

Fully-Connected Neural Nets

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

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

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

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

    return out, cache
```

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

```
def layer_backward(dout, cache):
    """
    Receive dout (derivative of loss with respect to outputs) and cache,
    and compute derivative with respect to inputs.
    """
    # Unpack cache values
    x, w, z, out = cache

    # Use values in cache to compute derivatives
    dx = # Derivative of loss with respect to x
    dw = # Derivative of loss with respect to w

    return dx, dw
```

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

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

```

In [31]: # As usual, a bit of setup
from __future__ import print_function
import time
import numpy as np
import matplotlib.pyplot as plt
from cs682.classifiers.fc_net import *
from cs682.data_utils import get_CIFAR10_data
from cs682.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
from cs682.solver import Solver

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

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2

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

```

The autoreload extension is already loaded. To reload it, use:

```
%reload_ext autoreload
```

```

In [61]: # Load the (preprocessed) CIFAR10 data.

```

```

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

X_train = data['X_train']
y_train = data['y_train']
X_val = data['X_val']
y_val = data['y_val']
X_test = data['X_test']
y_test = data['y_test']

```

```

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

```

Affine layer: forward

Open the file `cs682/layers.py` and implement the `affine_forward` function.

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

```
In [33]: # 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 [34]: # Test the affine_backward function
np.random.seed(231)
x = np.random.randn(10, 2, 3)
w = np.random.randn(6, 5)
b = np.random.randn(5)
dout = np.random.randn(10, 5)

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

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

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

```
Testing affine_backward function:
dx error: 5.399100368651805e-11
dw error: 9.904211865398145e-11
db error: 2.4122867568119087e-11
```

ReLU activation: forward

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

```
In [35]: # 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 [36]: 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:

Sigmoid and ReLU will have this problem. For sigmoid, the inputs that are small. For ReLU, it is inputs that are negative.

"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 `cs682/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 [37]: from cs682.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 `cs682/layers.py`.

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

```

In [38]: 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 `cs682/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 [39]: np.random.seed(231)
N, D, H, C = 3, 5, 50, 7
X = np.random.randn(N, D)
y = np.random.randint(C, size=N)

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

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

print('Testing test-time forward pass ... ')
model.params['W1'] = np.linspace(-0.7, 0.3, num=D*H).reshape(D, H)
model.params['b1'] = np.linspace(-0.1, 0.9, num=H)
model.params['W2'] = np.linspace(-0.3, 0.4, num=H*C).reshape(H, C)
model.params['b2'] = np.linspace(-0.9, 0.1, num=C)
X = np.linspace(-5.5, 4.5, num=N*D).reshape(D, N).T
scores = model.loss(X)
correct_scores = np.asarray(
    [[11.53165108, 12.2917344, 13.05181771, 13.81190102, 14.57198434, 15.33206765, 16.09215096],
     [12.05769098, 12.74614105, 13.43459113, 14.1230412, 14.81149128, 15.49994135, 16.18839143],
     [12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.05099822, 15.66781506, 16.2846319 ]])
scores_diff = np.abs(scores - correct_scores).sum()
assert scores_diff < 1e-6, 'Problem with test-time forward pass'

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

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

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

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

```



```

Testing initialization ...
Testing test-time forward pass ...
Testing training loss (no regularization)
Running numeric gradient check with reg = 0.0
W1 relative error: 1.52e-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: 7.98e-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 `cs682/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 [58]: model = TwoLayerNet()
solver = None

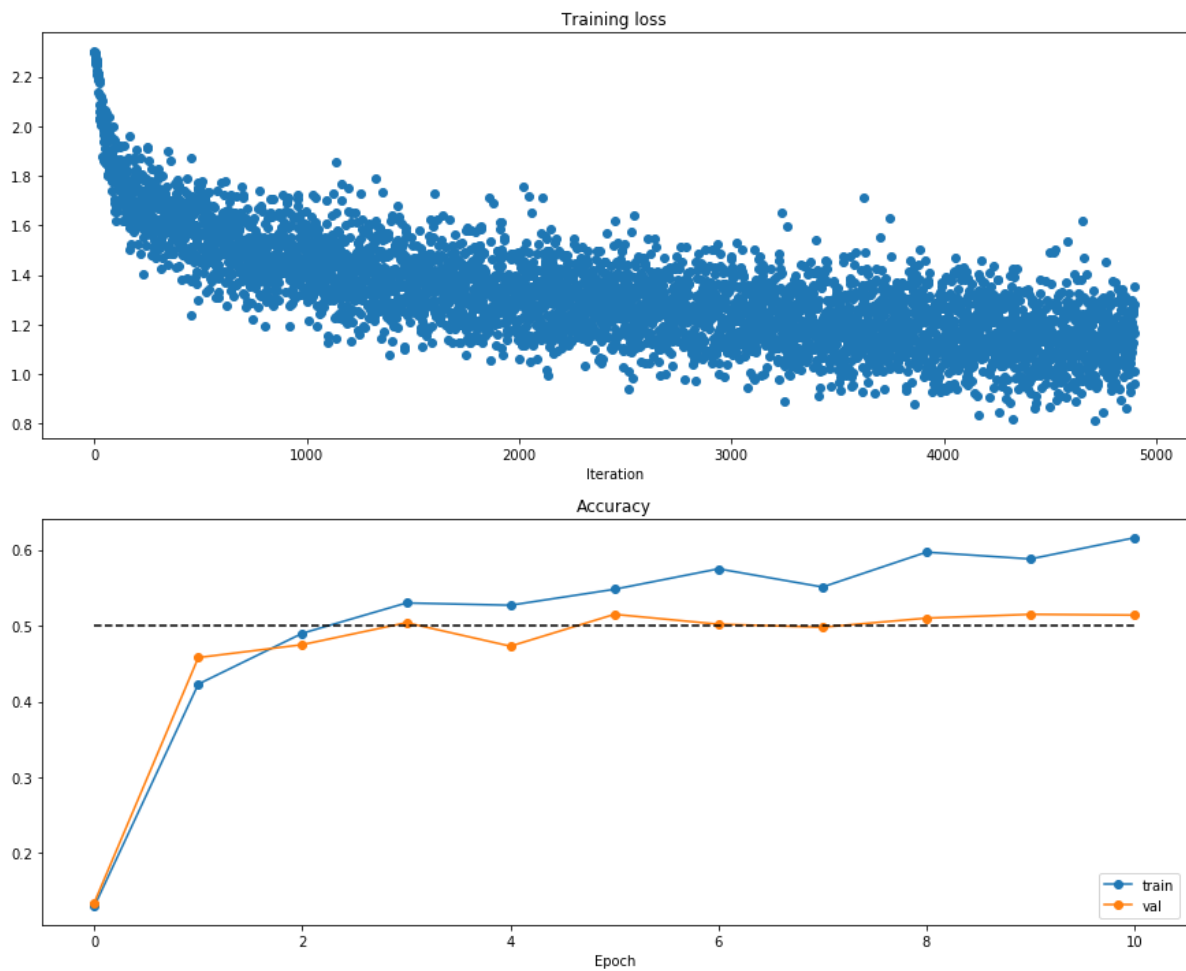
#####
# TODO: Use a Solver instance to train a TwoLayerNet that achieves at least #
# 50% accuracy on the validation set.                                     #
#####
data = {
    'X_train': X_train,
    'y_train': y_train,
    'X_val': X_val,
    'y_val': y_val,
}
solver = Solver(model, data,
                update_rule='sgd',
                optim_config={
                    'learning_rate': 1e-3,
                },
                verbose = False,
                lr_decay=0.95,
                num_epochs=10, batch_size=100,
                print_every=100)
solver.train()
#####
#                               END OF YOUR CODE                               #
#####

```

In [59]: *# 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 `cs682/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 [43]: np.random.seed(231)
N, D, H1, H2, C = 2, 15, 20, 30, 10
X = np.random.randn(N, D)
y = np.random.randint(C, size=(N,))

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

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

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

```

```

Running check with reg = 0
Initial loss: 2.3004790897684924
W1 relative error: 1.48e-07
W2 relative error: 2.21e-05
W3 relative error: 3.53e-07
b1 relative error: 5.38e-09
b2 relative error: 2.09e-09
b3 relative error: 5.80e-11
Running check with reg = 3.14
Initial loss: 7.052114776533016
W1 relative error: 7.36e-09
W2 relative error: 6.87e-08
W3 relative error: 3.48e-08
b1 relative error: 1.48e-08
b2 relative error: 1.72e-09
b3 relative error: 1.80e-10

```

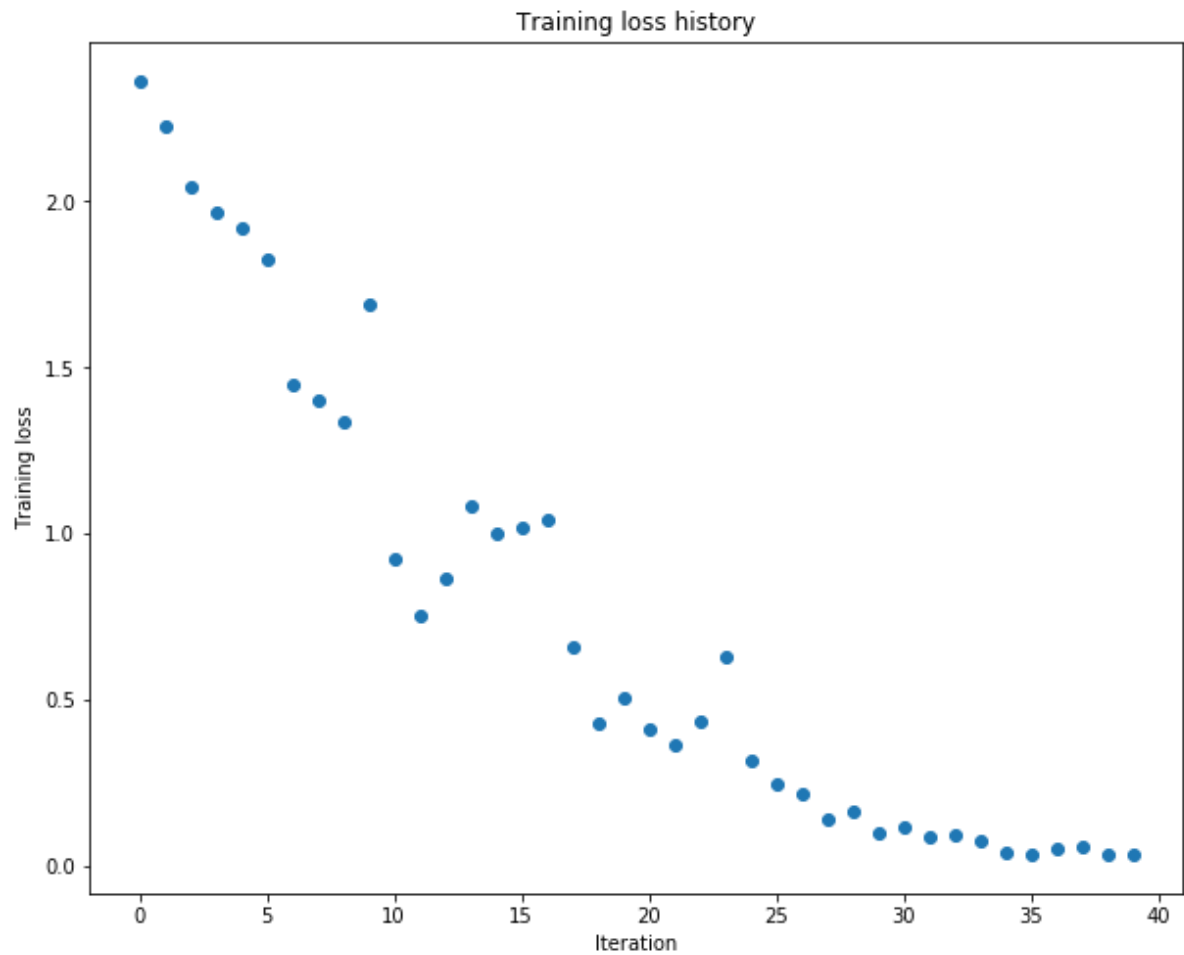
As another sanity check, make sure you can overfit a small dataset of 50 images. First we will try a three-layer network with 100 units in each hidden layer. In the following cell, tweak the learning rate and initialization scale to overfit and achieve 100% training accuracy within 20 epochs.

