - momentum) * sample_mean
215     running_var = momentum * running_var + (1 - momentum) * sample_var
216     sample_normalized = (x - sample_mean) / np.sqrt(sample_var + eps)
217     out = gamma * sample_normalized + beta
218     cache = (sample_normalized, gamma, beta, sample_mean, sample_var, x, eps)
219
220     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
221     #####
222     #                               END OF YOUR CODE                       #
223     #####
224 elif mode == 'test':
225     #####

```

```

226 # TODO: Implement the test-time forward pass for batch normalization. #
227 # Use the running mean and variance to normalize the incoming data, #
228 # then scale and shift the normalized data using gamma and beta. #
229 # Store the result in the out variable. #
230 #####
231 # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
232
233 sample_normalized = (x - running_mean) / np.sqrt(running_var + eps)
234 out = gamma * sample_normalized + beta
235
236 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
237 #####
238 #                               END OF YOUR CODE                               #
239 #####
240 else:
241     raise ValueError('Invalid forward batchnorm mode "%s"' % mode)
242
243 # Store the updated running means back into bn_param
244 bn_param['running_mean'] = running_mean
245 bn_param['running_var'] = running_var
246
247 return out, cache
248
249
250 def batchnorm_backward(dout, cache):
251     """
252     Backward pass for batch normalization.
253
254     For this implementation, you should write out a computation graph for
255     batch normalization on paper and propagate gradients backward through
256     intermediate nodes.
257
258     Inputs:
259     - dout: Upstream derivatives, of shape (N, D)
260     - cache: Variable of intermediates from batchnorm_forward.
261
262     Returns a tuple of:
263     - dx: Gradient with respect to inputs x, of shape (N, D)
264     - dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
265     - dbeta: Gradient with respect to shift parameter beta, of shape (D,)
266     """
267     dx, dgamma, dbeta = None, None, None
268     #####
269     # TODO: Implement the backward pass for batch normalization. Store the #
270     # results in the dx, dgamma, and dbeta variables. #
271     # Referencing the original paper (https://arxiv.org/abs/1502.03167) #
272     # might prove to be helpful. #
273     #####
274     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
275
276     sample_normalized, gamma, beta, sample_mean, sample_var, x, eps = cache
277     N = x.shape[0]
278     dx_hat = dout * gamma
279     dvar = np.sum(dx_hat * (x - sample_mean) * (-1 / 2) * (sample_var + eps)**(-3 / 2), axis = 0)
280     dmean = np.sum(np.divide(-dx_hat, np.sqrt(sample_var + eps)), axis = 0) + dvar * np.sum(-2 * (x -
281     sample_mean), axis = 0) / N
282     dx = dx_hat / np.sqrt(sample_var + eps) + dvar * 2 * (x - sample_mean) / N + dmean / N
283     dgamma = np.sum(dout * sample_normalized, axis = 0)
284     dbeta = np.sum(dout, axis = 0)
285
286     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
287     #####
288     #                               END OF YOUR CODE                               #
289     #####
290
291     return dx, dgamma, dbeta
292
293
294 def batchnorm_backward_alt(dout, cache):
295     """
296     Alternative backward pass for batch normalization.
297
298     For this implementation you should work out the derivatives for the batch
299     normalizaton backward pass on paper and simplify as much as possible. You
300     should be able to derive a simple expression for the backward pass.

```

```

301 See the jupyter notebook for more hints.
302
303 Note: This implementation should expect to receive the same cache variable
304 as batchnorm_backward, but might not use all of the values in the cache.
305
306 Inputs / outputs: Same as batchnorm_backward
307 """
308 dx, dgamma, dbeta = None, None, None
309 #####
310 # TODO: Implement the backward pass for batch normalization. Store the #
311 # results in the dx, dgamma, and dbeta variables. #
312 # #
313 # After computing the gradient with respect to the centered inputs, you #
314 # should be able to compute gradients with respect to the inputs in a #
315 # single statement; our implementation fits on a single 80-character line.#
316 #####
317 # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
318
319 sample_normalized, gamma, beta, sample_mean, sample_var, x, eps = cache
320 N = x.shape[0]
321 sigma = np.sqrt(sample_var + eps)
322
323 dgamma = np.sum(dout * sample_normalized, axis = 0)
324 dbeta = np.sum(dout, axis = 0)
325
326 dx = (1 / N) * gamma * 1/sigma * ((N * dout) - np.sum(dout, axis=0) -
327                                     (x - sample_mean) * np.square(1/sigma) * np.sum(dout * (x -
328 sample_mean), axis=0))
329
330
331 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
332 #####
333 # END OF YOUR CODE #
334 #####
335
336 return dx, dgamma, dbeta
337
338
339 def layernorm_forward(x, gamma, beta, ln_param):
340     """
341     Forward pass for layer normalization.
342
343     During both training and test-time, the incoming data is normalized per data-point,
344     before being scaled by gamma and beta parameters identical to that of batch normalization.
345
346     Note that in contrast to batch normalization, the behavior during train and test-time for
347     layer normalization are identical, and we do not need to keep track of running averages
348     of any sort.
349
350     Input:
351     - x: Data of shape (N, D)
352     - gamma: Scale parameter of shape (D,)
353     - beta: Shift parameter of shape (D,)
354     - ln_param: Dictionary with the following keys:
355         - eps: Constant for numeric stability
356
357     Returns a tuple of:
358     - out: of shape (N, D)
359     - cache: A tuple of values needed in the backward pass
360     """
361     out, cache = None, None
362     eps = ln_param.get('eps', 1e-5)
363     #####
364     # TODO: Implement the training-time forward pass for layer norm. #
365     # Normalize the incoming data, and scale and shift the normalized data #
366     # using gamma and beta. #
367     # HINT: this can be done by slightly modifying your training-time #
368     # implementation of batch normalization, and inserting a line or two of #
369     # well-placed code. In particular, can you think of any matrix #
370     # transformations you could perform, that would enable you to copy over #
371     # the batch norm code and leave it almost unchanged? #
372     #####
373     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
374
375     # Transpose x so that the dimension of x becomes (D, N), following calculation can be the same as

```

