ECE C147, Winter 2020
Homework 5
Submitted by Guy Ohayon

All of the relevant code is included at the end of this document.

# 1   Implement convolutional neural network layers

# Convolutional neural network layers

In this notebook, we will build the convolutional neural network layers. This will be followed by a spatial batchnorm, and then in the final notebook of this assignment, we will train a CNN to further improve the validation accuracy on CIFAR-10.

CS231n has built a solid API for building these modular frameworks and training them, and we will use their very well implemented framework as opposed to "reinventing the wheel." This includes using their Solver, various utility functions, their layer structure, and their implementation of fast CNN layers. This also includes nndl.fc_net, nndl.layers, and nndl.layer_utils. As in prior assignments, we thank Serena Yeung & Justin Johnson for permission to use code written for the CS 231n class (cs231n.stanford.edu).

```
In [1]:  ## Import and setups

         import time
         import numpy as np
         import matplotlib.pyplot as plt
         from nndl.conv_layers import *
         from cs231n.data_utils import get_CIFAR10_data
         from cs231n.gradient_check import eval_numerical_gradient, eval_numeri
         cal_gradient_array
         from cs231n.solver import Solver

         %matplotlib inline
         plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plo
         ts
         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
         ))))
```

## Implementing CNN layers

Just as we implemented modular layers for fully connected networks, batch normalization, and dropout, we'll want to implement modular layers for convolutional neural networks. These layers are in
nndl/conv_layers.py.

## Convolutional forward pass

Begin by implementing a naive version of the forward pass of the CNN that uses `for` loops. This function is `conv_forward_naive` in `nndl/conv_layers.py`. Don't worry about efficiency of implementation. Later on, we provide a fast implementation of these layers. This version ought to test your understanding of convolution. In our implementation, there is a triple `for` loop.

After you implement `conv_forward_naive`, test your implementation by running the cell below.

```
In [2]:  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 1e-8
         print('Testing conv_forward_naive')
         print('difference: ', rel_error(out, correct_out))
```

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

## Convolutional backward pass

Now, implement a naive version of the backward pass of the CNN. The function is `conv_backward_naive` in `nndl/conv_layers.py`. Don't worry about efficiency of implementation. Later on, we provide a fast implementation of these layers. This version ought to test your understanding of convolution. In our implementation, there is a quadruple `for` loop.

After you implement `conv_backward_naive`, test your implementation by running the cell below.

```
In [8]:  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}

         out, cache = conv_forward_naive(x,w,b,conv_param)

         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 1e-9'
         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:  1.7424735785934504e-09
dw error:  2.90493911988149e-10
db error:  1.3148447007113847e-11
```

## Max pool forward pass

In this section, we will implement the forward pass of the max pool. The function is `max_pool_forward_naive` in `nndl/conv_layers.py`. Do not worry about the efficiency of implementation.

After you implement `max_pool_forward_naive`, test your implementation by running the cell below.

```
In [13]: 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 around 1e-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 pool backward pass

In this section, you will implement the backward pass of the max pool. The function is
`max_pool_backward_naive` in `nndl/conv_layers.py`. Do not worry about the efficiency of
implementation.

After you implement `max_pool_backward_naive`, test your implementation by running the cell below.

```
In [14]: 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_naiv
         e(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 around 1e-12
         print('Testing max_pool_backward_naive function:')
         print('dx error: ', rel_error(dx, dx_num))
```

```
Testing max_pool_backward_naive function:
dx error:  3.2756214023732326e-12
```

# Fast implementation of the CNN layers

Implementing fast versions of the CNN layers can be difficult. We will provide you with the fast layers implemented by cs231n. They are provided in `cs231n/fast_layers.py`.

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

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

**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 cell below.

You should see pretty drastic speedups in the implementation of these layers. On our machine, the forward pass speeds up by 17x and the backward pass speeds up by 840x. Of course, these numbers will vary from machine to machine, as well as on your precise implementation of the naive layers.

In [15]:
```python
from cs231n.fast_layers import conv_forward_fast, conv_backward_fast
from time import time

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: 4.174019s
Fast: 0.025160s
Speedup: 165.900093x
Difference:  5.7276727339241154e-11

Testing conv_backward_fast:
Naive: 6.863882s
Fast: 0.010698s
Speedup: 641.628028x
dx difference:  1.4309080970108594e-11
dw difference:  7.877815815390147e-11
db difference:  0.0
```

In [16]:
```python
from cs231n.fast_layers import max_pool_forward_fast, max_pool_backward_fast

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('speedup: %fx' % ((t1 - t0) / (t2 - t1)))
print('dx difference: ', rel_error(dx_naive, dx_fast))
```

```
Testing pool_forward_fast:
Naive: 0.299417s
fast: 0.002297s
speedup: 130.355616x
difference:  0.0