In [45]: *# 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
learning_rate = 9e-3
model = FullyConnectedNet([100, 100],
    weight_scale=weight_scale, dtype=np.float64)
solver = Solver(model, small_data,
    print_every=10, num_epochs=20, batch_size=25,
    update_rule='sgd',
    optim_config={
        'learning_rate': learning_rate,
    },
    verbose = False
)
solver.train()

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



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

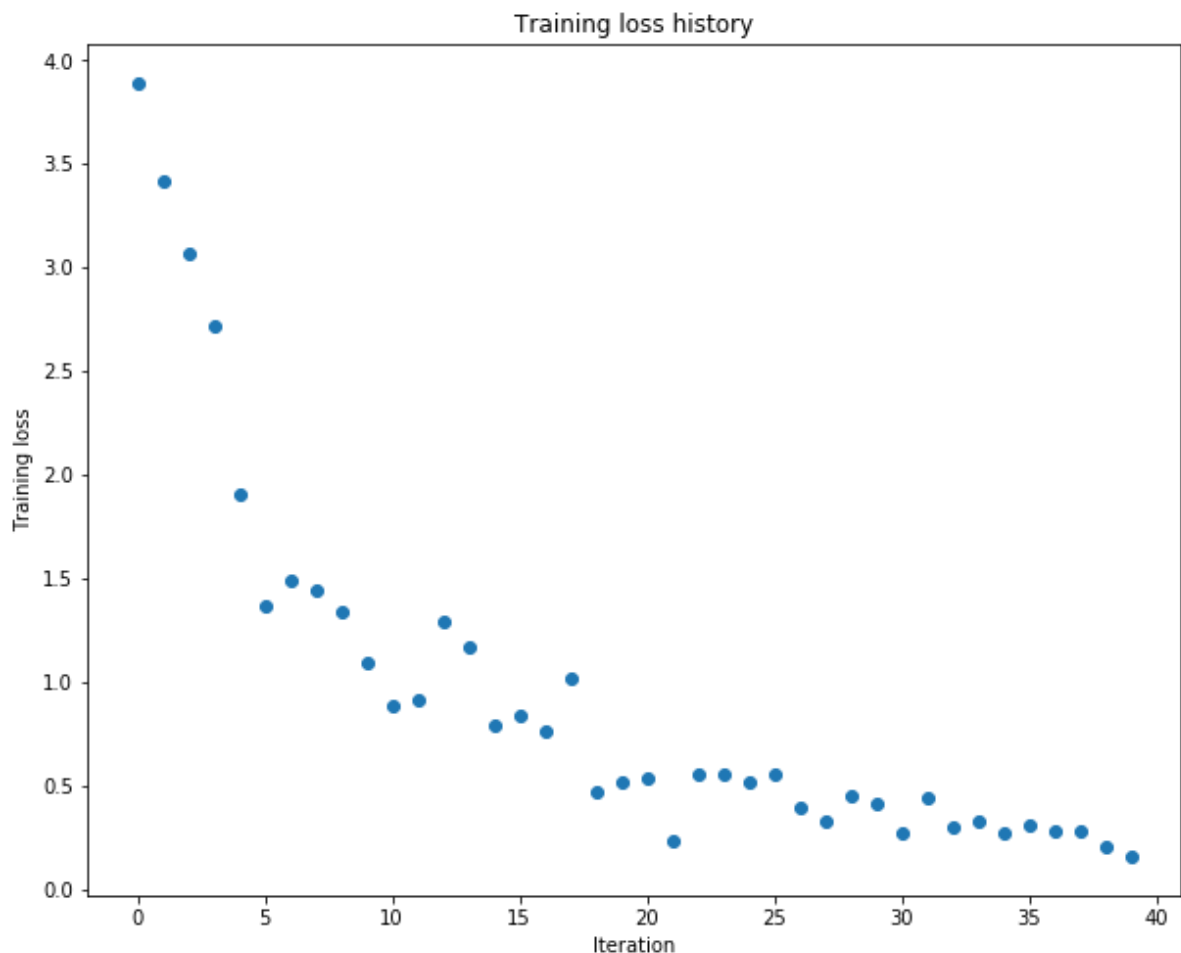
In [49]: *# TODO: Use a five-layer Net to overfit 50 training examples by
tweaking just the learning rate and initialization scale.*

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

learning_rate = 2e-3
weight_scale = 5e-2
model = FullyConnectedNet([100, 100, 100, 100],
                           weight_scale=weight_scale, dtype=np.float64)
solver = Solver(model, small_data,
                print_every=10, num_epochs=20, batch_size=25,
                update_rule='sgd',
                optim_config={
                    'learning_rate': learning_rate,
                },
                verbose = True
            )
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: 3.892097
(Epoch 0 / 20) train acc: 0.260000; val_acc: 0.140000
(Epoch 1 / 20) train acc: 0.300000; val_acc: 0.118000
(Epoch 2 / 20) train acc: 0.400000; val_acc: 0.105000
(Epoch 3 / 20) train acc: 0.500000; val_acc: 0.122000
(Epoch 4 / 20) train acc: 0.560000; val_acc: 0.129000
(Epoch 5 / 20) train acc: 0.660000; val_acc: 0.126000
(Iteration 11 / 40) loss: 0.886108
(Epoch 6 / 20) train acc: 0.660000; val_acc: 0.125000
(Epoch 7 / 20) train acc: 0.780000; val_acc: 0.125000
(Epoch 8 / 20) train acc: 0.820000; val_acc: 0.120000
(Epoch 9 / 20) train acc: 0.920000; val_acc: 0.147000
(Epoch 10 / 20) train acc: 0.920000; val_acc: 0.126000
(Iteration 21 / 40) loss: 0.538654
(Epoch 11 / 20) train acc: 0.940000; val_acc: 0.128000
(Epoch 12 / 20) train acc: 0.940000; val_acc: 0.131000
(Epoch 13 / 20) train acc: 0.960000; val_acc: 0.142000
(Epoch 14 / 20) train acc: 0.960000; val_acc: 0.136000
(Epoch 15 / 20) train acc: 0.960000; val_acc: 0.137000
(Iteration 31 / 40) loss: 0.272727
(Epoch 16 / 20) train acc: 0.980000; val_acc: 0.141000
(Epoch 17 / 20) train acc: 0.980000; val_acc: 0.140000
(Epoch 18 / 20) train acc: 0.980000; val_acc: 0.148000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.150000
(Epoch 20 / 20) train acc: 0.980000; val_acc: 0.149000



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 neural network needed to be initialize with smaller weights. I think this is because there are more layers and more backpropagation happening, and larger weights may cause exploding gradients.

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 <https://compsci682-fa18.github.io/notes/neural-networks-3/#sgd> (<https://compsci682-fa18.github.io/notes/neural-networks-3/#sgd>) for more information.

Open the file `cs682/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 [50]: from cs682.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 [52]: num_train = 4000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

solvers = {}

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

    solver = Solver(model, small_data,
                    num_epochs=5, batch_size=100,
                    update_rule=update_rule,
                    optim_config={
                        'learning_rate': 1e-2,
                    },
                    verbose=False)
    solvers[update_rule] = solver
    solver.train()
    print()

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

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

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

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

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

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

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

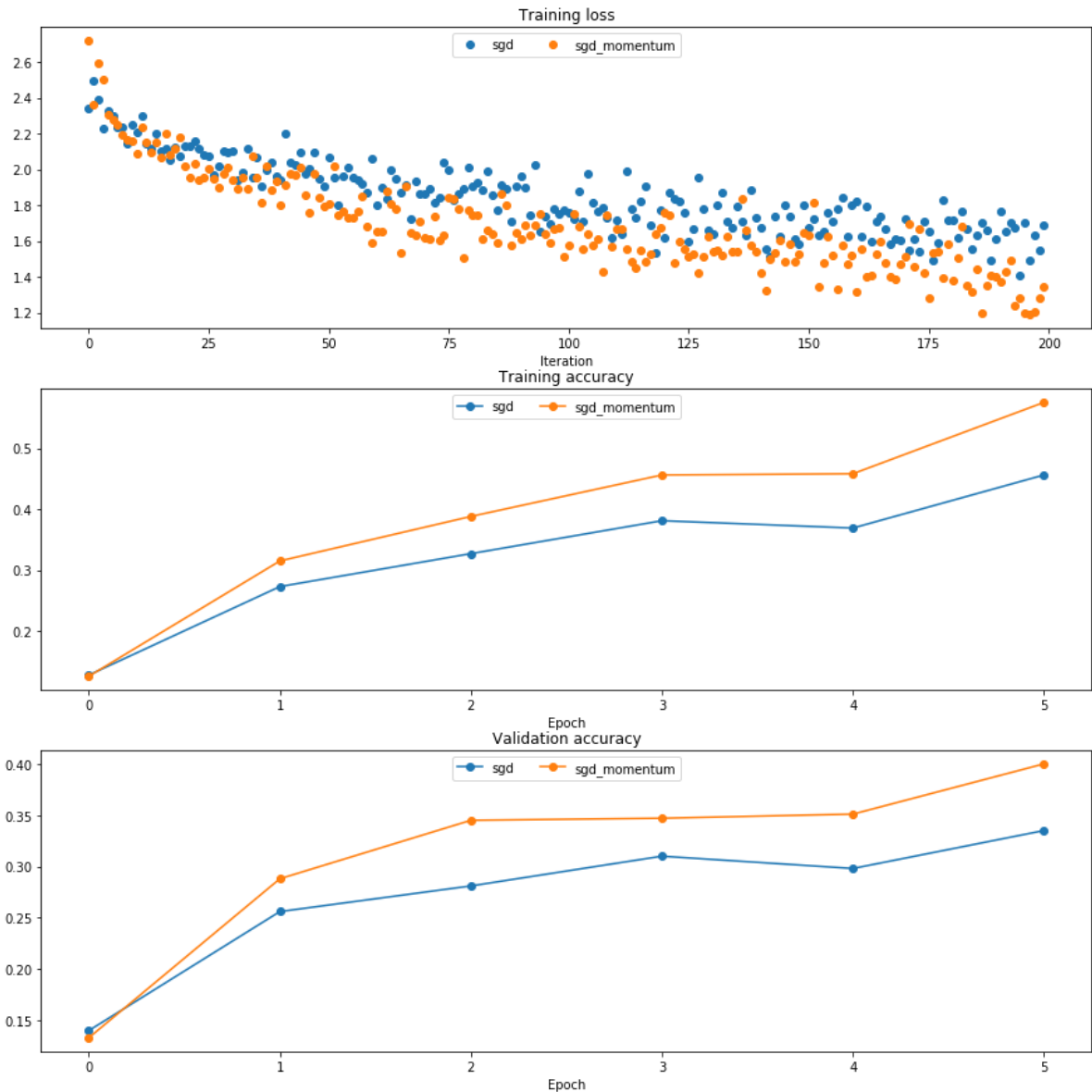
```

running with `sgd`

running with `sgd_momentum`

C:\Users\kucharskib\AppData\Roaming\Python\Python36\site-packages\matplotlib\cbook\deprecation.py:106: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

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



RMSProp and Adam

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

In the file `cs682/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 [53]: # Test RMSProp implementation
from cs682.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 [54]: # Test Adam implementation
from cs682.optim import adam

N, D = 4, 5
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
m = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)
v = np.linspace(0.7, 0.5, num=N*D).reshape(N, D)

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

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

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

next_w error: 1.1395691798535431e-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 [55]: 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=False)
    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)

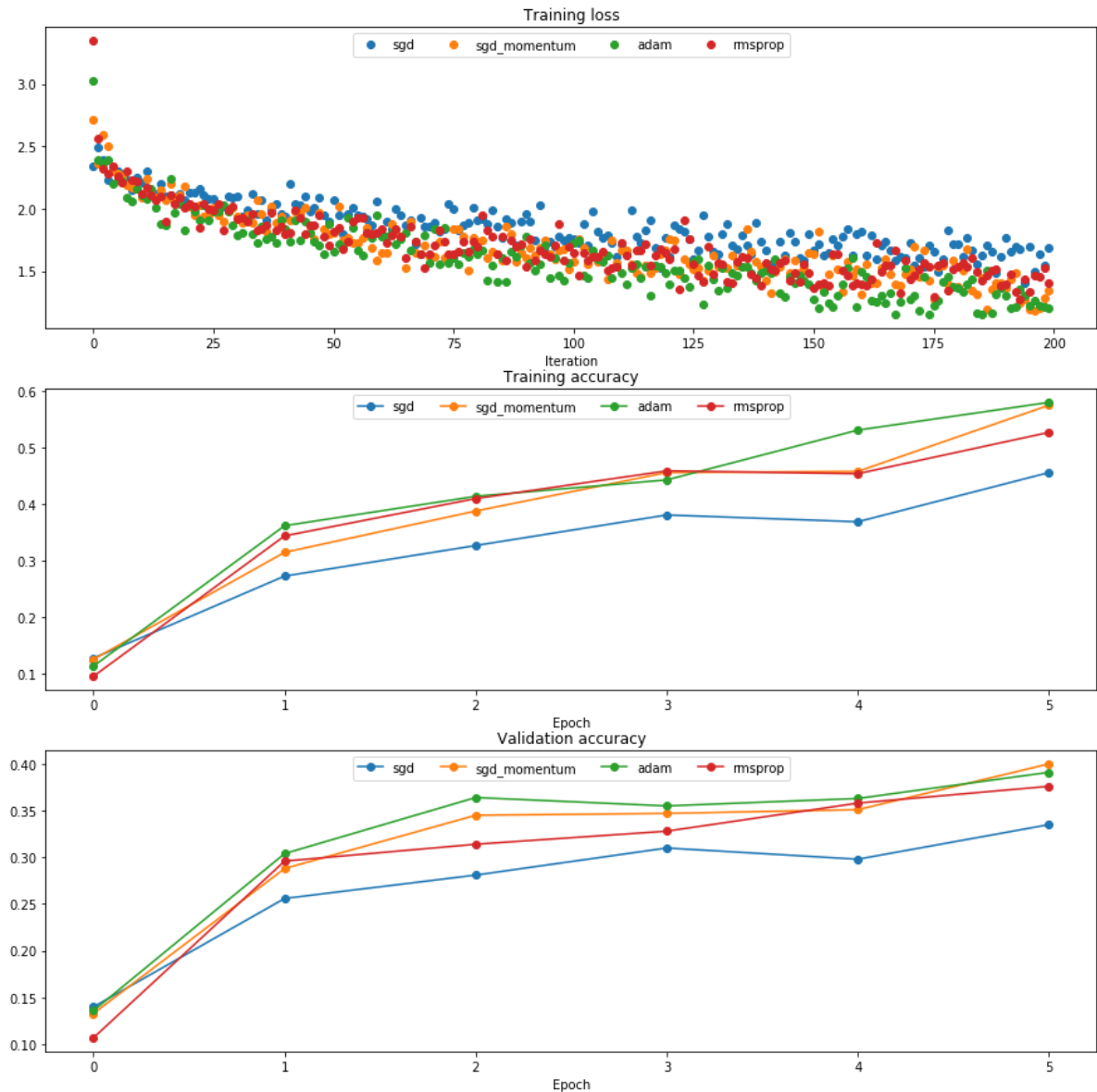
    for i in [1, 2, 3]:
        plt.subplot(3, 1, i)
        plt.legend(loc='upper center', ncol=4)
    plt.gcf().set_size_inches(15, 15)
    plt.show()
```

running with adam

running with rmsprop

C:\Users\kucharskib\AppData\Roaming\Python\Python36\site-packages\matplotlib\cbook\deprecation.py:106: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

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



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:

The updates become very small because of the squared dw term. Since the dw is squared, the cache will also be positive, which will lead to the division of a larger and larger number over time. This leads to smaller updates over time.