```

376 batch_forward
377 x = x.T
378
379 sample_mean = np.mean(x, axis = 0)
380 sample_var = np.var(x, axis = 0)
381 sample_normalized = (x - sample_mean) / np.sqrt(sample_var + eps)
382
383 # Transpose sample_normalized so that the result can has the correct dimension (N, D)
384 sample_normalized = sample_normalized.T
385 out = gamma * sample_normalized + beta
386
387 # Transpose x again so that x is restored
388 x = x.T
389
390 cache = (sample_normalized, gamma, beta, sample_mean, sample_var, x, eps)
391
392 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
393 #####
394 #                               END OF YOUR CODE                               #
395 #####
396 return out, cache
397
398
399 def layernorm_backward(dout, cache):
400     """
401     Backward pass for layer normalization.
402
403     For this implementation, you can heavily rely on the work you've done already
404     for batch normalization.
405
406     Inputs:
407     - dout: Upstream derivatives, of shape (N, D)
408     - cache: Variable of intermediates from layernorm_forward.
409
410     Returns a tuple of:
411     - dx: Gradient with respect to inputs x, of shape (N, D)
412     - dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
413     - dbeta: Gradient with respect to shift parameter beta, of shape (D,)
414     """
415     dx, dgamma, dbeta = None, None, None
416     #####
417     # TODO: Implement the backward pass for layer norm.
418     #
419     # HINT: this can be done by slightly modifying your training-time
420     # implementation of batch normalization. The hints to the forward pass
421     # still apply!
422     #####
423     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
424
425     sample_normalized, gamma, beta, sample_mean, sample_var, x, eps = cache
426
427     # The calculation of dgamma and dbeta remain the same
428     dgamma = np.sum(dout * sample_normalized, axis = 0)
429     dbeta = np.sum(dout, axis = 0)
430
431     dx_hat = dout * gamma
432
433     # At first transpose sample_normalized, x and dx_hat so that their dimensions are all (D, N) now
434     sample_normalized = sample_normalized.T
435     x = x.T
436     dx_hat = dx_hat.T
437
438     # Actually x.shape[0] should be D now, but I still use N so that the code below don't have to be
439     # changed.
440     N = x.shape[0]
441
442     # The following calculation can be the same as they are in batchnorm_backward
443     dvar = np.sum(dx_hat * (x - sample_mean) * (-1 / 2) * (sample_var + eps)**(-3 / 2), axis = 0)
444     dmean = np.sum(np.divide(-dx_hat, np.sqrt(sample_var + eps)), axis = 0) + dvar * np.sum(-2 * (x -
445     sample_mean), axis = 0) / N
446     dx = dx_hat / np.sqrt(sample_var + eps) + dvar * 2 * (x - sample_mean) / N + dmean / N
447
448     # Transpose dx so that dx can have the correct dimension (N, D) now
449     dx = dx.T

```

```

449
450
451
452 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
453 #####
454 #                                     END OF YOUR CODE                                     #
455 #####
456 return dx, dgamma, dbeta
457
458
459 def dropout_forward(x, dropout_param):
460     """
461     Performs the forward pass for (inverted) dropout.
462
463     Inputs:
464     - x: Input data, of any shape
465     - dropout_param: A dictionary with the following keys:
466       - p: Dropout parameter. We keep each neuron output with probability p.
467       - mode: 'test' or 'train'. If the mode is train, then perform dropout;
468         if the mode is test, then just return the input.
469       - seed: Seed for the random number generator. Passing seed makes this
470         function deterministic, which is needed for gradient checking but not
471         in real networks.
472
473     Outputs:
474     - out: Array of the same shape as x.
475     - cache: tuple (dropout_param, mask). In training mode, mask is the dropout
476       mask that was used to multiply the input; in test mode, mask is None.
477
478     NOTE: Please implement **inverted** dropout, not the vanilla version of dropout.
479     See http://cs231n.github.io/neural-networks-2/#reg for more details.
480
481     NOTE 2: Keep in mind that p is the probability of **keep** a neuron
482     output; this might be contrary to some sources, where it is referred to
483     as the probability of dropping a neuron output.
484     """
485     p, mode = dropout_param['p'], dropout_param['mode']
486     if 'seed' in dropout_param:
487         np.random.seed(dropout_param['seed'])
488
489     mask = None
490     out = None
491
492     if mode == 'train':
493         #####
494         # TODO: Implement training phase forward pass for inverted dropout. #
495         # Store the dropout mask in the mask variable. #
496         #####
497         # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
498
499         mask = (np.random.rand(*x.shape) < p) / p
500         out = x * mask
501
502         # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
503         #####
504         #                                     END OF YOUR CODE                                     #
505         #####
506     elif mode == 'test':
507         #####
508         # TODO: Implement the test phase forward pass for inverted dropout. #
509         #####
510         # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
511
512         out = x
513
514         # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
515         #####
516         #                                     END OF YOUR CODE                                     #
517         #####
518
519     cache = (dropout_param, mask)
520     out = out.astype(x.dtype, copy=False)
521
522     return out, cache
523
524

```

```

525 def dropout_backward(dout, cache):
526     """
527     Perform the backward pass for (inverted) dropout.
528
529     Inputs:
530     - dout: Upstream derivatives, of any shape
531     - cache: (dropout_param, mask) from dropout_forward.
532     """
533     dropout_param, mask = cache
534     mode = dropout_param['mode']
535
536     dx = None
537     if mode == 'train':
538         #####
539         # TODO: Implement training phase backward pass for inverted dropout #
540         #####
541         # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
542
543         dx = dout * mask
544
545         # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
546         #####
547         #                               END OF YOUR CODE                               #
548         #####
549     elif mode == 'test':
550         dx = dout
551     return dx
552
553 def conv_forward_naive(x, w, b, conv_param):
554     """
555     A naive implementation of the forward pass for a convolutional layer.
556
557     The input consists of N data points, each with C channels, height H and
558     width W. We convolve each input with F different filters, where each filter
559     spans all C channels and has height HH and width WW.
560
561     Input:
562     - x: Input data of shape (N, C, H, W)
563     - w: Filter weights of shape (F, C, HH, WW)
564     - b: Biases, of shape (F,)
565     - conv_param: A dictionary with the following keys:
566       - 'stride': The number of pixels between adjacent receptive fields in the
567         horizontal and vertical directions.
568       - 'pad': The number of pixels that will be used to zero-pad the input.
569
570
571     During padding, 'pad' zeros should be placed symmetrically (i.e equally on both sides)
572     along the height and width axes of the input. Be careful not to modify the original
573     input x directly.
574
575     Returns a tuple of:
576     - out: Output data, of shape (N, F, H', W') where H' and W' are given by
577        $H' = 1 + (H + 2 * pad - HH) / stride$ 
578        $W' = 1 + (W + 2 * pad - WW) / stride$ 
579     - cache: (x, w, b, conv_param)
580     """
581
582     out = None
583     #####
584     # TODO: Implement the convolutional forward pass. #
585     # Hint: you can use the function np.pad for padding. #
586     #####
587     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
588
589     stride = conv_param['stride']
590     pad = conv_param['pad']
591     x_pad = np.pad(x, ((0, 0), (0, 0), (pad, pad), (pad, pad)), 'constant')
592     N, C, H, W = x.shape
593     F, C, HH, WW = w.shape
594
595     H_out = int(1 + (H + 2 * pad - HH) / stride)
596     W_out = int(1 + (W + 2 * pad - WW) / stride)
597
598     out = np.zeros((N, F, H_out, W_out))
599
600     for n in range(N):

```