Testing pool_backward_fast:
Naive: 0.386648s
speedup: 45.100339x
dx difference:  0.0
```

# Implementation of cascaded layers

We've provided the following functions in `nndl/conv_layer_utils.py`:

```
- conv_relu_forward
- conv_relu_backward
- conv_relu_pool_forward
- conv_relu_pool_backward
```

These use the fast implementations of the conv net layers. You can test them below:

In [19]:
```python
from nndl.conv_layer_utils import conv_relu_pool_forward, conv_relu_pool_backward

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)

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:  2.0024267129946835e-08
dw error:  1.268998646853029e-09
db error:  4.518439410456892e-11
```

In [20]:
```python
from nndl.conv_layer_utils import conv_relu_forward, conv_relu_backward

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)

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:   3.2510365523829036e-09
dw error:   4.535278548200337e-09
db error:   8.921161938841727e-11
```

# What next?

We saw how helpful batch normalization was for training FC nets. In the next notebook, we'll implement a batch normalization for convolutional neural networks, and then finish off by implementing a CNN to improve our validation accuracy on CIFAR-10.

# 2 Implement spatial normalization for CNNs

# Spatial batch normalization

In fully connected networks, we performed batch normalization on the activations. To do something equivalent on CNNs, we modify batch normalization slightly.

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 accepts inputs of shape `(N, C, H, W)` and produces 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.

How do we calculate the spatial averages? First, notice that for the `C` feature maps we have (i.e., the layer has `C` filters) that each of these ought to have its own batch norm statistics, since each feature map may be picking out very different features in the images. However, within a feature map, we may assume that across all inputs and across all locations in the feature map, there ought to be relatively similar first and second order statistics. Hence, one way to think of spatial batch-normalization is to reshape the `(N, C, H, W)` array as an `(N*H*W, C)` array and perform batch normalization on this array.

Since spatial batch norm and batch normalization are similar, it'd be good to at this point also copy and paste our prior implemented layers from HW #4. Please copy and paste your prior implemented code from HW #4 to start this assignment. If you did not correctly implement the layers in HW #4, you may collaborate with a classmate to use their implementations from HW #4. You may also visit TA or Prof OH to correct your implementation.

You'll want to copy and paste from HW #4:

```
- layers.py for your FC network layers, as well as batchnorm and dropout.
- layer_utils.py for your combined FC network layers.
- optim.py for your optimizers.
```

Be sure to place these in the `nndl/` directory so they're imported correctly. Note, as announced in class, we will not be releasing our solutions.

If you use your prior implementations of the batchnorm, then your spatial batchnorm implementation may be very short. Our implementations of the forward and backward pass are each 6 lines of code.

CS231n has built a solid API for building these modular frameworks and training them, and we will use their very well implemented framework as opposed to "reinventing the wheel." This includes using their Solver, various utility functions, their layer structure, and their implementation of fast CNN layers. This also includes nndl.fc_net, nndl.layers, and nndl.layer_utils. As in prior assignments, we thank Serena Yeung & Justin Johnson for permission to use code written for the CS 231n class (cs231n.stanford.edu).

In [5]:
```python
## Import and setups

import time
import numpy as np
import matplotlib.pyplot as plt
from nndl.conv_layers import *
from cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
from cs231n.solver import Solver

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

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

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

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

## Spatial batch normalization forward pass

Implement the forward pass, spatial_batchnorm_forward in nndl/conv_layers.py. Test your implementation by running the cell below.

In [34]:
```python
# 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:  [ 9.45617032  8.28015132 10.72071444]
  Stds:  [4.63858299 3.64970889 4.94905998]
After spatial batch normalization:
  Shape:  (2, 3, 4, 5)
  Means:  [ 1.44328993e-16  2.22044605e-17 -1.77635684e-16]
  Stds:  [0.99999977 0.99999962 0.9999998 ]
After spatial batch normalization (nontrivial gamma, beta):
  Shape:  (2, 3, 4, 5)
  Means:  [6. 7. 8.]
  Stds:  [2.9999993  3.9999985  4.99999898]
```

## Spatial batch normalization backward pass

Implement the backward pass, `spatial_batchnorm_backward` in `nndl/conv_layers.py`. Test your implementation by running the cell below.

In [35]:
```python
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)

_, cache = spatial_batchnorm_forward(x, gamma, beta, bn_param)
dx, dgamma, dbeta = spatial_batchnorm_backward(dout, cache)
print('dx error: ', rel_error(dx_num, dx))
print('dgamma error: ', rel_error(da_num, dgamma))
print('dbeta error: ', rel_error(db_num, dbeta))
```

```
dx error:  2.4036691259599493e-07
dgamma error:  7.065051488624474e-10
dbeta error:  2.4721900718328043e-11
```

# 3    Optimize your CNN for CIFAR-10

# Convolutional neural networks

In this notebook, we'll put together our convolutional layers to implement a 3-layer CNN. Then, we'll ask you to implement a CNN that can achieve > 65% validation error on CIFAR-10.