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 [29]: 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. #
#####
weight_scale = 5e-2
learning_rate = 1e-3
model = FullyConnectedNet([200, 100, 50, 25], normalization='batchnorm', weight_scale=weight_scale, dtype=np.float64)
solver = Solver(model, data,
                 print_every=100,
                 num_epochs=5,
                 batch_size=64,
                 update_rule='adam',
                 optim_config={
                     'learning_rate': learning_rate,
                 },
                 verbose = False
            )
solver.train()

best_model = model
#####
#                               END OF YOUR CODE                               #
#####

```

Test your model!

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

```

In [62]: 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: 0.523

Test set accuracy: 0.488

Dropout

Dropout [1] is a technique for regularizing neural networks by randomly setting some features 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 (<https://arxiv.org/abs/1207.0580>).

```
In [6]: # As usual, a bit of setup
from __future__ import print_function
import time
import numpy as np
import matplotlib.pyplot as plt
from cs682.classifiers.fc_net import *
from cs682.data_utils import get_CIFAR10_data
from cs682.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
from cs682.solver import Solver

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

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2

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

The autoreload extension is already loaded. To reload it, use:

```
%reload_ext autoreload
```

```
In [7]: # Load the (preprocessed) CIFAR10 data.
```

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

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

Dropout forward pass

In the file `cs682/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 [8]: 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 `cs682/layers.py`, implement the backward pass for dropout. After doing so, run the following cell to numerically gradient-check your implementation.

```
In [9]: 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:

If you do not divide by p then the next layer activation values will be reduced by a factor of p . Dividing "bumps" the values back up to what they're supposed to be.

Fully-connected nets with Dropout

In the file `cs682/classifiers/fc_net.py`, modify your implementation to use dropout. Specifically, if the constructor of the net receives a value that is not 1 for the dropout parameter, then the net should add dropout immediately after every ReLU nonlinearity. After doing so, run the following to numerically gradient-check your implementation.

```

In [10]: 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 [11]: # 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=False, print_every=100)
    solver.train()
    solvers[dropout] = solver
```

1

0.25

In [12]: *# Plot train and validation accuracies of the two models*

```

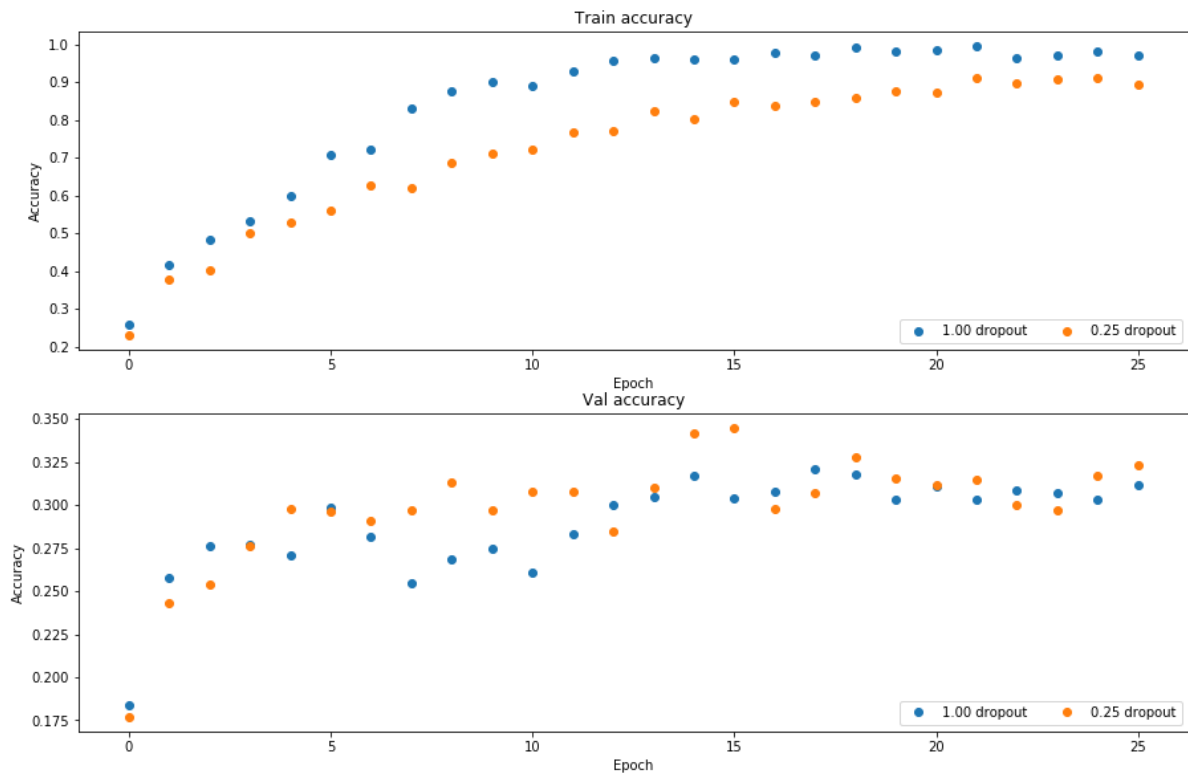
train_accs = []
val_accs = []
for dropout in dropout_choices:
    solver = solvers[dropout]
    train_accs.append(solver.train_acc_history[-1])
    val_accs.append(solver.val_acc_history[-1])

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 is lower with dropout compared to no dropout. The testing accuracies are similar, but dropout seems to be a little higher in some places. This does prove the theory that dropout reduces overfitting, because the gap between the test/train accuracy is smaller when compared with no dropout.

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). How should we modify p , if at all, if we decide to decrease the size of the hidden layers (that is, the number of nodes in each layer)?

Answer:

You should keep more hidden nodes as you decrease the layer size. For image classification, you have many features in your beginning nodes where there is higher chance of overfitting. This means you want to dropout more nodes to prevent overfitting in the early layers. For the smaller layers, you want to a smaller chance of dropping nodes because there arent as many nodes and less of a chance to overfit.

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 [3] 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 [3] hypothesize that the shifting distribution of features inside deep neural networks may make training deep networks more difficult. To overcome this problem, [3] 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.

[3] Sergey Ioffe and Christian Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift", ICML 2015. (<https://arxiv.org/abs/1502.03167>)

```
In [119]: # As usual, a bit of setup
import time
import numpy as np
import matplotlib.pyplot as plt
from cs682.classifiers.fc_net import *
from cs682.data_utils import get_CIFAR10_data
from cs682.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
from cs682.solver import Solver

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

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2

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

def print_mean_std(x, axis=0):
    print(' means: ', x.mean(axis=axis))
    print(' stds: ', x.std(axis=axis))
    print()
```

The autoreload extension is already loaded. To reload it, use:

```
%reload_ext autoreload
```

```
In [120]: # Load the (preprocessed) CIFAR10 data.
data = get_CIFAR10_data()
for k, v in data.items():
    print('%s: ' % k, v.shape)
```

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

Batch normalization: forward

In the file `cs682/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 would be helpful!

In [121]: *# 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: [5.32907052e-17 7.04991621e-17 4.22578639e-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 [122]: *# 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 [123]: # 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.

Given a set of inputs $X = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_N \end{bmatrix}$, we first calculate the mean $\mu = \frac{1}{N} \sum_{k=1}^N x_k$ and variance

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

With μ and v calculated, we can calculate the standard deviation $\sigma = \sqrt{v + \epsilon}$ and normalized data Y with

$$y_i = \frac{x_i - \mu}{\sigma}.$$

The meat of our problem is to get $\frac{\partial L}{\partial X}$ from the upstream gradient $\frac{\partial L}{\partial Y}$. It might be challenging to directly reason about the gradients over X and Y - try reasoning about it in terms of x_i and y_i first.

You will need to come up with the derivations for $\frac{\partial L}{\partial x_i}$, by relying on the Chain Rule to first calculate the intermediate $\frac{\partial \mu}{\partial x_i}$, $\frac{\partial v}{\partial x_i}$, $\frac{\partial \sigma}{\partial x_i}$, then assemble these pieces to calculate $\frac{\partial y_i}{\partial x_i}$. You should make sure each of the intermediary steps are all as simple as possible.

After doing so, implement the simplified batch normalization backward pass in the function `batchnorm_backward_alt` and compare the two implementations by running the following. Your two implementations should compute nearly identical results, but the alternative implementation should be a bit faster.

```

In [125]: 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: 8.41670880961546e-13
dgamma difference: 0.0
dbeta difference: 0.0
speedup: 1.00x

```

Fully Connected Nets with Batch Normalization

Now that you have a working implementation for batch normalization, go back to your FullyConnectedNet in the file `cs682/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 `cs682/layer_utils.py`. If you decide to do so, do it in the file `cs682/classifiers/fc_net.py`.


```

In [126]: 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: 2.21e-07
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: 2.28e-06
W3 relative error: 1.11e-08
b1 relative error: 1.12e-08
b2 relative error: 2.27e-08
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 [128]: 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'],
}

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)

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

solver = Solver(model, small_data,
    num_epochs=10, batch_size=50,
    update_rule='adam',
    optim_config={
        'learning_rate': 1e-3,
    },
    verbose=False, print_every=20)
solver.train()
```

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.

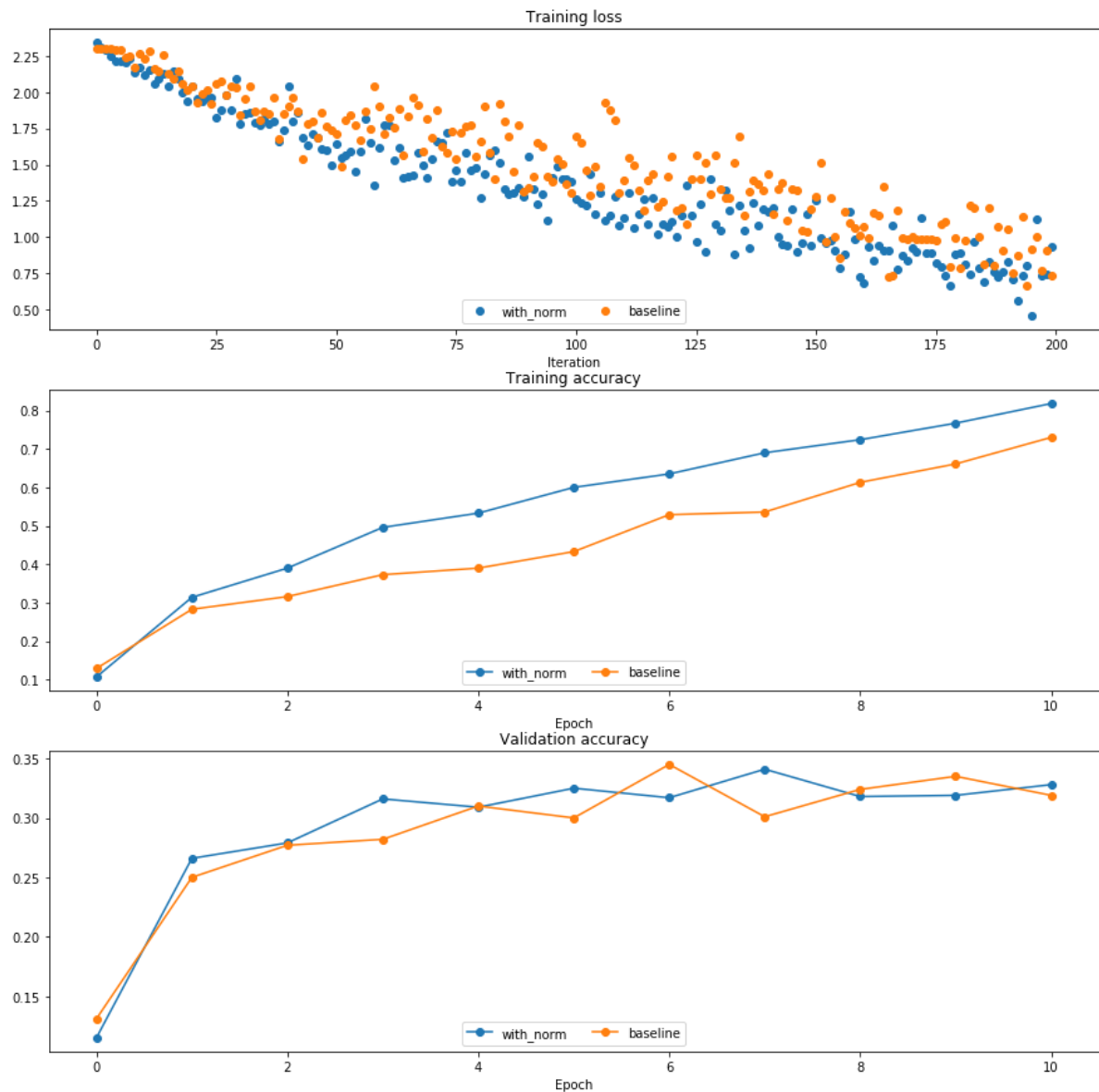
```

In [111]: def plot_training_history(title, label, baseline, bn_solvers, plot_fn, bl_marker='.', bn_marker='.', labels=None):
            """utility function for plotting training history"""
            plt.title(title)
            plt.xlabel(label)
            bn_plots = [plot_fn(bn_solver) for bn_solver in bn_solvers]
            bl_plot = plot_fn(baseline)
            num_bn = len(bn_plots)
            for i in range(num_bn):
                label='with_norm'
                if labels is not None:
                    label += str(labels[i])
                plt.plot(bn_plots[i], bn_marker, label=label)
            label='baseline'
            if labels is not None:
                label += str(labels[0])
            plt.plot(bl_plot, bl_marker, label=label)
            plt.legend(loc='lower center', ncol=num_bn+1)

            plt.subplot(3, 1, 1)
            plot_training_history('Training loss','Iteration', solver, [bn_solver], \
                                lambda x: x.loss_history, bl_marker='o', bn_marker='o')
            plt.subplot(3, 1, 2)
            plot_training_history('Training accuracy','Epoch', solver, [bn_solver], \
                                lambda x: x.train_acc_history, bl_marker='-o', bn_marker='-o')
            plt.subplot(3, 1, 3)
            plot_training_history('Validation accuracy','Epoch', solver, [bn_solver], \
                                lambda x: x.val_acc_history, bl_marker='-o', bn_marker='-o')

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

```



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.