```

601         for f in range(F):
602             for i in range(H_out):
603                 for j in range(W_out):
604                     out[n, f, i, j] = np.sum(x_pad[n, :, i * stride: i * stride + HH, j * stride: j *
stride + WW] * w[f]) + b[f]
605
606     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
607     #####
608     #                                     END OF YOUR CODE                                     #
609     #####
610     cache = (x, w, b, conv_param)
611     return out, cache
612
613
614 def conv_backward_naive(dout, cache):
615     """
616     A naive implementation of the backward pass for a convolutional layer.
617
618     Inputs:
619     - dout: Upstream derivatives.
620     - cache: A tuple of (x, w, b, conv_param) as in conv_forward_naive
621
622     Returns a tuple of:
623     - dx: Gradient with respect to x
624     - dw: Gradient with respect to w
625     - db: Gradient with respect to b
626     """
627     dx, dw, db = None, None, None
628     #####
629     # TODO: Implement the convolutional backward pass.                                     #
630     #####
631     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
632
633     x, w, b, conv_param = cache
634     pad = conv_param['pad']
635     stride = conv_param['stride']
636     F, C, HH, WW = w.shape
637     N, C, H, W = x.shape
638     H_out = int(1 + (H + 2 * pad - HH) / stride)
639     W_out = int(1 + (W + 2 * pad - WW) / stride)
640     x_pad = np.pad(x, ((0, 0), (0, 0), (pad, pad), (pad, pad)), 'constant')
641
642     dx_pad = np.zeros_like(x_pad)
643     dw = np.zeros_like(w)
644     db = np.zeros_like(b)
645
646     # To calculate db, just sum up all the upstream gradients for each filters bias.
647     for f in range(F):
648         db[f] = np.sum(dout[:, f, :, :])
649
650     for n in range(N):
651         for f in range(F):
652             for i in range(H_out):
653                 for j in range(W_out):
654                     # According to chain rule, dw = dout * x, dx = dout * w. Be careful about the
dimension
655                     dw[f] += dout[n, f, i, j] * x_pad[n, :, i * stride: i * stride + HH, j * stride: j *
stride + WW]
656                     dx_pad[n, :, i * stride: i * stride + HH, j * stride: j * stride + WW] += dout[n, f, i
, j] * w[f]
657
658     # Get rid of the pad around dx
659     dx = dx_pad[:, :, pad: pad+H, pad: pad+W]
660
661     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
662     #####
663     #                                     END OF YOUR CODE                                     #
664     #####
665     return dx, dw, db
666
667
668 def max_pool_forward_naive(x, pool_param):
669     """
670     A naive implementation of the forward pass for a max-pooling layer.
671
672     Inputs:

```

```

673 - x: Input data, of shape (N, C, H, W)
674 - pool_param: dictionary with the following keys:
675   - 'pool_height': The height of each pooling region
676   - 'pool_width': The width of each pooling region
677   - 'stride': The distance between adjacent pooling regions
678
679 No padding is necessary here. Output size is given by
680
681 Returns a tuple of:
682 - out: Output data, of shape (N, C, H', W') where H' and W' are given by
683   H' = 1 + (H - pool_height) / stride
684   W' = 1 + (W - pool_width) / stride
685 - cache: (x, pool_param)
686 """
687 out = None
688 #####
689 # TODO: Implement the max-pooling forward pass #
690 #####
691 # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
692
693 pool_height = pool_param['pool_height']
694 pool_width = pool_param['pool_width']
695 stride = pool_param['stride']
696 N, C, H, W = x.shape
697
698 H_out = int (1 + (H - pool_height) / stride)
699 W_out = int (1 + (W - pool_width) / stride)
700
701 out = np.zeros((N, C, H_out, W_out))
702
703 for n in range(N):
704     for c in range(C):
705         for i in range(H_out):
706             for j in range(W_out):
707                 out[n, c, i, j] = np.max(x[n, c, i * stride: i * stride + pool_height, j * stride: j *
stride + pool_width])
708
709
710 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
711 #####
712 #                               END OF YOUR CODE                               #
713 #####
714 cache = (x, pool_param)
715 return out, cache
716
717
718 def max_pool_backward_naive(dout, cache):
719     """
720     A naive implementation of the backward pass for a max-pooling layer.
721
722     Inputs:
723     - dout: Upstream derivatives
724     - cache: A tuple of (x, pool_param) as in the forward pass.
725
726     Returns:
727     - dx: Gradient with respect to x
728     """
729     dx = None
730     #####
731     # TODO: Implement the max-pooling backward pass #
732     #####
733     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
734
735     x, pool_param = cache
736     pool_height = pool_param['pool_height']
737     pool_width = pool_param['pool_width']
738     stride = pool_param['stride']
739     N, C, H, W = x.shape
740
741     H_out = int (1 + (H - pool_height) / stride)
742     W_out = int (1 + (W - pool_width) / stride)
743
744     dx = np.zeros_like(x)
745
746     for n in range(N):
747         for c in range(C):

```

```

748         for i in range(H_out):
749             for j in range(W_out):
750                 block = x[n, c, i * stride: i * stride + pool_height, j * stride: j * stride +
pool_width]
751                 maximum = np.max(block)
752                 dx[n, c, i * stride: i * stride + pool_height, j * stride: j * stride + pool_width] =
\
753                                                                 (block == maximum) * dout[
n, c, i, j]
754
755     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
756     #####
757     #                                     END OF YOUR CODE                                     #
758     #####
759     return dx
760
761 def spatial_batchnorm_forward(x, gamma, beta, bn_param):
762     """
763     Computes the forward pass for spatial batch normalization.
764
765     Inputs:
766     - x: Input data of shape (N, C, H, W)
767     - gamma: Scale parameter, of shape (C,)
768     - beta: Shift parameter, of shape (C,)
769     - bn_param: Dictionary with the following keys:
770         - mode: 'train' or 'test'; required
771         - eps: Constant for numeric stability
772         - momentum: Constant for running mean / variance. momentum=0 means that
773           old information is discarded completely at every time step, while
774           momentum=1 means that new information is never incorporated. The
775           default of momentum=0.9 should work well in most situations.
776         - running_mean: Array of shape (D,) giving running mean of features
777         - running_var: Array of shape (D,) giving running variance of features
778
779     Returns a tuple of:
780     - out: Output data, of shape (N, C, H, W)
781     - cache: Values needed for the backward pass
782     """
783     out, cache = None, None
784
785     #####
786     # TODO: Implement the forward pass for spatial batch normalization.           #
787     #                                                                           #
788     # HINT: You can implement spatial batch normalization by calling the         #
789     # vanilla version of batch normalization you implemented above.             #
790     # Your implementation should be very short; ours is less than five lines.   #
791     #####
792     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
793
794     N, C, H, W = x.shape
795
796     # Transpose x so that the shape is (N, H, W, C), then reshape x into (N * H * W, C)
797     x_new = np.reshape(x.transpose(0, 2, 3, 1), (N * H * W, C))
798     out, cache = batchnorm_forward(x_new, gamma, beta, bn_param)
799
800     # Modify the final output
801     out = np.transpose(out.reshape(N, H, W, C), (0, 3, 1, 2))
802
803     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
804     #####
805     #                                     END OF YOUR CODE                                     #
806     #####
807
808     return out, cache
809
810
811
812 def spatial_batchnorm_backward(dout, cache):
813     """
814     Computes the backward pass for spatial batch normalization.
815
816     Inputs:
817     - dout: Upstream derivatives, of shape (N, C, H, W)
818     - cache: Values from the forward pass
819
820

```

```

821 Returns a tuple of:
822 - dx: Gradient with respect to inputs, of shape (N, C, H, W)
823 - dgamma: Gradient with respect to scale parameter, of shape (C,)
824 - dbeta: Gradient with respect to shift parameter, of shape (C,)
825 """
826 dx, dgamma, dbeta = None, None, None
827
828 #####
829 # TODO: Implement the backward pass for spatial batch normalization. #
830 # #
831 # HINT: You can implement spatial batch normalization by calling the #
832 # vanilla version of batch normalization you implemented above. #
833 # Your implementation should be very short; ours is less than five lines. #
834 #####
835 # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
836
837
838 N, C, H, W = dout.shape
839 dout_new = np.reshape(dout.transpose(0, 2, 3, 1), (N * H * W, C))
840 dx, dgamma, dbeta = batchnorm_backward(dout_new, cache)
841 dx = np.transpose(dx.reshape(N, H, W, C), (0, 3, 1, 2))
842
843 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
844 #####
845 # END OF YOUR CODE #
846 #####
847
848 return dx, dgamma, dbeta
849
850
851 def spatial_groupnorm_forward(x, gamma, beta, G, gn_param):
852     """
853     Computes the forward pass for spatial group normalization.
854     In contrast to layer normalization, group normalization splits each entry
855     in the data into G contiguous pieces, which it then normalizes independently.
856     Per feature shifting and scaling are then applied to the data, in a manner identical to that of batch
857     normalization and layer normalization.
858
859     Inputs:
860     - x: Input data of shape (N, C, H, W)
861     - gamma: Scale parameter, of shape (C,)
862     - beta: Shift parameter, of shape (C,)
863     - G: Integer number of groups to split into, should be a divisor of C
864     - gn_param: Dictionary with the following keys:
865         - eps: Constant for numeric stability
866
867     Returns a tuple of:
868     - out: Output data, of shape (N, C, H, W)
869     - cache: Values needed for the backward pass
870     """
871     out, cache = None, None
872     eps = gn_param.get('eps', 1e-5)
873     #####
874     # TODO: Implement the forward pass for spatial group normalization. #
875     # This will be extremely similar to the layer norm implementation. #
876     # In particular, think about how you could transform the matrix so that #
877     # the bulk of the code is similar to both train-time batch normalization #
878     # and layer normalization! #
879     #####
880     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
881
882     N, C, H, W = x.shape
883
884     # Just reshape x so that the number of group is multiplied by G while C is divided by G
885     x = np.reshape(x, (N * G, C // G * H * W))
886
887     # Other code are basically copied from layer norm implementation
888     x = x.T
889     sample_mean = np.mean(x, axis = 0)
890     sample_var = np.var(x, axis = 0)
891     sample_normalized = (x - sample_mean) / np.sqrt(sample_var + eps)
892
893     sample_normalized = np.reshape(sample_normalized.T, (N, C, H, W))
894     out = gamma * sample_normalized + beta
895     x = np.reshape(x.T, (N, C, H, W))

```

```

896 cache = (sample_normalized, gamma, beta, sample_mean, sample_var, x, eps, G)
897
898
899
900
901 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
902 #####
903 #                                     END OF YOUR CODE                                     #
904 #####
905 return out, cache
906
907
908 def spatial_groupnorm_backward(dout, cache):
909     """
910     Computes the backward pass for spatial group normalization.
911
912     Inputs:
913     - dout: Upstream derivatives, of shape (N, C, H, W)
914     - cache: Values from the forward pass
915
916     Returns a tuple of:
917     - dx: Gradient with respect to inputs, of shape (N, C, H, W)
918     - dgamma: Gradient with respect to scale parameter, of shape (C,)
919     - dbeta: Gradient with respect to shift parameter, of shape (C,)
920     """
921     dx, dgamma, dbeta = None, None, None
922
923     #####
924     # TODO: Implement the backward pass for spatial group normalization.      #
925     # This will be extremely similar to the layer norm implementation.        #
926     #####
927     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
928
929     sample_normalized, gamma, beta, sample_mean, sample_var, x, eps, G = cache
930     N, C, H, W = x.shape
931
932     # Remember to keep dimensions of dgamma and dbeta unchanged
933     dgamma = np.sum(dout * sample_normalized, axis = (0, 2, 3), keepdims=True)
934     dbeta = np.sum(dout, axis = (0, 2, 3), keepdims=True)
935     dx_hat = dout * gamma
936
937     # Reshape those matrices at first, then transpose them
938     sample_normalized = np.reshape(sample_normalized, (N * G, C // G * H * W))
939     x = np.reshape(x, (N * G, C // G * H * W))
940     dx_hat = np.reshape(dx_hat, (N * G, C // G * H * W))
941
942     sample_normalized = sample_normalized.T
943     x = x.T
944     dx_hat = dx_hat.T
945
946     # The following calculation is quite similar to layer norm backward
947     N_new = x.shape[0]
948
949     dvar = np.sum(dx_hat * (x - sample_mean) * (-1 / 2) * (sample_var + eps)**(-3 / 2), axis = 0)
950     dmean = np.sum(np.divide(-dx_hat, np.sqrt(sample_var + eps)), axis = 0) + \
951         dvar * np.sum(-2 * (x - sample_mean), axis = 0) / N_new
952     dx = dx_hat / np.sqrt(sample_var + eps) + dvar * 2 * (x - sample_mean) / N_new + dmean / N_new
953
954     dx = np.reshape(dx.T, (N, C, H, W))
955
956     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
957     #####
958     #                                     END OF YOUR CODE                                     #
959     #####
960     return dx, dgamma, dbeta
961
962
963 def svm_loss(x, y):
964     """
965     Computes the loss and gradient using for multiclass SVM classification.
966
967     Inputs:
968     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth
969         class for the ith input.
970     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
971         0 <= y[i] < C

```

```

972     Returns a tuple of:
973     - loss: Scalar giving the loss
974     - dx: Gradient of the loss with respect to x
975     """
976     N = x.shape[0]
977     correct_class_scores = x[np.arange(N), y]
978     margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
979     margins[np.arange(N), y] = 0
980     loss = np.sum(margins) / N
981     num_pos = np.sum(margins > 0, axis=1)
982     dx = np.zeros_like(x)
983     dx[margins > 0] = 1
984     dx[np.arange(N), y] -= num_pos
985     dx /= N
986     return loss, dx
987
988
989 def softmax_loss(x, y):
990     """
991     Computes the loss and gradient for softmax classification.
992
993     Inputs:
994     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth
995         class for the ith input.
996     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
997         0 <= y[i] < C
998
999     Returns a tuple of:
1000     - loss: Scalar giving the loss
1001     - dx: Gradient of the loss with respect to x
1002     """
1003     shifted_logits = x - np.max(x, axis=1, keepdims=True)
1004     Z = np.sum(np.exp(shifted_logits), axis=1, keepdims=True)
1005     log_probs = shifted_logits - np.log(Z)
1006     probs = np.exp(log_probs)
1007     N = x.shape[0]
1008     loss = -np.sum(log_probs[np.arange(N), y]) / N
1009     dx = probs.copy()
1010     dx[np.arange(N), y] -= 1
1011     dx /= N
1012     return loss, dx
1013

```

2 fc_net.py