```

In [129]: np.random.seed(231)
           # Try training a very deep net with batchnorm
           hidden_dims = [50, 50, 50, 50, 50, 50, 50]
           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'],
           }

           bn_solvers_ws = {}
           solvers_ws = {}
           weight_scales = np.logspace(-4, 0, num=20)
           for i, weight_scale in enumerate(weight_scales):
               print('Running weight scale %d / %d' % (i + 1, len(weight_scales)))
               bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, normalization='batchnorm')
               model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, normalization=None)

               bn_solver = Solver(bn_model, small_data,
                                   num_epochs=10, batch_size=50,
                                   update_rule='adam',
                                   optim_config={
                                       'learning_rate': 1e-3,
                                   },
                                   verbose=False, print_every=200)
               bn_solver.train()
               bn_solvers_ws[weight_scale] = bn_solver

               solver = Solver(model, small_data,
                                   num_epochs=10, batch_size=50,
                                   update_rule='adam',
                                   optim_config={
                                       'learning_rate': 1e-3,
                                   },
                                   verbose=False, print_every=200)
               solver.train()
               solvers_ws[weight_scale] = solver

```

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

```
In [130]: # Plot results of weight scale experiment
best_train_accs, bn_best_train_accs = [], []
best_val_accs, bn_best_val_accs = [], []
final_train_loss, bn_final_train_loss = [], []

for ws in weight_scales:
    best_train_accs.append(max(solvers_ws[ws].train_acc_history))
    bn_best_train_accs.append(max(bn_solvers_ws[ws].train_acc_history))

    best_val_accs.append(max(solvers_ws[ws].val_acc_history))
    bn_best_val_accs.append(max(bn_solvers_ws[ws].val_acc_history))

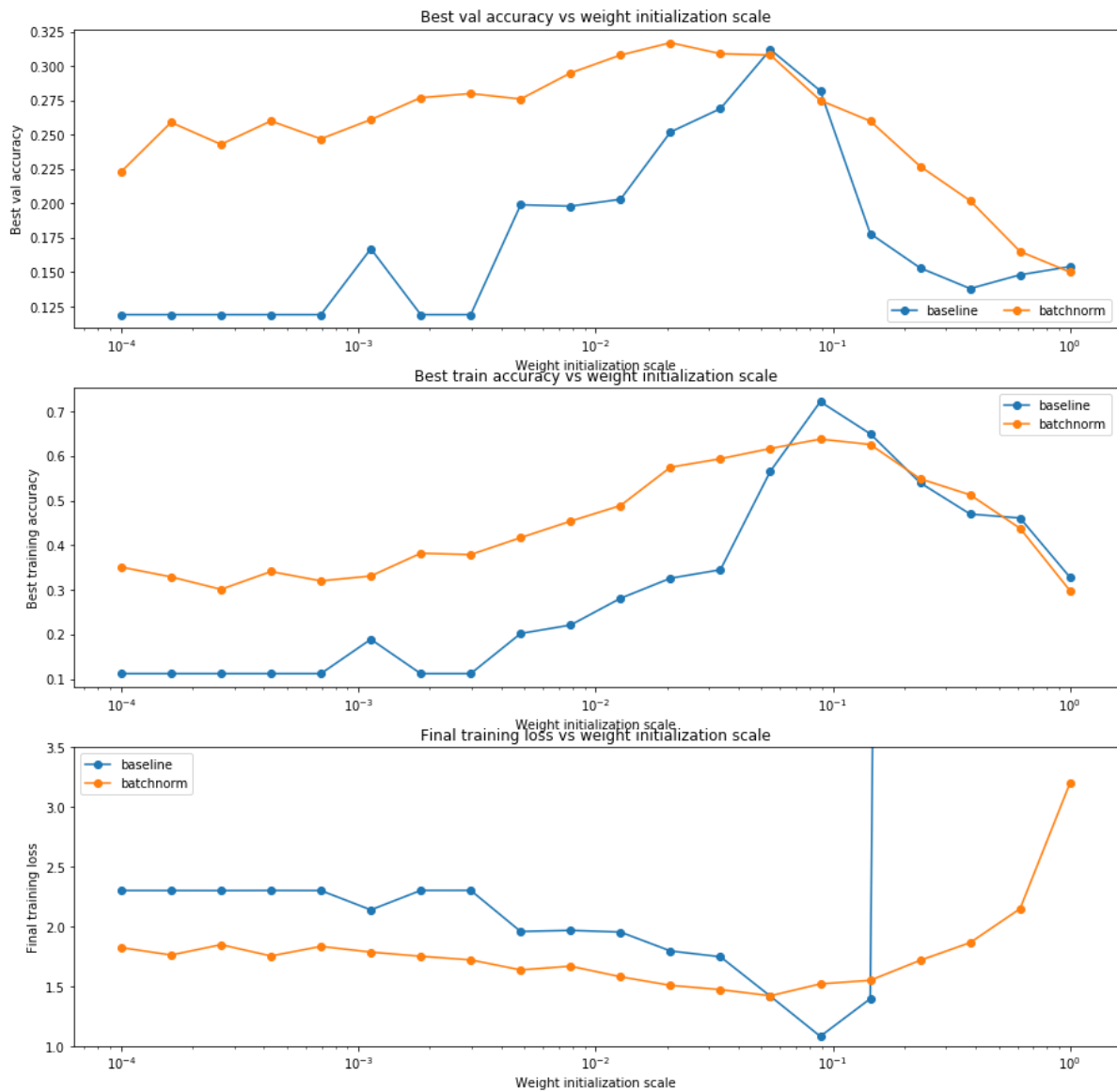
    final_train_loss.append(np.mean(solvers_ws[ws].loss_history[-100:]))
    bn_final_train_loss.append(np.mean(bn_solvers_ws[ws].loss_history[-100:]))

plt.subplot(3, 1, 1)
plt.title('Best val accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Best val accuracy')
plt.semilogx(weight_scales, best_val_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_val_accs, '-o', label='batchnorm')
plt.legend(ncol=2, loc='lower right')

plt.subplot(3, 1, 2)
plt.title('Best train accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Best training accuracy')
plt.semilogx(weight_scales, best_train_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_train_accs, '-o', label='batchnorm')
plt.legend()

plt.subplot(3, 1, 3)
plt.title('Final training loss vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Final training loss')
plt.semilogx(weight_scales, final_train_loss, '-o', label='baseline')
plt.semilogx(weight_scales, bn_final_train_loss, '-o', label='batchnorm')
plt.legend()
plt.gca().set_ylim(1.0, 3.5)

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



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:

Batch normalization always has a higher training and test accuracy and also a lower training loss if the weights are not too large. Batch norm outperforms the baseline because it ensures that the mean and variance of each layer is roughly 0 and 1. Typically, you want to preprocess the image dataset to be mean 0 variance 1 to the input layer because this will normalize the entire training set and make the network less sensitive to features that are on different scales. This same idea is applied to each layer of the network via batch norm.

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.

```

In [94]: def run_batchsize_experiments(normalization_mode):
    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'],
    }
    n_epochs=10
    weight_scale = 2e-2
    batch_sizes = [5,10,50]
    lr = 10**(-3.5)
    solver_bsize = batch_sizes[0]

    print('No normalization: batch size = ',solver_bsize)
    model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, normalization=None)
    solver = Solver(model, small_data,
        num_epochs=n_epochs, batch_size=solver_bsize,
        update_rule='adam',
        optim_config={
            'learning_rate': lr,
        },
        verbose=False)
    solver.train()

    bn_solvers = []
    for i in range(len(batch_sizes)):
        b_size=batch_sizes[i]
        print('Normalization: batch size = ',b_size)
        bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, normalization=normalization_mode)
        bn_solver = Solver(bn_model, small_data,
            num_epochs=n_epochs, batch_size=b_size,
            update_rule='adam',
            optim_config={
                'learning_rate': lr,
            },
            verbose=False)
        bn_solver.train()
        bn_solvers.append(bn_solver)

    return bn_solvers, solver, batch_sizes

batch_sizes = [5,10,50]
bn_solvers_bsize, solver_bsize, batch_sizes = run_batchsize_experiments('batchnorm')