```
1 from builtins import range
2 from builtins import object
3 import numpy as np
4
5 from cs231n.layers import *
6 from cs231n.layer_utils import *
7
8
9 class TwoLayerNet(object):
10     """
11     A two-layer fully-connected neural network with ReLU nonlinearity and
12     softmax loss that uses a modular layer design. We assume an input dimension
13     of D, a hidden dimension of H, and perform classification over C classes.
14
15     The architecture should be affine - relu - affine - softmax.
16
17     Note that this class does not implement gradient descent; instead, it
18     will interact with a separate Solver object that is responsible for running
19     optimization.
20
21     The learnable parameters of the model are stored in the dictionary
22     self.params that maps parameter names to numpy arrays.
23     """
24
25     def __init__(self, input_dim=3*32*32, hidden_dim=100, num_classes=10,
26                 weight_scale=1e-3, reg=0.0):
27         """
28         Initialize a new network.
29
30         Inputs:
31         - input_dim: An integer giving the size of the input
32         - hidden_dim: An integer giving the size of the hidden layer
33         - num_classes: An integer giving the number of classes to classify
34         - weight_scale: Scalar giving the standard deviation for random
35           initialization of the weights.
36         - reg: Scalar giving L2 regularization strength.
37         """
38         self.params = {}
39         self.reg = reg
40
41         #####
42         # TODO: Initialize the weights and biases of the two-layer net. Weights #
43         # should be initialized from a Gaussian centered at 0.0 with #
44         # standard deviation equal to weight_scale, and biases should be #
45         # initialized to zero. All weights and biases should be stored in the #
46         # dictionary self.params, with first layer weights #
47         # and biases using the keys 'W1' and 'b1' and second layer #
48         # weights and biases using the keys 'W2' and 'b2'. #
49         #####
50         # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
51
52         # Weights is initialized from a Gaussian centered at 0.0 with standard deviation equal to
53         # weight_scale
54         # Use np.random.normal function
55         self.params['W1'] = np.random.normal(0.0, weight_scale, (input_dim, hidden_dim))
56         self.params['b1'] = np.zeros((1, hidden_dim))
57         self.params['W2'] = np.random.normal(0.0, weight_scale, (hidden_dim, num_classes))
58         self.params['b2'] = np.zeros((1, num_classes))
59
60         # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
61         #####
62         # END OF YOUR CODE #
63         #####
64
65     def loss(self, X, y=None):
66         """
67         Compute loss and gradient for a minibatch of data.
68
69         Inputs:
70         - X: Array of input data of shape (N, d_1, ..., d_k)
71         - y: Array of labels, of shape (N,). y[i] gives the label for X[i].
72         """
```

```

73 Returns:
74 If y is None, then run a test-time forward pass of the model and return:
75 - scores: Array of shape (N, C) giving classification scores, where
76   scores[i, c] is the classification score for X[i] and class c.
77
78 If y is not None, then run a training-time forward and backward pass and
79 return a tuple of:
80 - loss: Scalar value giving the loss
81 - grads: Dictionary with the same keys as self.params, mapping parameter
82   names to gradients of the loss with respect to those parameters.
83 """
84 scores = None
85 #####
86 # TODO: Implement the forward pass for the two-layer net, computing the #
87 # class scores for X and storing them in the scores variable. #
88 #####
89 # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
90
91 W1, b1 = self.params['W1'], self.params['b1']
92 W2, b2 = self.params['W2'], self.params['b2']
93
94 # Firstly, do the affine_relu forward pass to get the hidden layer
95 hidden1, cache1 = affine_relu_forward(X, W1, b1)
96
97 # Secondly, do the affine forward pass
98 out, cache2 = affine_forward(hidden1, W2, b2)
99
100 scores = out
101
102 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
103 #####
104 #                               END OF YOUR CODE                               #
105 #####
106
107 # If y is None then we are in test mode so just return scores
108 if y is None:
109     return scores
110
111 loss, grads = 0, {}
112 #####
113 # TODO: Implement the backward pass for the two-layer net. Store the loss #
114 # in the loss variable and gradients in the grads dictionary. Compute data #
115 # loss using softmax, and make sure that grads[k] holds the gradients for #
116 # self.params[k]. Don't forget to add L2 regularization! #
117 # # #
118 # NOTE: To ensure that your implementation matches ours and you pass the #
119 # automated tests, make sure that your L2 regularization includes a factor #
120 # of 0.5 to simplify the expression for the gradient. #
121 #####
122 # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
123
124 # Use softmax to calculate the loss and gradient dout.
125 loss, dout = softmax_loss(scores, y)
126
127 # Based on the upstream gradient dout, use affine_backward to get the first downstream gradient.
128 dX2, dW2, db2 = affine_backward(dout, cache2)
129
130 # Based on the first downstream gradient dX2, use affine_relu_backward to get the second
131 # downstream gradient.
132 dX1, dW1, db1 = affine_relu_backward(dX2, cache1)
133
134 loss += 0.5 * self.reg * (np.sum(W1 * W1) + np.sum(W2 * W2))
135
136 dW2 += self.reg * W2
137 dW1 += self.reg * W1
138
139 grads['W1'] = dW1
140 grads['b1'] = db1
141 grads['W2'] = dW2
142 grads['b2'] = db2
143
144
145 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
146 #####
147 #                               END OF YOUR CODE                               #

```