```

No normalization: batch size = 5

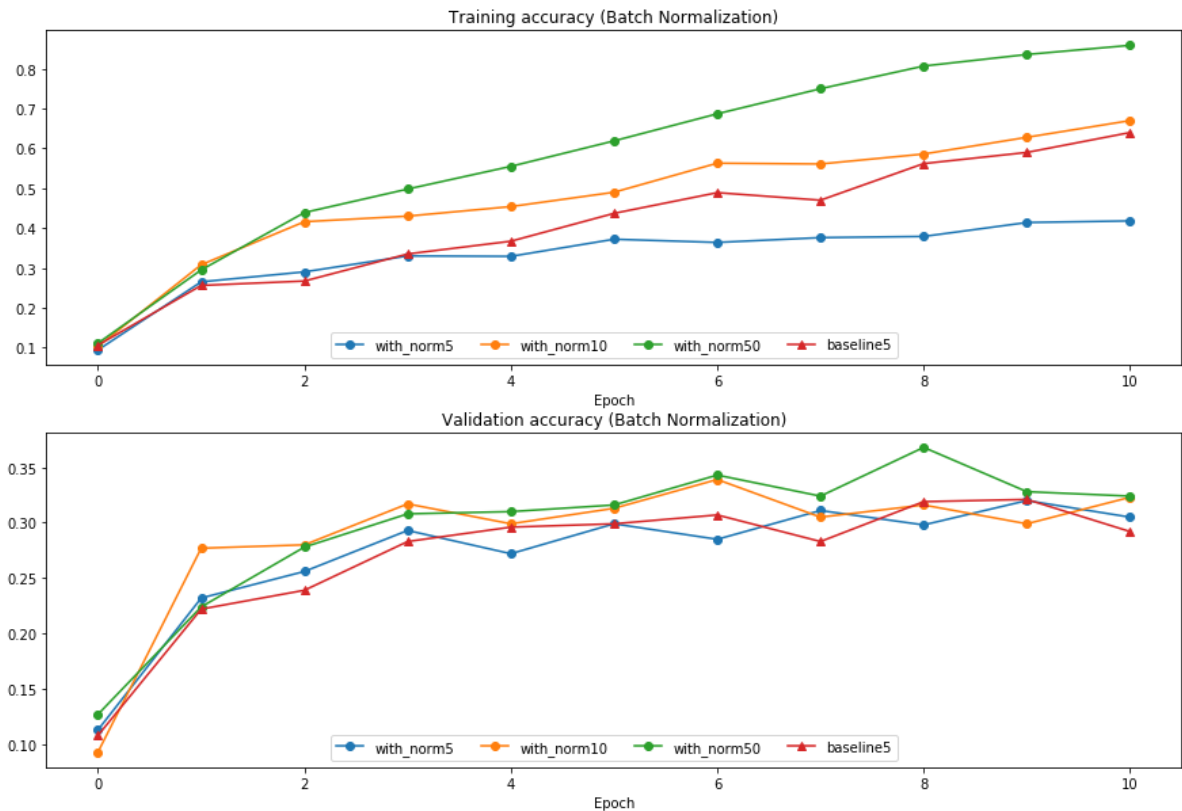
Normalization: batch size = 5

Normalization: batch size = 10

Normalization: batch size = 50

```
In [95]: 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_sizes)
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_sizes)

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:

In general, batch norm is better with larger batch sizes. I think this is because there is a larger sample to get a better accuracy for the mean and variance for each batch. With few batches, there may be a lot larger variance in the distribution. The more minibatches you will have a better mean of the data.

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 [4]. 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.

[4] Ba, Jimmy Lei, Jamie Ryan Kiros, and Geoffrey E. Hinton. "Layer Normalization." *stat 1050 (2016)*: 21. (<https://arxiv.org/pdf/1607.06450.pdf>).

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. is analogous to batch normalization
2. is analogous to layer 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 `cs682/layers.py`, implement the forward pass for layer normalization in the function `layernorm_backward`.

Run the cell below to check your results.

- In `cs682/layers.py`, implement the backward pass for layer normalization in the function `layernorm_backward`.

Run the second cell below to check your results.

- Modify `cs682/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 [115]: *# 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 [116]: # 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: 1.433615146847572e-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 [118]: 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_sizes)
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_sizes)

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

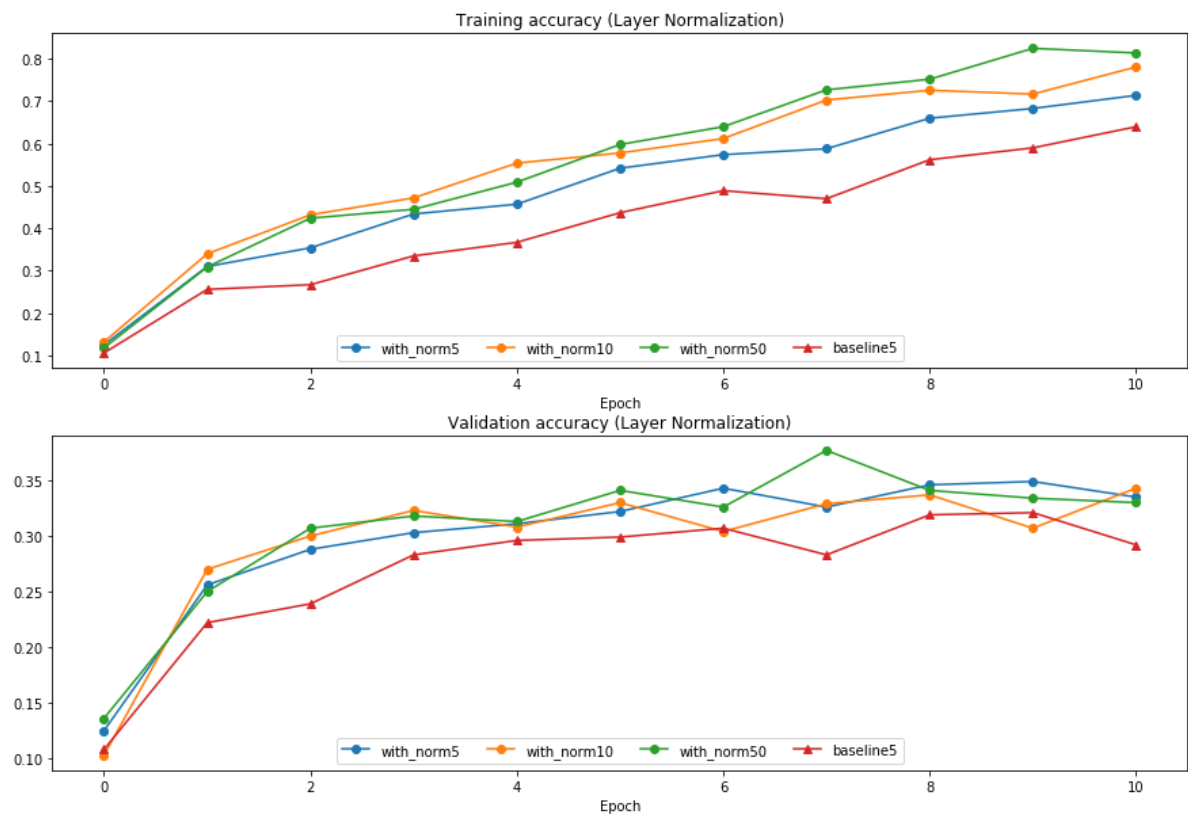
```

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:

1. because you get an average over all of the features, so if you do not have a lot of features you will not have a good average

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.

```
In [104]: # As usual, a bit of setup
import numpy as np
import matplotlib.pyplot as plt
from cs682.classifiers.cnn import *
from cs682.data_utils import get_CIFAR10_data
from cs682.gradient_check import eval_numerical_gradient_array, eval_numerical_gradient
from cs682.layers import *
from cs682.fast_layers import *
from cs682.solver import Solver

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

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2

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

The autoreload extension is already loaded. To reload it, use:

```
%reload_ext autoreload
```

```
In [105]: # Load the (preprocessed) CIFAR10 data.
```

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

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

Convolution: Naive forward pass

The core of a convolutional network is the convolution operation. In the file `cs682/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 [106]: 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.

```

In [107]: from scipy.misc import imread, imresize

kitten, puppy = imread('kitten.jpg'), imread('puppy.jpg')
# kitten is wide, and puppy is already square
d = kitten.shape[1] - kitten.shape[0]
kitten_cropped = kitten[:, d//2:-d//2, :]

img_size = 200 # Make this smaller if it runs too slow
x = np.zeros((2, 3, img_size, img_size))
x[0, :, :, :] = imresize(puppy, (img_size, img_size)).transpose((2, 0, 1))
x[1, :, :, :] = imresize(kitten_cropped, (img_size, img_size)).transpose((2, 0, 1))

# Set up a convolutional weights holding 2 filters, each 3x3
w = np.zeros((2, 3, 3, 3))

# The first filter converts the image to grayscale.
# Set up the red, green, and blue channels of the filter.
w[0, 0, :, :] = [[0, 0, 0], [0, 0.3, 0], [0, 0, 0]]
w[0, 1, :, :] = [[0, 0, 0], [0, 0.6, 0], [0, 0, 0]]
w[0, 2, :, :] = [[0, 0, 0], [0, 0.1, 0], [0, 0, 0]]

# Second filter detects horizontal edges in the blue channel.
w[1, 2, :, :] = [[1, 2, 1], [0, 0, 0], [-1, -2, -1]]

# Vector of biases. We don't need any bias for the grayscale
# filter; but for the edge detection filter we want to add 128
# to each output so that nothing is negative.
b = np.array([0, 128])

# Compute the result of convolving each input in x with each filter in w,
# offsetting by b, and storing the results in out.
out, _ = conv_forward_naive(x, w, b, {'stride': 1, 'pad': 1})

def imshow_noax(img, normalize=True):
    """ Tiny helper to show images as uint8 and remove axis labels """
    if normalize:
        img_max, img_min = np.max(img), np.min(img)
        img = 255.0 * (img - img_min) / (img_max - img_min)
        plt.imshow(img.astype('uint8'))
        plt.gca().axis('off')

# Show the original images and the results of the conv operation
plt.subplot(2, 3, 1)
imshow_noax(puppy, normalize=False)
plt.title('Original image')
plt.subplot(2, 3, 2)
imshow_noax(out[0, 0])
plt.title('Grayscale')
plt.subplot(2, 3, 3)
imshow_noax(out[0, 1])
plt.title('Edges')
plt.subplot(2, 3, 4)
imshow_noax(kitten_cropped, normalize=False)
plt.subplot(2, 3, 5)
imshow_noax(out[1, 0])
plt.subplot(2, 3, 6)

```

```
imshow_noax(out[1, 1])
plt.show()
```

C:\Users\kucharskib\AppData\Roaming\Python\Python36\site-packages\ipykernel_launcher.py:3: DeprecationWarning: `imread` is deprecated!

`imread` is deprecated in SciPy 1.0.0, and will be removed in 1.2.0.

Use ``imageio.imread`` instead.

This is separate from the ipykernel package so we can avoid doing imports until

C:\Users\kucharskib\AppData\Roaming\Python\Python36\site-packages\ipykernel_launcher.py:10: DeprecationWarning: `imresize` is deprecated!

`imresize` is deprecated in SciPy 1.0.0, and will be removed in 1.2.0.

Use ``skimage.transform.resize`` instead.

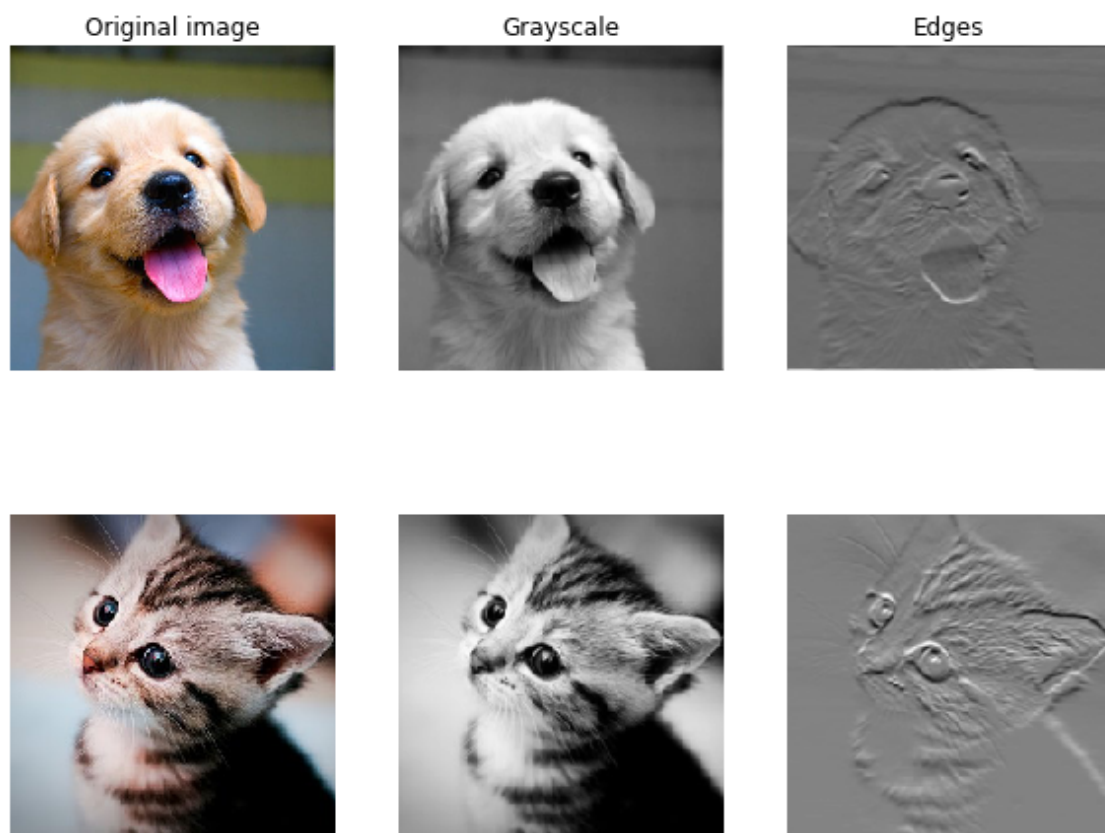
Remove the CWD from sys.path while we load stuff.

C:\Users\kucharskib\AppData\Roaming\Python\Python36\site-packages\ipykernel_launcher.py:11: DeprecationWarning: `imresize` is deprecated!

`imresize` is deprecated in SciPy 1.0.0, and will be removed in 1.2.0.

Use ``skimage.transform.resize`` instead.

This is added back by InteractiveShellApp.init_path()



Convolution: Naive backward pass

Implement the backward pass for the convolution operation in the function `conv_backward_naive` in the file `cs682/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 [108]: 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: 3.37264006649648e-11

```

Max-Pooling: Naive forward

Implement the forward pass for the max-pooling operation in the function `max_pool_forward_naive` in the file `cs682/layers.py`. Again, don't worry too much about computational efficiency.

Check your implementation by running the following:

```
In [109]: 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 `cs682/layers.py`. You don't need to worry about computational efficiency.

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

```
In [110]: 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 `cs682/fast_layers.py`.

The fast convolution implementation depends on a Cython extension; to compile it you need to run the following from the `cs682` 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 [111]: # Rel errors should be around e-9 or less
from cs682.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('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: 3.849738s
Fast: 0.022903s
Speedup: 168.090828x
Difference: 4.926407851494105e-11

```

```

Testing conv_backward_fast:
Naive: 7.799168s
Fast: 0.015932s
Speedup: 489.540570x
dx difference: 1.949764775345631e-11
dw difference: 5.188375174206562e-13
db difference: 3.481354613192702e-14

```

```
In [112]: # Relative errors should be close to 0.0
from cs682.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: 0.330285s
fast: 0.006981s
speedup: 47.309508x
difference: 0.0
```

```
Testing pool_backward_fast:
Naive: 0.805879s
fast: 0.015952s
speedup: 50.517897x
dx difference: 0.0
```

Convolutional "sandwich" layers

Previously we introduced the concept of "sandwich" layers that combine multiple operations into commonly used patterns. In the file `cs682/layer_utils.py` you will find sandwich layers that implement a few commonly used patterns for convolutional networks.

```
In [113]: from cs682.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: 5.828178746516271e-09
dw error: 8.443628091870788e-09
db error: 3.57960501324485e-10
```

```
In [114]: from cs682.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 `cs682/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 this should go up.

```
In [115]: 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): 2.7122622217792203
```

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 [116]: 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 [117]: 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)