```

148 #####
149
150     return loss, grads
151
152
153 class FullyConnectedNet(object):
154     """
155     A fully-connected neural network with an arbitrary number of hidden layers,
156     ReLU nonlinearities, and a softmax loss function. This will also implement
157     dropout and batch/layer normalization as options. For a network with L layers,
158     the architecture will be
159
160     {affine - [batch/layer norm] - relu - [dropout]} x (L - 1) - affine - softmax
161
162     where batch/layer normalization and dropout are optional, and the {...} block is
163     repeated L - 1 times.
164
165     Similar to the TwoLayerNet above, learnable parameters are stored in the
166     self.params dictionary and will be learned using the Solver class.
167     """
168
169     def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,
170                 dropout=1, normalization=None, reg=0.0,
171                 weight_scale=1e-2, dtype=np.float32, seed=None):
172         """
173         Initialize a new FullyConnectedNet.
174
175         Inputs:
176         - hidden_dims: A list of integers giving the size of each hidden layer.
177         - input_dim: An integer giving the size of the input.
178         - num_classes: An integer giving the number of classes to classify.
179         - dropout: Scalar between 0 and 1 giving dropout strength. If dropout=1 then
180           the network should not use dropout at all.
181         - normalization: What type of normalization the network should use. Valid values
182           are "batchnorm", "layernorm", or None for no normalization (the default).
183         - reg: Scalar giving L2 regularization strength.
184         - weight_scale: Scalar giving the standard deviation for random
185           initialization of the weights.
186         - dtype: A numpy datatype object; all computations will be performed using
187           this datatype. float32 is faster but less accurate, so you should use
188           float64 for numeric gradient checking.
189         - seed: If not None, then pass this random seed to the dropout layers. This
190           will make the dropout layers deterministic so we can gradient check the
191           model.
192         """
193         self.normalization = normalization
194         self.use_dropout = dropout != 1
195         self.reg = reg
196         self.num_layers = 1 + len(hidden_dims)
197         self.dtype = dtype
198         self.params = {}
199
200         #####
201         # TODO: Initialize the parameters of the network, storing all values in #
202         # the self.params dictionary. Store weights and biases for the first layer #
203         # in W1 and b1; for the second layer use W2 and b2, etc. Weights should be #
204         # initialized from a normal distribution centered at 0 with standard #
205         # deviation equal to weight_scale. Biases should be initialized to zero. #
206         # #
207         # When using batch normalization, store scale and shift parameters for the #
208         # first layer in gamma1 and beta1; for the second layer use gamma2 and #
209         # beta2, etc. Scale parameters should be initialized to ones and shift #
210         # parameters should be initialized to zeros. #
211         #####
212         # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
213
214         all_dims = [input_dim] + hidden_dims + [num_classes]
215         for i in range(len(all_dims) - 1):
216             self.params['W' + str(i + 1)] = np.random.normal(0.0, weight_scale, (all_dims[i], all_dims[i +
1]))
217             self.params['b' + str(i + 1)] = np.zeros((1, all_dims[i + 1]))
218
219         # If we haven't reached the final output layer, there may be a normalization layer
220         if i != self.num_layers - 1:
221             if self.normalization != None:
222                 self.params['gamma' + str(i + 1)] = np.ones((1, all_dims[i + 1]))

```

```

223         self.params['beta' + str(i + 1)] = np.zeros((1, all_dims[i + 1]))
224
225     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
226     #####
227     #                               END OF YOUR CODE                               #
228     #####
229
230     # When using dropout we need to pass a dropout_param dictionary to each
231     # dropout layer so that the layer knows the dropout probability and the mode
232     # (train / test). You can pass the same dropout_param to each dropout layer.
233     self.dropout_param = {}
234     if self.use_dropout:
235         self.dropout_param = {'mode': 'train', 'p': dropout}
236         if seed is not None:
237             self.dropout_param['seed'] = seed
238
239     # With batch normalization we need to keep track of running means and
240     # variances, so we need to pass a special bn_param object to each batch
241     # normalization layer. You should pass self.bn_params[0] to the forward pass
242     # of the first batch normalization layer, self.bn_params[1] to the forward
243     # pass of the second batch normalization layer, etc.
244     self.bn_params = []
245     if self.normalization=='batchnorm':
246         self.bn_params = [{'mode': 'train'} for i in range(self.num_layers - 1)]
247     if self.normalization=='layernorm':
248         self.bn_params = [{} for i in range(self.num_layers - 1)]
249
250     # Cast all parameters to the correct datatype
251     for k, v in self.params.items():
252         self.params[k] = v.astype(dtype)
253
254
255 def loss(self, X, y=None):
256     """
257     Compute loss and gradient for the fully-connected net.
258
259     Input / output: Same as TwoLayerNet above.
260     """
261     X = X.astype(self.dtype)
262     mode = 'test' if y is None else 'train'
263
264     # Set train/test mode for batchnorm params and dropout param since they
265     # behave differently during training and testing.
266     if self.use_dropout:
267         self.dropout_param['mode'] = mode
268     if self.normalization=='batchnorm':
269         for bn_param in self.bn_params:
270             bn_param['mode'] = mode
271     scores = None
272     #####
273     # TODO: Implement the forward pass for the fully-connected net, computing #
274     # the class scores for X and storing them in the scores variable.         #
275     #                                                                           #
276     # When using dropout, you'll need to pass self.dropout_param to each    #
277     # dropout forward pass.                                                  #
278     #                                                                           #
279     # When using batch normalization, you'll need to pass self.bn_params[0] to #
280     # the forward pass for the first batch normalization layer, pass         #
281     # self.bn_params[1] to the forward pass for the second batch normalization #
282     # layer, etc.                                                            #
283     #####
284     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
285
286     affine_cache = {}
287     norm_cache = {}
288     relu_cache = {}
289     dropout_cache = {}
290
291     input_param = X
292
293     for i in range(self.num_layers - 1):
294         W, b = self.params['W' + str(i + 1)], self.params['b' + str(i + 1)]
295
296         if self.normalization == 'batchnorm':
297             # the first part is affine - batchnorm - relu
298             gamma, beta = self.params['gamma' + str(i + 1)], self.params['beta' + str(i + 1)]

```

```

299         affine_outcome, affine_cache[i + 1] = affine_forward(input_param, W, b)
300         norm_outcome, norm_cache[i + 1] = batchnorm_forward(affine_outcome, gamma, beta, self.
bn_params[i])
301         relu_outcome, relu_cache[i + 1] = relu_forward(norm_outcome)
302
303     elif self.normalization == 'layernorm':
304         # the first part is affine - layernorm - relu
305         gamma, beta = self.params['gamma' + str(i + 1)], self.params['beta' + str(i + 1)]
306         affine_outcome, affine_cache[i + 1] = affine_forward(input_param, W, b)
307         norm_outcome, norm_cache[i + 1] = layernorm_forward(affine_outcome, gamma, beta, self.
bn_params[i])
308         relu_outcome, relu_cache[i + 1] = relu_forward(norm_outcome)
309
310     else:
311         # the first part is affine - relu
312         relu_outcome, (affine_cache[i + 1], relu_cache[i + 1]) = affine_relu_forward(input_param,
W, b)
313
314     if self.use_dropout:
315         dropout_outcome, dropout_cache[i + 1] = dropout_forward(relu_outcome, self.dropout_param)
316
317     # Update input_param
318     input_param = dropout_outcome if self.use_dropout else relu_outcome
319
320     # Get the last layer
321     scores, last_cache = affine_forward(input_param,
322                                         self.params['W' + str(self.num_layers)],
323                                         self.params['b' + str(self.num_layers)])
324
325     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
326     #####
327     #                               END OF YOUR CODE                               #
328     #####
329
330     # If test mode return early
331     if mode == 'test':
332         return scores
333
334     loss, grads = 0.0, {}
335     #####
336     # TODO: Implement the backward pass for the fully-connected net. Store the #
337     # loss in the loss variable and gradients in the grads dictionary. Compute #
338     # data loss using softmax, and make sure that grads[k] holds the gradients #
339     # for self.params[k]. Don't forget to add L2 regularization!                #
340     #                                                                           #
341     # When using batch/layer normalization, you don't need to regularize the scale #
342     # and shift parameters.                                                    #
343     #                                                                           #
344     # NOTE: To ensure that your implementation matches ours and you pass the #
345     # automated tests, make sure that your L2 regularization includes a factor #
346     # of 0.5 to simplify the expression for the gradient.                    #
347     #####
348     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
349
350     # Use softmax to calculate the loss and gradient dout.
351     loss, dout = softmax_loss(scores, y)
352
353     reg_loss = 0
354     for i in range(self.num_layers):
355         reg_loss += 0.5 * self.reg * np.sum(np.square(self.params['W' + str(i+1)]))
356
357     loss += reg_loss
358
359     # Based on the upstream gradient dout, use affine_backward to get the first downstream gradient.
360     dX, dW, db = affine_backward(dout, last_cache)
361     dW += self.reg * self.params['W' + str(self.num_layers)]
362     grads['W' + str(self.num_layers)] = dW
363     grads['b' + str(self.num_layers)] = db
364
365     for i in range(self.num_layers - 1, 0, -1):
366         if self.use_dropout:
367             # If there is a dropout at the end
368             dX = dropout_backward(dX, dropout_cache[i])
369
370         dX = relu_backward(dX, relu_cache[i])
371

```

```

372     if self.normalization == 'batchnorm':
373         # If there is a batchnorm in the middle
374         dX, dgamma, dbeta = batchnorm_backward(dX, norm_cache[i])
375         grads['gamma'+str(i)] = dgamma
376         grads['beta'+str(i)] = dbeta
377     elif self.normalization == 'layernorm':
378         # If there is a layernorm in the middle
379         dX, dgamma, dbeta = layernorm_backward(dX, norm_cache[i])
380         grads['gamma'+str(i)] = dgamma
381         grads['beta'+str(i)] = dbeta
382
383     dX, dW, db = affine_backward(dX, affine_cache[i])
384
385     dW += self.reg * self.params['W' + str(i)]
386
387     grads['W'+str(i)] = dW
388     grads['b'+str(i)] = db
389
390
391     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
392     #####
393     #                                     END OF YOUR CODE                                     #
394     #####
395
396     return loss, grads

```

3 optim.py

```
1 import numpy as np
2
3 """
4 This file implements various first-order update rules that are commonly used
5 for training neural networks. Each update rule accepts current weights and the
6 gradient of the loss with respect to those weights and produces the next set of
7 weights. Each update rule has the same interface:
8
9 def update(w, dw, config=None):
10
11 Inputs:
12 - w: A numpy array giving the current weights.
13 - dw: A numpy array of the same shape as w giving the gradient of the
14     loss with respect to w.
15 - config: A dictionary containing hyperparameter values such as learning
16     rate, momentum, etc. If the update rule requires caching values over many
17     iterations, then config will also hold these cached values.
18
19 Returns:
20 - next_w: The next point after the update.
21 - config: The config dictionary to be passed to the next iteration of the
22     update rule.
23
24 NOTE: For most update rules, the default learning rate will probably not
25 perform well; however the default values of the other hyperparameters should
26 work well for a variety of different problems.
27
28 For efficiency, update rules may perform in-place updates, mutating w and
29 setting next_w equal to w.
30 """
31
32 def sgd(w, dw, config=None):
33     """
34     Performs vanilla stochastic gradient descent.
35
36     config format:
37     - learning_rate: Scalar learning rate.
38     """
39     if config is None: config = {}
40     config.setdefault('learning_rate', 1e-2)
41
42     w -= config['learning_rate'] * dw
43     return w, config
44
45
46 def sgd_momentum(w, dw, config=None):
47     """
48     Performs stochastic gradient descent with momentum.
49
50     config format:
51     - learning_rate: Scalar learning rate.
52     - momentum: Scalar between 0 and 1 giving the momentum value.
53         Setting momentum = 0 reduces to sgd.
54     - velocity: A numpy array of the same shape as w and dw used to store a
55         moving average of the gradients.
56     """
57     if config is None: config = {}
58     config.setdefault('learning_rate', 1e-2)
59     config.setdefault('momentum', 0.9)
60     v = config.get('velocity', np.zeros_like(w))
61
62     next_w = None
63     #####
64     # TODO: Implement the momentum update formula. Store the updated value in #
65     # the next_w variable. You should also use and update the velocity v.      #
66     #####
67     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
68
69
70     v = config['momentum'] * v - config['learning_rate'] * dw
71     w += v
72     next_w = w
73
```

```

74 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
75 #####
76 #                                     END OF YOUR CODE                                     #
77 #####
78 config['velocity'] = v
79
80 return next_w, config
81
82
83
84 def rmsprop(w, dw, config=None):
85     """
86     Uses the RMSProp update rule, which uses a moving average of squared
87     gradient values to set adaptive per-parameter learning rates.
88
89     config format:
90     - learning_rate: Scalar learning rate.
91     - decay_rate: Scalar between 0 and 1 giving the decay rate for the squared
92       gradient cache.
93     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
94     - cache: Moving average of second moments of gradients.
95     """
96     if config is None: config = {}
97     config.setdefault('learning_rate', 1e-2)
98     config.setdefault('decay_rate', 0.99)
99     config.setdefault('epsilon', 1e-8)
100    config.setdefault('cache', np.zeros_like(w))
101
102    next_w = None
103    #####
104    # TODO: Implement the RMSprop update formula, storing the next value of w #
105    # in the next_w variable. Don't forget to update cache value stored in #
106    # config['cache']. #
107    #####
108    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
109
110    config['cache'] = config['decay_rate'] * config['cache'] + (1 - config['decay_rate']) * dw**2
111    w += - config['learning_rate'] * dw / (np.sqrt(config['cache']) + config['epsilon'])
112    next_w = w
113
114    # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
115    #####
116    #                                     END OF YOUR CODE                                     #
117    #####
118
119    return next_w, config
120
121
122 def adam(w, dw, config=None):
123     """
124     Uses the Adam update rule, which incorporates moving averages of both the
125     gradient and its square and a bias correction term.
126
127     config format:
128     - learning_rate: Scalar learning rate.
129     - beta1: Decay rate for moving average of first moment of gradient.
130     - beta2: Decay rate for moving average of second moment of gradient.
131     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
132     - m: Moving average of gradient.
133     - v: Moving average of squared gradient.
134     - t: Iteration number.
135     """
136     if config is None: config = {}
137     config.setdefault('learning_rate', 1e-3)
138     config.setdefault('beta1', 0.9)
139     config.setdefault('beta2', 0.999)
140     config.setdefault('epsilon', 1e-8)
141     config.setdefault('m', np.zeros_like(w))
142     config.setdefault('v', np.zeros_like(w))
143     config.setdefault('t', 0)
144
145     next_w = None
146     #####
147     # TODO: Implement the Adam update formula, storing the next value of w in #
148     # the next_w variable. Don't forget to update the m, v, and t variables #
149     # stored in config. #

```