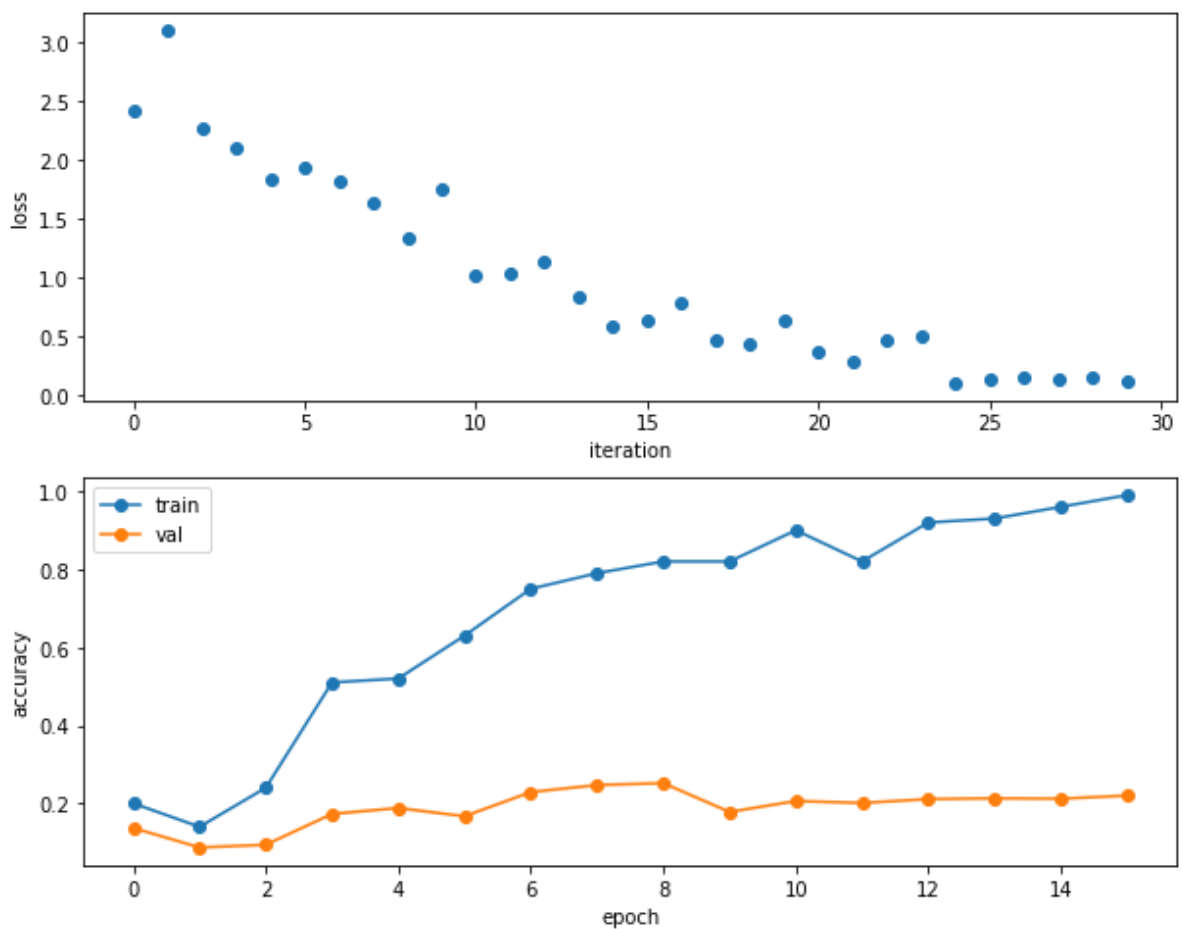
solver = Solver(model, small_data,
                num_epochs=15, batch_size=50,
                update_rule='adam',
                optim_config={
                    'learning_rate': 1e-3,
                },
                verbose=False, print_every=1)
solver.train()

```

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

```
In [118]: 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 [121]: 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()
```

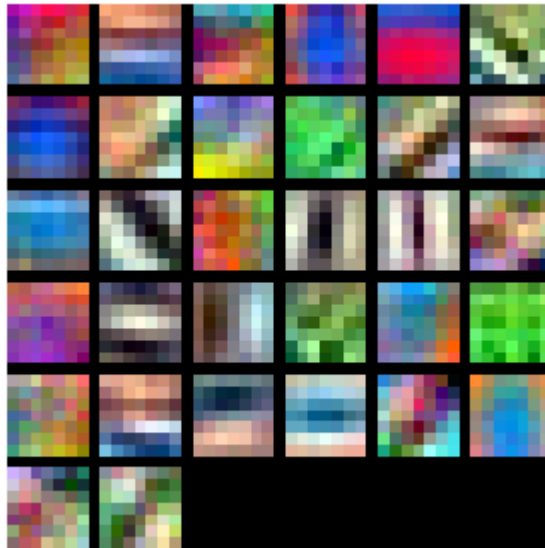
(Iteration 1 / 980) loss: 2.306477
(Epoch 0 / 1) train acc: 0.113000; val_acc: 0.119000
(Iteration 21 / 980) loss: 2.159847
(Iteration 41 / 980) loss: 2.321392
(Iteration 61 / 980) loss: 1.954295
(Iteration 81 / 980) loss: 1.923235
(Iteration 101 / 980) loss: 2.070575
(Iteration 121 / 980) loss: 1.749879
(Iteration 141 / 980) loss: 1.989133
(Iteration 161 / 980) loss: 1.840851
(Iteration 181 / 980) loss: 1.719386
(Iteration 201 / 980) loss: 2.036195
(Iteration 221 / 980) loss: 2.194792
(Iteration 241 / 980) loss: 1.786065
(Iteration 261 / 980) loss: 1.939395
(Iteration 281 / 980) loss: 1.660571
(Iteration 301 / 980) loss: 2.130397
(Iteration 321 / 980) loss: 1.918348
(Iteration 341 / 980) loss: 1.738603
(Iteration 361 / 980) loss: 1.694989
(Iteration 381 / 980) loss: 1.856295
(Iteration 401 / 980) loss: 1.448928
(Iteration 421 / 980) loss: 1.738722
(Iteration 441 / 980) loss: 1.725295
(Iteration 461 / 980) loss: 1.512955
(Iteration 481 / 980) loss: 1.753725
(Iteration 501 / 980) loss: 1.683960
(Iteration 521 / 980) loss: 1.602517
(Iteration 541 / 980) loss: 1.649211
(Iteration 561 / 980) loss: 1.476181
(Iteration 581 / 980) loss: 1.851032
(Iteration 601 / 980) loss: 1.949298
(Iteration 621 / 980) loss: 1.651621
(Iteration 641 / 980) loss: 1.566318
(Iteration 661 / 980) loss: 2.102920
(Iteration 681 / 980) loss: 1.563032
(Iteration 701 / 980) loss: 1.509992
(Iteration 721 / 980) loss: 1.620093
(Iteration 741 / 980) loss: 1.737529
(Iteration 761 / 980) loss: 1.515969
(Iteration 781 / 980) loss: 1.496814
(Iteration 801 / 980) loss: 1.546206
(Iteration 821 / 980) loss: 1.928535
(Iteration 841 / 980) loss: 1.346990
(Iteration 861 / 980) loss: 1.646753
(Iteration 881 / 980) loss: 1.736880
(Iteration 901 / 980) loss: 1.695742
(Iteration 921 / 980) loss: 1.672395
(Iteration 941 / 980) loss: 1.578137
(Iteration 961 / 980) loss: 1.703719
(Epoch 1 / 1) train acc: 0.448000; val_acc: 0.438000

Visualize Filters

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

```
In [122]: from cs682.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 [3], 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 the statistics of each feature channel to be relatively consistent both between different images and different locations within the same image. Therefore spatial batch normalization computes a mean and variance for each of the C feature channels by computing statistics over both the minibatch dimension N and the spatial dimensions H and W .

[3] Sergey Ioffe and Christian Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift", ICML 2015. (<https://arxiv.org/abs/1502.03167>)

Spatial batch normalization: forward

In the file `cs682/layers.py`, implement the forward pass for spatial batch normalization in the function `spatial_batchnorm_forward`. Check your implementation by running the following:

```
In [123]: 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: [1.38777878e-16 1.94289029e-17 -9.43689571e-17]

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 [124]: 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 `cs682/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 [125]: 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: 3.423838616073709e-07
dgamma error: 7.0963199356067174e-12
dbeta error: 3.275380797385891e-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 [4] 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 [5] 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.



Comparison of normalization techniques discussed so far

****Visual comparison of the normalization techniques discussed so far (image edited from [5])****

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 [6]-- 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 arXiv *less than a month ago* -- this truly is still an ongoing and excitingly active field of research!

[4] Ba, Jimmy Lei, Jamie Ryan Kiros, and Geoffrey E. Hinton. "Layer Normalization." *stat* 1050 (2016): 21. (<https://arxiv.org/pdf/1607.06450.pdf>).

[5] Wu, Yuxin, and Kaiming He. "Group Normalization." *arXiv preprint arXiv:1803.08494* (2018). (<https://arxiv.org/abs/1803.08494>)

[6] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *Computer Vision and Pattern Recognition (CVPR)*, 2005. (<https://ieeexplore.ieee.org/abstract/document/1467360/>).

Group normalization: forward

In the file `cs682/layers.py`, implement the forward pass for group normalization in the function `spatial_groupnorm_forward`. Check your implementation by running the following:

```

In [126]: 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 `cs682/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 [127]: 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: 7.413109437563619e-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 switch over to that notebook).

What is PyTorch?

PyTorch is a system for executing dynamic computational graphs over Tensor objects that behave similarly as numpy ndarray. It comes with a powerful automatic differentiation engine that removes the need for manual back-propagation.

Why?

- Our code will now run on GPUs! Much faster training. When using a framework like PyTorch or TensorFlow you can harness the power of the GPU for your own custom neural network architectures without having to write CUDA code directly (which is beyond the scope of this class).
- We want you to be ready to use one of these frameworks for your project so you can experiment more efficiently than if you were writing every feature you want to use by hand.
- We want you to stand on the shoulders of giants! TensorFlow and PyTorch are both excellent frameworks that will make your lives a lot easier, and now that you understand their guts, you are free to use them :)
- We want you to be exposed to the sort of deep learning code you might run into in academia or industry.

PyTorch versions

This notebook assumes that you are using **PyTorch version 0.4**. Prior to this version, Tensors had to be wrapped in Variable objects to be used in autograd; however Variables have now been deprecated. In addition 0.4 also separates a Tensor's datatype from its device, and uses numpy-style factories for constructing Tensors rather than directly invoking Tensor constructors.

How will I learn PyTorch?

Justin Johnson has made an excellent [tutorial](https://github.com/jcjohnson/pytorch-examples) (<https://github.com/jcjohnson/pytorch-examples>) for PyTorch.

You can also find the detailed [API doc](http://pytorch.org/docs/stable/index.html) (<http://pytorch.org/docs/stable/index.html>) here. If you have other questions that are not addressed by the API docs, the [PyTorch forum](https://discuss.pytorch.org/) (<https://discuss.pytorch.org/>) is a much better place to ask than StackOverflow.

Table of Contents

This assignment has 5 parts. You will learn PyTorch on different levels of abstractions, which will help you understand it better and prepare you for the final project.

1. Preparation: we will use CIFAR-10 dataset.
2. Barebones PyTorch: we will work directly with the lowest-level PyTorch Tensors.
3. PyTorch Module API: we will use `nn.Module` to define arbitrary neural network architecture.
4. PyTorch Sequential API: we will use `nn.Sequential` to define a linear feed-forward network very conveniently.
5. CIFAR-10 open-ended challenge: please implement your own network to get as high accuracy as possible on CIFAR-10. You can experiment with any layer, optimizer, hyperparameters or other advanced features.

Here is a table of comparison:

API	Flexibility	Convenience
Barebone	High	Low
<code>nn.Module</code>	High	Medium
<code>nn.Sequential</code>	Low	High

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.

```
In [91]: import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader
from torch.utils.data import sampler
import torchvision
import torchvision.datasets as dset
import torchvision.transforms as T

import numpy as np
```

```
In [92]: NUM_TRAIN = 49000

# The torchvision.transforms package provides tools for preprocessing data
# and for performing data augmentation; here we set up a transform to
# preprocess the data by subtracting the mean RGB value and dividing by the
# standard deviation of each RGB value; we've hardcoded the mean and std.
transform = T.Compose([
    T.ToTensor(),
    T.Normalize((0.4914, 0.4822, 0.4465), (0.2023, 0.1994, 0.2010))
])

# We set up a Dataset object for each split (train / val / test); Datasets load
# training examples one at a time, so we wrap each Dataset in a DataLoader which
# iterates through the Dataset and forms minibatches. We divide the CIFAR-10
# training set into train and val sets by passing a Sampler object to the
# DataLoader telling how it should sample from the underlying Dataset.
cifar10_train = dset.CIFAR10('./cs682/datasets', train=True, download=True,
                             transform=transform)
loader_train = DataLoader(cifar10_train, batch_size=64,
                          sampler=sampler.SubsetRandomSampler(range(NUM_TRAIN)))

cifar10_val = dset.CIFAR10('./cs682/datasets', train=True, download=True,
                           transform=transform)
loader_val = DataLoader(cifar10_val, batch_size=64,
                       sampler=sampler.SubsetRandomSampler(range(NUM_TRAIN, 50000)))

cifar10_test = dset.CIFAR10('./cs682/datasets', train=False, download=True,
                             transform=transform)
loader_test = DataLoader(cifar10_test, batch_size=64)
```

Files already downloaded and verified
Files already downloaded and verified
Files already downloaded and verified

You have an option to **use GPU by setting the flag to True below**. It is not necessary to use GPU for this assignment. Note that if your computer does not have CUDA enabled, `torch.cuda.is_available()` will return False and this notebook will fallback to CPU mode.

The global variables `dtype` and `device` will control the data types throughout this assignment.

```
In [93]: USE_GPU = True

dtype = torch.float32 # we will be using float throughout this tutorial

if USE_GPU and torch.cuda.is_available():
    device = torch.device('cuda')
else:
    device = torch.device('cpu')

# Constant to control how frequently we print train loss
print_every = 100

print('using device:', device)

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.

PyTorch Tensors: Flatten Function

A PyTorch Tensor is conceptionally similar to a numpy array: it is an n-dimensional grid of numbers, and like numpy PyTorch provides many functions to efficiently operate on Tensors. As a simple example, we provide a `flatten` function below which reshapes image data for use in a fully-connected neural network.

Recall that image data is typically stored in a Tensor of shape $N \times C \times H \times W$, where:

- N is the number of datapoints
- C is the number of channels
- H is the height of the intermediate feature map in pixels
- W is the width of the intermediate feature map in pixels

This is the right way to represent the data when we are doing something like a 2D convolution, that needs spatial understanding of where the intermediate features are relative to each other. When we use fully connected affine layers to process the image, however, we want each datapoint to be represented by a single vector -- it's no longer useful to segregate the different channels, rows, and columns of the data. So, we use a "flatten" operation to collapse the $C \times H \times W$ values per representation into a single long vector. The `flatten` function below first reads in the N , C , H , and W values from a given batch of data, and then returns a "view" of that data. "View" is analogous to numpy's "reshape" method: it reshapes x's dimensions to be $N \times ??$, where $??$ is allowed to be anything (in this case, it will be $C \times H \times W$, but we don't need to specify that explicitly).