```

150 #
151 # NOTE: In order to match the reference output, please modify t _before_ #
152 # using it in any calculations. #
153 #####
154 # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
155
156 m = config['m']
157 v = config['v']
158 beta1 = config['beta1']
159 beta2 = config['beta2']
160 learning_rate = config['learning_rate']
161 epsilon = config['epsilon']
162 t = config['t'] + 1
163
164 m = beta1*m + (1-beta1)*dw
165 mt = m / (1-beta1**t)
166 v = beta2*v + (1-beta2)*(dw**2)
167 vt = v / (1-beta2**t)
168 w += - learning_rate * mt / (np.sqrt(vt) + epsilon)
169
170 next_w = w
171 config['m'] = m
172 config['v'] = v
173 config['t'] = t
174
175 # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
176 #####
177 #                                     END OF YOUR CODE #
178 #####
179
180 return next_w, config

```

4 cnn.py

```
1 from builtins import object
2 import numpy as np
3
4 from cs231n.layers import *
5 from cs231n.fast_layers import *
6 from cs231n.layer_utils import *
7
8
9 class ThreeLayerConvNet(object):
10     """
11     A three-layer convolutional network with the following architecture:
12
13     conv - relu - 2x2 max pool - affine - relu - affine - softmax
14
15     The network operates on minibatches of data that have shape (N, C, H, W)
16     consisting of N images, each with height H and width W and with C input
17     channels.
18     """
19
20     def __init__(self, input_dim=(3, 32, 32), num_filters=32, filter_size=7,
21                 hidden_dim=100, num_classes=10, weight_scale=1e-3, reg=0.0,
22                 dtype=np.float32):
23         """
24         Initialize a new network.
25
26         Inputs:
27         - input_dim: Tuple (C, H, W) giving size of input data
28         - num_filters: Number of filters to use in the convolutional layer
29         - filter_size: Width/height of filters to use in the convolutional layer
30         - hidden_dim: Number of units to use in the fully-connected hidden layer
31         - num_classes: Number of scores to produce from the final affine layer.
32         - weight_scale: Scalar giving standard deviation for random initialization
33           of weights.
34         - reg: Scalar giving L2 regularization strength
35         - dtype: numpy datatype to use for computation.
36         """
37         self.params = {}
38         self.reg = reg
39         self.dtype = dtype
40
41         #####
42         # TODO: Initialize weights and biases for the three-layer convolutional #
43         # network. Weights should be initialized from a Gaussian centered at 0.0 #
44         # with standard deviation equal to weight_scale; biases should be #
45         # initialized to zero. All weights and biases should be stored in the #
46         # dictionary self.params. Store weights and biases for the convolutional #
47         # layer using the keys 'W1' and 'b1'; use keys 'W2' and 'b2' for the #
48         # weights and biases of the hidden affine layer, and keys 'W3' and 'b3' #
49         # for the weights and biases of the output affine layer. #
50         # #
51         # IMPORTANT: For this assignment, you can assume that the padding #
52         # and stride of the first convolutional layer are chosen so that #
53         # **the width and height of the input are preserved**. Take a look at #
54         # the start of the loss() function to see how that happens. #
55         #####
56         # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
57
58         C, H, W = input_dim
59         self.params['W1'] = weight_scale * np.random.randn(num_filters, C, filter_size, filter_size)
60         self.params['b1'] = np.zeros((1, num_filters))
61
62         # 2x2 max pool reduces the width and height by half
63         self.params['W2'] = weight_scale * np.random.randn(num_filters * H * W // (2 * 2), hidden_dim)
64         self.params['b2'] = np.zeros((1, hidden_dim))
65
66         self.params['W3'] = weight_scale * np.random.randn(hidden_dim, num_classes)
67         self.params['b3'] = np.zeros((1, num_classes))
68
69         # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
70         #####
71         # END OF YOUR CODE #
72         #####
73
```



```

74     for k, v in self.params.items():
75         self.params[k] = v.astype(dtype)
76
77
78 def loss(self, X, y=None):
79     """
80     Evaluate loss and gradient for the three-layer convolutional network.
81
82     Input / output: Same API as TwoLayerNet in fc_net.py.
83     """
84     W1, b1 = self.params['W1'], self.params['b1']
85     W2, b2 = self.params['W2'], self.params['b2']
86     W3, b3 = self.params['W3'], self.params['b3']
87
88     # pass conv_param to the forward pass for the convolutional layer
89     # Padding and stride chosen to preserve the input spatial size
90     filter_size = W1.shape[2]
91     conv_param = {'stride': 1, 'pad': (filter_size - 1) // 2}
92
93     # pass pool_param to the forward pass for the max-pooling layer
94     pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}
95
96     scores = None
97     #####
98     # TODO: Implement the forward pass for the three-layer convolutional net, #
99     # computing the class scores for X and storing them in the scores      #
100     # variable.                                                            #
101     #                                                                       #
102     # Remember you can use the functions defined in cs231n/fast_layers.py and #
103     # cs231n/layer_utils.py in your implementation (already imported).      #
104     #####
105     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
106
107     conv_relu_pool_outcome, conv_relu_pool_cache = conv_relu_pool_forward(X, W1, b1, conv_param,
pool_param)
108     affine_relu_outcome, affine_relu_cache = affine_relu_forward(conv_relu_pool_outcome, W2, b2)
109     scores, last_cache = affine_forward(affine_relu_outcome, W3, b3)
110
111
112     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
113     #####
114     #                               END OF YOUR CODE                       #
115     #####
116
117     if y is None:
118         return scores
119
120     loss, grads = 0, {}
121     #####
122     # TODO: Implement the backward pass for the three-layer convolutional net, #
123     # storing the loss and gradients in the loss and grads variables. Compute #
124     # data loss using softmax, and make sure that grads[k] holds the gradients #
125     # for self.params[k]. Don't forget to add L2 regularization!              #
126     #                                                                           #
127     # NOTE: To ensure that your implementation matches ours and you pass the  #
128     # automated tests, make sure that your L2 regularization includes a factor #
129     # of 0.5 to simplify the expression for the gradient.                    #
130     #####
131     # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
132
133     loss, dout = softmax_loss(scores, y)
134     reg_loss = 0.5 * self.reg * (np.sum(np.square(self.params['W1'])) + np.sum(np.square(self.params['
W2']))) + np.sum(np.square(self.params['W3']))
135     loss += reg_loss
136
137     dX3, dW3, db3 = affine_backward(dout, last_cache)
138     dX2, dW2, db2 = affine_relu_backward(dX3, affine_relu_cache)
139     dX1, dW1, db1 = conv_relu_pool_backward(dX2, conv_relu_pool_cache)
140
141     dW1 += self.reg * self.params['W1']
142     dW2 += self.reg * self.params['W2']
143     dW3 += self.reg * self.params['W3']
144
145     grads['W1'] = dW1
146     grads['b1'] = db1
147     grads['W2'] = dW2

```

```

148     grads[ 'b2' ] = db2
149     grads[ 'W3' ] = dW3
150     grads[ 'b3' ] = db3
151
152
153
154     # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
155     #####
156     #                               END OF YOUR CODE                               #
157     #####
158
159     return loss , grads

```