```
In [94]: def flatten(x):
          N = x.shape[0] # read in N, C, H, W
          return x.view(N, -1) # "flatten" the C * H * W values into a single vector per image

          def test_flatten():
              x = torch.arange(12).view(2, 1, 3, 2)
              print('Before flattening: ', x)
              print('After flattening: ', flatten(x))

          test_flatten()
```

```
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]])
```

Barebones PyTorch: Two-Layer Network

Here we define a function `two_layer_fc` which performs the forward pass of a two-layer fully-connected ReLU network on a batch of image data. After defining the forward pass we check that it doesn't crash and that it produces outputs of the right shape by running zeros through the network.

You don't have to write any code here, but it's important that you read and understand the implementation.

In [95]: **import torch.nn.functional as F** *# useful stateless functions*

def two_layer_fc(x, params):

"""

*A fully-connected neural networks; the architecture is:
 NN is fully connected -> ReLU -> fully connected layer.
 Note that this function only defines the forward pass;
 PyTorch will take care of the backward pass for us.*

*The input to the network will be a minibatch of data, of shape
 (N, d1, ..., dM) where $d1 * \dots * dM = D$. The hidden layer will have H units,
 and the output layer will produce scores for C classes.*

Inputs:

- x: A PyTorch Tensor of shape (N, d1, ..., dM) giving a minibatch of input data.
- params: A list [w1, w2] of PyTorch Tensors giving weights for the network; w1 has shape (D, H) and w2 has shape (H, C).

Returns:

- scores: A PyTorch Tensor of shape (N, C) giving classification scores for the input data x.

"""

first we flatten the image

x = flatten(x) *# shape: [batch_size, C x H x W]*

w1, w2 = params

Forward pass: compute predicted y using operations on Tensors. Since w1 and w2 have requires_grad=True, operations involving these Tensors will cause PyTorch to build a computational graph, allowing automatic computation of gradients. Since we are no longer implementing the backward pass by hand we don't need to keep references to intermediate values.

you can also use .clamp(min=0), equivalent to F.relu()

x = F.relu(x.mm(w1))

x = x.mm(w2)

return x

def two_layer_fc_test():

hidden_layer_size = 42

x = torch.zeros((64, 50), dtype=dtype) *# minibatch size 64, feature dimension 50*

w1 = torch.zeros((50, hidden_layer_size), dtype=dtype)

w2 = torch.zeros((hidden_layer_size, 10), dtype=dtype)

scores = two_layer_fc(x, [w1, w2])

print(scores.size()) *# you should see [64, 10]*

two_layer_fc_test()

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

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

In [96]: **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.

#####

z1 = F.conv2d(x, conv_w1, padding=2, bias = conv_b1)

a1 = F.relu(z1)

z2 = F.conv2d(a1, conv_w2, padding=1, bias = conv_b2)

a2 = F.relu(z2)

scores = flatten(a2).mm(fc_w) + fc_b

#####

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).


```
In [97]: def three_layer_convnet_test():
        x = torch.zeros((64, 3, 32, 32), dtype=dtype) # minibatch size 64, image size [3, 32, 32]

        conv_w1 = torch.zeros((6, 3, 5, 5), dtype=dtype) # [out_channel, in_channel, kernel_H, kernel_W]
        conv_b1 = torch.zeros((6,)) # out_channel
        conv_w2 = torch.zeros((9, 6, 3, 3), dtype=dtype) # [out_channel, in_channel, kernel_H, kernel_W]
        conv_b2 = torch.zeros((9,)) # out_channel

        # you must calculate the shape of the tensor after two conv layers, before the fully-connected layer
        fc_w = torch.zeros((9 * 32 * 32, 10))
        fc_b = torch.zeros(10)

        scores = three_layer_convnet(x, [conv_w1, conv_b1, conv_w2, conv_b2, fc_w, fc_b])
        print(scores.size()) # you should see [64, 10]
        three_layer_convnet_test()

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:

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

```
In [98]: def random_weight(shape):
        """
        Create random Tensors for weights; setting requires_grad=True means that we
        want to compute gradients for these Tensors during the backward pass.
        We use Kaiming normalization: sqrt(2 / fan_in)
        """
        if len(shape) == 2: # FC weight
            fan_in = shape[0]
        else:
            fan_in = np.prod(shape[1:]) # conv weight [out_channel, in_channel, kH, kW]
            # randn is standard normal distribution generator.
        w = torch.randn(shape, device=device, dtype=dtype) * np.sqrt(2. / fan_in)
        w.requires_grad = True
        return w

def zero_weight(shape):
    return torch.zeros(shape, device=device, dtype=dtype, requires_grad=True)

# create a weight of shape [3 x 5]
# you should see the type `torch.cuda.FloatTensor` if you use GPU.
# Otherwise it should be `torch.FloatTensor`
random_weight((3, 5))
```

```
Out[98]: tensor([[ 0.4842,  1.1535,  0.1829, -0.2879, -0.5325],
                 [-0.6539,  0.2046, -1.5711,  0.0469, -0.7790],
                 [ 0.1059, -0.5773,  1.2528,  0.0794,  0.3104]],
                 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.

```
In [99]: def check_accuracy_part2(loader, model_fn, params):
        """
        Check the accuracy of a classification model.

        Inputs:
        - loader: A DataLoader for the data split we want to check
        - model_fn: A function that performs the forward pass of the model,
          with the signature scores = model_fn(x, params)
        - params: List of PyTorch Tensors giving parameters of the model

        Returns: Nothing, but prints the accuracy of the model
        """
        split = 'val' if loader.dataset.train else 'test'
        print('Checking accuracy on the %s set' % split)
        num_correct, num_samples = 0, 0
        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.int64)
                scores = model_fn(x, params)
                _, 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))
```

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](http://pytorch.org/docs/stable/nn.html#cross-entropy) (<http://pytorch.org/docs/stable/nn.html#cross-entropy>).

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

```

In [100]: def train_part2(model_fn, params, learning_rate):
          """
          Train a model on CIFAR-10.

          Inputs:
          - model_fn: A Python function that performs the forward pass of the model.
            It should have the signature scores = model_fn(x, params) where x is a
            PyTorch Tensor of image data, params is a list of PyTorch Tensors giving
            model weights, and scores is a PyTorch Tensor of shape (N, C) giving
            scores for the elements in x.
          - params: List of PyTorch Tensors giving weights for the model
          - learning_rate: Python scalar giving the learning rate to use for SGD

          Returns: Nothing
          """
          for t, (x, y) in enumerate(loader_train):
              # Move the data to the proper device (GPU or CPU)
              x = x.to(device=device, dtype=dtype)
              y = y.to(device=device, dtype=torch.long)

              # Forward pass: compute scores and loss
              scores = model_fn(x, params)
              loss = F.cross_entropy(scores, y)

              # Backward pass: PyTorch figures out which Tensors in the computational
              # graph has requires_grad=True and uses backpropagation to compute the
              # gradient of the loss with respect to these Tensors, and stores the
              # gradients in the .grad attribute of each Tensor.
              loss.backward()

              # Update parameters. We don't want to backpropagate through the
              # parameter updates, so we scope the updates under a torch.no_grad()
              # context manager to prevent a computational graph from being built.
              with torch.no_grad():
                  for w in params:
                      w -= learning_rate * w.grad

                      # Manually zero the gradients after running the backward pass
                      w.grad.zero_()

              if t % print_every == 0:
                  print('Iteration %d, loss = %.4f' % (t, loss.item()))
                  check_accuracy_part2(loader_val, model_fn, params)
                  print()

```

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, w_1 and w_2 .

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 w_1 . The second dimension of w_1 is the hidden layer size, which will also be the first dimension of w_2 .

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 [101]: 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 = 3.4224
Checking accuracy on the val set
Got 151 / 1000 correct (15.10%)

Iteration 100, loss = 2.4432
Checking accuracy on the val set
Got 335 / 1000 correct (33.50%)

Iteration 200, loss = 1.9713
Checking accuracy on the val set
Got 365 / 1000 correct (36.50%)

Iteration 300, loss = 1.8364
Checking accuracy on the val set
Got 433 / 1000 correct (43.30%)

Iteration 400, loss = 2.0489
Checking accuracy on the val set
Got 422 / 1000 correct (42.20%)

Iteration 500, loss = 1.8577
Checking accuracy on the val set
Got 418 / 1000 correct (41.80%)

Iteration 600, loss = 2.0871
Checking accuracy on the val set
Got 461 / 1000 correct (46.10%)

Iteration 700, loss = 1.7865
Checking accuracy on the val set
Got 456 / 1000 correct (45.60%)

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 [102]: 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.          #
#####
conv_w1 = random_weight((channel_1, 3, 5, 5)) # [out_channel, in_channel, kernel_H, kernel_W]
conv_b1 = zero_weight((32,)) # out_channel
conv_w2 = random_weight((channel_2, channel_1, 3, 3)) # [out_channel, in_channel, kernel_H, kernel_W]
conv_b2 = zero_weight((16,)) # out_channel
fc_w = random_weight((channel_2 * 32 * 32, 10))
fc_b = zero_weight((10,))

#####
#                               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 = 3.0632
Checking accuracy on the val set
Got 91 / 1000 correct (9.10%)

Iteration 100, loss = 1.9420
Checking accuracy on the val set
Got 344 / 1000 correct (34.40%)

Iteration 200, loss = 1.8680
Checking accuracy on the val set
Got 399 / 1000 correct (39.90%)

Iteration 300, loss = 1.7798
Checking accuracy on the val set
Got 415 / 1000 correct (41.50%)

Iteration 400, loss = 1.7605
Checking accuracy on the val set
Got 448 / 1000 correct (44.80%)

Iteration 500, loss = 1.5281
Checking accuracy on the val set
Got 456 / 1000 correct (45.60%)

Iteration 600, loss = 1.4988
Checking accuracy on the val set
Got 459 / 1000 correct (45.90%)

Iteration 700, loss = 1.4402
Checking accuracy on the val set
Got 456 / 1000 correct (45.60%)

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](http://pytorch.org/docs/master/optim.html) (<http://pytorch.org/docs/master/optim.html>) 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](http://pytorch.org/docs/master/nn.html) (<http://pytorch.org/docs/master/nn.html>) 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 [103]: 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> (<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 [104]: 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. #
    #####
    self.conv2d_1 = nn.Conv2d(in_channel, channel_1, 5, padding = 2)
    nn.init.kaiming_normal_(self.conv2d_1.weight)
    self.conv2d_2 = nn.Conv2d(channel_1, channel_2, 3, padding = 1)
    nn.init.kaiming_normal_(self.conv2d_2.weight)
    self.fc1 = nn.Linear(channel_2 * 32 * 32, num_classes)
    nn.init.kaiming_normal_(self.fc1.weight)

    #####
    #                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() #
    #####
    z1 = self.conv2d_1(x)
    a1 = F.relu(z1)
    z2 = self.conv2d_2(a1)
    a2 = F.relu(z2)
    scores = self.fc1(flatten(a2))
    #####
    #                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 [105]: 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 [106]: 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):
              print("current epochL: " + str(e))
              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.
                  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 [107]: 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)
```

```
current epochL: 0
Iteration 0, loss = 3.4510
Checking accuracy on validation set
Got 142 / 1000 correct (14.20)
```

```
Iteration 100, loss = 2.5083
Checking accuracy on validation set
Got 312 / 1000 correct (31.20)
```

```
Iteration 200, loss = 2.0687
Checking accuracy on validation set
Got 350 / 1000 correct (35.00)
```

```
Iteration 300, loss = 1.8207
Checking accuracy on validation set
Got 383 / 1000 correct (38.30)
```

```
Iteration 400, loss = 1.4958
Checking accuracy on validation set
Got 428 / 1000 correct (42.80)
```

```
Iteration 500, loss = 1.6498
Checking accuracy on validation set
Got 401 / 1000 correct (40.10)
```

```
Iteration 600, loss = 1.5736
Checking accuracy on validation set
Got 427 / 1000 correct (42.70)
```

```
Iteration 700, loss = 1.6715
Checking accuracy on validation set
Got 424 / 1000 correct (42.40)
```

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 [108]: learning_rate = 3e-3
          channel_1 = 32
          channel_2 = 16

          model = None
          optimizer = None
          #####
          # TODO: Instantiate your ThreeLayerConvNet model and a corresponding optimizer #
          #####
          model = ThreeLayerConvNet(3, channel_1, channel_2, 10)
          optimizer = optim.SGD(model.parameters(), lr=learning_rate)
          #####
          #                                     END OF YOUR CODE
          #####

          train_part34(model, optimizer)

```

current epochL: 0
 Iteration 0, loss = 2.7229
 Checking accuracy on validation set
 Got 103 / 1000 correct (10.30)

Iteration 100, loss = 1.9029
 Checking accuracy on validation set
 Got 282 / 1000 correct (28.20)

Iteration 200, loss = 1.6833
 Checking accuracy on validation set
 Got 381 / 1000 correct (38.10)

Iteration 300, loss = 1.6309
 Checking accuracy on validation set
 Got 401 / 1000 correct (40.10)

Iteration 400, loss = 1.7831
 Checking accuracy on validation set
 Got 429 / 1000 correct (42.90)

Iteration 500, loss = 1.5737
 Checking accuracy on validation set
 Got 460 / 1000 correct (46.00)

Iteration 600, loss = 1.7391
 Checking accuracy on validation set
 Got 469 / 1000 correct (46.90)

Iteration 700, loss = 1.6367
 Checking accuracy on validation set
 Got 464 / 1000 correct (46.40)

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 [109]: # 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)
```

```
current epochL: 0
Iteration 0, loss = 2.3330
Checking accuracy on validation set
Got 141 / 1000 correct (14.10)
```

```
Iteration 100, loss = 1.6306
Checking accuracy on validation set
Got 394 / 1000 correct (39.40)
```

```
Iteration 200, loss = 2.2703
Checking accuracy on validation set
Got 400 / 1000 correct (40.00)
```

```
Iteration 300, loss = 1.7348
Checking accuracy on validation set
Got 430 / 1000 correct (43.00)
```

```
Iteration 400, loss = 1.7541
Checking accuracy on validation set
Got 410 / 1000 correct (41.00)
```

```
Iteration 500, loss = 1.3861
Checking accuracy on validation set
Got 453 / 1000 correct (45.30)
```

```
Iteration 600, loss = 1.5783
Checking accuracy on validation set
Got 437 / 1000 correct (43.70)
```

```
Iteration 700, loss = 1.4816
Checking accuracy on validation set
Got 417 / 1000 correct (41.70)
```

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 [110]: 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. #
#####
model = nn.Sequential(
    nn.Conv2d(3,channel_1,5,padding = 2, bias=True),
    nn.ReLU(),
    nn.Conv2d(channel_1,channel_2,3,padding = 1, bias=True),
    nn.ReLU(),
    Flatten(),
    nn.Linear(channel_2* 32 * 32, 10) ,
)

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

train_part34(model, optimizer)

```

current epochL: 0
Iteration 0, loss = 2.3237
Checking accuracy on validation set
Got 131 / 1000 correct (13.10)

Iteration 100, loss = 1.8639
Checking accuracy on validation set
Got 474 / 1000 correct (47.40)

Iteration 200, loss = 1.2983
Checking accuracy on validation set
Got 481 / 1000 correct (48.10)

Iteration 300, loss = 1.2639
Checking accuracy on validation set
Got 514 / 1000 correct (51.40)

Iteration 400, loss = 1.4102
Checking accuracy on validation set
Got 532 / 1000 correct (53.20)

Iteration 500, loss = 1.1298
Checking accuracy on validation set
Got 555 / 1000 correct (55.50)

Iteration 600, loss = 1.2510
Checking accuracy on validation set
Got 568 / 1000 correct (56.80)

Iteration 700, loss = 1.1365
Checking accuracy on validation set
Got 577 / 1000 correct (57.70)

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> (<http://pytorch.org/docs/stable/nn.html>)
- Activations: <http://pytorch.org/docs/stable/nn.html#non-linear-activations> (<http://pytorch.org/docs/stable/nn.html#non-linear-activations>)
- Loss functions: <http://pytorch.org/docs/stable/nn.html#loss-functions> (<http://pytorch.org/docs/stable/nn.html#loss-functions>)
- Optimizers: <http://pytorch.org/docs/stable/optim.html> (<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 \(https://arxiv.org/abs/1512.00567\)](https://arxiv.org/abs/1512.00567). (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 (<https://arxiv.org/abs/1512.03385>) where the input from the previous layer is added to the output.
 - DenseNets (<https://arxiv.org/abs/1608.06993>) where inputs into previous layers are concatenated together.
 - This blog has an in-depth overview (<https://chatbotlife.com/resnets-highwaynets-and-densenets-oh-my-9bb15918ee32>)

Have fun and happy training!

```

In [111]: #####
# 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

model = nn.Sequential(
    nn.Conv2d(3, 64, kernel_size=3, stride=2, padding=1),
    nn.ReLU(inplace=True),
    nn.MaxPool2d(kernel_size=2),
    nn.Conv2d(64, 192, kernel_size=3, padding=1),
    nn.ReLU(inplace=True),
    nn.MaxPool2d(kernel_size=2),
    nn.Conv2d(192, 384, kernel_size=3, padding=1),
    nn.ReLU(inplace=True),
    nn.Conv2d(384, 256, kernel_size=3, padding=1),
    nn.ReLU(inplace=True),
    nn.Conv2d(256, 256, kernel_size=3, padding=1),
    nn.ReLU(inplace=True),
    nn.MaxPool2d(kernel_size=2),
    Flatten(),
    nn.Dropout(),
    nn.Linear(256 * 2 * 2, 4096),
    nn.ReLU(inplace=True),
    nn.Dropout(),
    nn.Linear(4096, 4096),
    nn.ReLU(inplace=True),
    nn.Linear(4096, 10),
)

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

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

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

```


current epochL: 0
Iteration 0, loss = 2.3010
Checking accuracy on validation set
Got 71 / 1000 correct (7.10)

Iteration 100, loss = 2.2993
Checking accuracy on validation set
Got 132 / 1000 correct (13.20)

Iteration 200, loss = 2.2895
Checking accuracy on validation set
Got 179 / 1000 correct (17.90)

Iteration 300, loss = 2.0655
Checking accuracy on validation set
Got 184 / 1000 correct (18.40)

Iteration 400, loss = 1.8086
Checking accuracy on validation set
Got 283 / 1000 correct (28.30)

Iteration 500, loss = 1.7071
Checking accuracy on validation set
Got 359 / 1000 correct (35.90)

Iteration 600, loss = 1.6815
Checking accuracy on validation set
Got 390 / 1000 correct (39.00)

Iteration 700, loss = 1.5303
Checking accuracy on validation set
Got 417 / 1000 correct (41.70)

current epochL: 1
Iteration 0, loss = 1.3924
Checking accuracy on validation set
Got 447 / 1000 correct (44.70)

Iteration 100, loss = 1.4552
Checking accuracy on validation set
Got 462 / 1000 correct (46.20)

Iteration 200, loss = 1.4698
Checking accuracy on validation set
Got 491 / 1000 correct (49.10)

Iteration 300, loss = 1.4218
Checking accuracy on validation set
Got 531 / 1000 correct (53.10)

Iteration 400, loss = 1.1817
Checking accuracy on validation set
Got 513 / 1000 correct (51.30)

Iteration 500, loss = 1.0783
Checking accuracy on validation set
Got 578 / 1000 correct (57.80)

Iteration 600, loss = 1.3941
Checking accuracy on validation set
Got 593 / 1000 correct (59.30)

Iteration 700, loss = 1.2384
Checking accuracy on validation set
Got 587 / 1000 correct (58.70)

current epochL: 2
Iteration 0, loss = 1.3248
Checking accuracy on validation set
Got 571 / 1000 correct (57.10)

Iteration 100, loss = 1.1010
Checking accuracy on validation set
Got 621 / 1000 correct (62.10)

Iteration 200, loss = 1.2086
Checking accuracy on validation set
Got 643 / 1000 correct (64.30)

Iteration 300, loss = 1.1690
Checking accuracy on validation set
Got 642 / 1000 correct (64.20)

Iteration 400, loss = 0.9294
Checking accuracy on validation set
Got 648 / 1000 correct (64.80)

Iteration 500, loss = 0.9897
Checking accuracy on validation set
Got 658 / 1000 correct (65.80)

Iteration 600, loss = 1.1792
Checking accuracy on validation set
Got 656 / 1000 correct (65.60)

Iteration 700, loss = 0.9562
Checking accuracy on validation set
Got 667 / 1000 correct (66.70)

current epochL: 3
Iteration 0, loss = 1.2297
Checking accuracy on validation set
Got 669 / 1000 correct (66.90)

Iteration 100, loss = 0.9282
Checking accuracy on validation set
Got 670 / 1000 correct (67.00)

Iteration 200, loss = 0.7887
Checking accuracy on validation set
Got 683 / 1000 correct (68.30)

Iteration 300, loss = 0.8488
Checking accuracy on validation set

Got 663 / 1000 correct (66.30)

Iteration 400, loss = 0.9115
Checking accuracy on validation set
Got 675 / 1000 correct (67.50)

Iteration 500, loss = 0.9091
Checking accuracy on validation set
Got 697 / 1000 correct (69.70)

Iteration 600, loss = 0.7985
Checking accuracy on validation set
Got 704 / 1000 correct (70.40)

Iteration 700, loss = 0.8518
Checking accuracy on validation set
Got 690 / 1000 correct (69.00)

current epochL: 4
Iteration 0, loss = 0.5638
Checking accuracy on validation set
Got 732 / 1000 correct (73.20)

Iteration 100, loss = 0.6395
Checking accuracy on validation set
Got 718 / 1000 correct (71.80)

Iteration 200, loss = 0.6961
Checking accuracy on validation set
Got 711 / 1000 correct (71.10)

Iteration 300, loss = 0.4976
Checking accuracy on validation set
Got 688 / 1000 correct (68.80)

Iteration 400, loss = 0.7171
Checking accuracy on validation set
Got 713 / 1000 correct (71.30)

Iteration 500, loss = 0.6211
Checking accuracy on validation set
Got 709 / 1000 correct (70.90)

Iteration 600, loss = 0.6227
Checking accuracy on validation set
Got 717 / 1000 correct (71.70)

Iteration 700, loss = 0.8699
Checking accuracy on validation set
Got 720 / 1000 correct (72.00)

current epochL: 5
Iteration 0, loss = 0.4783
Checking accuracy on validation set
Got 710 / 1000 correct (71.00)

Iteration 100, loss = 0.5556

Checking accuracy on validation set
Got 739 / 1000 correct (73.90)

Iteration 200, loss = 0.3349
Checking accuracy on validation set
Got 727 / 1000 correct (72.70)

Iteration 300, loss = 0.7452
Checking accuracy on validation set
Got 723 / 1000 correct (72.30)

Iteration 400, loss = 0.5503
Checking accuracy on validation set
Got 728 / 1000 correct (72.80)

Iteration 500, loss = 0.5096
Checking accuracy on validation set
Got 728 / 1000 correct (72.80)

Iteration 600, loss = 0.3231
Checking accuracy on validation set
Got 737 / 1000 correct (73.70)

Iteration 700, loss = 0.4366
Checking accuracy on validation set
Got 749 / 1000 correct (74.90)

current epochL: 6
Iteration 0, loss = 0.5000
Checking accuracy on validation set
Got 759 / 1000 correct (75.90)

Iteration 100, loss = 0.7219
Checking accuracy on validation set
Got 746 / 1000 correct (74.60)

Iteration 200, loss = 0.3430
Checking accuracy on validation set
Got 752 / 1000 correct (75.20)

Iteration 300, loss = 0.3824
Checking accuracy on validation set
Got 731 / 1000 correct (73.10)

Iteration 400, loss = 0.5620
Checking accuracy on validation set
Got 736 / 1000 correct (73.60)

Iteration 500, loss = 0.6034
Checking accuracy on validation set
Got 738 / 1000 correct (73.80)

Iteration 600, loss = 0.6595
Checking accuracy on validation set
Got 749 / 1000 correct (74.90)

Iteration 700, loss = 0.5449

Checking accuracy on validation set
Got 744 / 1000 correct (74.40)

current epochL: 7
Iteration 0, loss = 0.4442
Checking accuracy on validation set
Got 732 / 1000 correct (73.20)

Iteration 100, loss = 0.3310
Checking accuracy on validation set
Got 742 / 1000 correct (74.20)

Iteration 200, loss = 0.3595
Checking accuracy on validation set
Got 729 / 1000 correct (72.90)

Iteration 300, loss = 0.5537
Checking accuracy on validation set
Got 730 / 1000 correct (73.00)

Iteration 400, loss = 0.3462
Checking accuracy on validation set
Got 735 / 1000 correct (73.50)

Iteration 500, loss = 0.3496
Checking accuracy on validation set
Got 763 / 1000 correct (76.30)

Iteration 600, loss = 0.5785
Checking accuracy on validation set
Got 747 / 1000 correct (74.70)

Iteration 700, loss = 0.2291
Checking accuracy on validation set
Got 747 / 1000 correct (74.70)

current epochL: 8
Iteration 0, loss = 0.2554
Checking accuracy on validation set
Got 751 / 1000 correct (75.10)

Iteration 100, loss = 0.3989
Checking accuracy on validation set
Got 763 / 1000 correct (76.30)

Iteration 200, loss = 0.4799
Checking accuracy on validation set
Got 757 / 1000 correct (75.70)

Iteration 300, loss = 0.3244
Checking accuracy on validation set
Got 751 / 1000 correct (75.10)

Iteration 400, loss = 0.2785
Checking accuracy on validation set
Got 764 / 1000 correct (76.40)

Iteration 500, loss = 0.6507
Checking accuracy on validation set
Got 762 / 1000 correct (76.20)

Iteration 600, loss = 0.3024
Checking accuracy on validation set
Got 745 / 1000 correct (74.50)

Iteration 700, loss = 0.5283
Checking accuracy on validation set
Got 732 / 1000 correct (73.20)

current epochL: 9
Iteration 0, loss = 0.2400
Checking accuracy on validation set
Got 746 / 1000 correct (74.60)

Iteration 100, loss = 0.2353
Checking accuracy on validation set
Got 740 / 1000 correct (74.00)

Iteration 200, loss = 0.1778
Checking accuracy on validation set
Got 753 / 1000 correct (75.30)

Iteration 300, loss = 0.1559
Checking accuracy on validation set
Got 756 / 1000 correct (75.60)

Iteration 400, loss = 0.3302
Checking accuracy on validation set
Got 745 / 1000 correct (74.50)

Iteration 500, loss = 0.3191
Checking accuracy on validation set
Got 747 / 1000 correct (74.70)

Iteration 600, loss = 0.1905
Checking accuracy on validation set
Got 751 / 1000 correct (75.10)

Iteration 700, loss = 0.2807
Checking accuracy on validation set
Got 751 / 1000 correct (75.10)

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.

I implemented the AlexNet architecture. I did this because I wanted to see how easy it was to implement a state of the art CNN model using only the same concepts from this assignment (Max Pooling, Dropout, etc).

The AlexNet Paper mentions an input size of 256 with a kernel size of 11 for the first convolution. Since the input size for our dataset is 32, I changed the kernel size and the padding to account for this.

The architecture is exactly the same with 5 convolution layers and 3 FC layers

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 [112]: best_model = model  
          check_accuracy_part34(loader_test, best_model)
```

```
Checking accuracy on test set  
Got 7461 / 10000 correct (74.61)
```