CS231n has built a solid API for building these modular frameworks and training them, and we will use their very well implemented framework as opposed to "reinventing the wheel." This includes using their Solver, various utility functions, their layer structure, and their implementation of fast CNN layers. This also includes using nndl.fc_net, nndl.layers, and nndl.layer_utils. As in prior assignments, we thank Serena Yeung & Justin Johnson for permission to use code written for the CS 231n class (cs231n.stanford.edu).

If you have not completed the Spatial BatchNorm Notebook, please see the following description from that notebook:

Please copy and paste your prior implemented code from HW #4 to start this assignment. If you did not correctly implement the layers in HW #4, you may collaborate with a classmate to use their layer implementations from HW #4. You may also visit TA or Prof OH to correct your implementation.

You'll want to copy and paste from HW #4:

```
- layers.py for your FC network layers, as well as batchnorm and dropout.
- layer_utils.py for your combined FC network layers.
- optim.py for your optimizers.
```

Be sure to place these in the `nndl/` directory so they're imported correctly. Note, as announced in class, we will not be releasing our solutions.

2/20/2020

In [1]:
```python
# As usual, a bit of setup

import numpy as np
import matplotlib.pyplot as plt
from nndl.cnn import *
from cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient_array, eval_
numerical_gradient
from nndl.layers import *
from nndl.conv_layers import *
from cs231n.fast_layers import *
from cs231n.solver import Solver

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

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

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

In [2]:
```python
# Load the (preprocessed) CIFAR10 data.

data = get_CIFAR10_data()
for k in data.keys():
  print('{}: {} '.format(k, data[k].shape))
```

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

# Three layer CNN

In this notebook, you will implement a three layer CNN. The `ThreeLayerConvNet` class is in `nndl/cnn.py`. You'll need to modify that code for this section, including the initialization, as well as the calculation of the loss and gradients. You should be able to use the building blocks you have either earlier coded or that we have provided. Be sure to use the fast layers.

The architecture of this CNN will be:

conv - relu - 2x2 max pool - affine - relu - affine - softmax

We won't use batchnorm yet. You've also done enough of these to know how to debug; use the cells below.

Note: As we are implementing several layers CNN networks. The gradient error can be expected for the `eval_numerical_gradient()` function. If your `W1 max relative error` and `W2 max relative error` are around or below 0.01, they should be acceptable. Other errors should be less than 1e-5.

```
In [10]:  num_inputs = 2
          input_dim = (3, 16, 16)
          reg = 0.0
          num_classes = 10
          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)
          for param_name in sorted(grads):
              f = lambda _: model.loss(X, y)[0]
              param_grad_num = eval_numerical_gradient(f, model.params[param_nam
          e], verbose=False, h=1e-6)
              e = rel_error(param_grad_num, grads[param_name])
              print('{} max relative error: {}'.format(param_name, rel_error(par
          am_grad_num, grads[param_name])))
```

```
W1 max relative error: 0.003339638093857734
W2 max relative error: 0.007733551671560727
W3 max relative error: 0.00013319376741324182
b1 max relative error: 2.912627580018027e-06
b2 max relative error: 6.225815960689804e-07
b3 max relative error: 8.210797326525663e-10
```

## Overfit small dataset

To check your CNN implementation, let's overfit a small dataset.

In [11]:
```python
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=10, batch_size=50,
                update_rule='adam',
                optim_config={
                  'learning_rate': 1e-3,
                },
                verbose=True, print_every=1)
solver.train()
```

```
(Iteration 1 / 20) loss: 2.338184
(Epoch 0 / 10) train acc: 0.180000; val_acc: 0.088000
(Iteration 2 / 20) loss: 3.145421
(Epoch 1 / 10) train acc: 0.220000; val_acc: 0.114000
(Iteration 3 / 20) loss: 4.220575
(Iteration 4 / 20) loss: 2.789140
(Epoch 2 / 10) train acc: 0.290000; val_acc: 0.122000
(Iteration 5 / 20) loss: 2.198183
(Iteration 6 / 20) loss: 2.334225
(Epoch 3 / 10) train acc: 0.290000; val_acc: 0.110000
(Iteration 7 / 20) loss: 2.026942
(Iteration 8 / 20) loss: 1.948826
(Epoch 4 / 10) train acc: 0.340000; val_acc: 0.165000
(Iteration 9 / 20) loss: 1.859220
(Iteration 10 / 20) loss: 2.146719
(Epoch 5 / 10) train acc: 0.550000; val_acc: 0.164000
(Iteration 11 / 20) loss: 1.968799
(Iteration 12 / 20) loss: 1.587138
(Epoch 6 / 10) train acc: 0.580000; val_acc: 0.140000
(Iteration 13 / 20) loss: 1.415540
(Iteration 14 / 20) loss: 1.577199
(Epoch 7 / 10) train acc: 0.640000; val_acc: 0.186000
(Iteration 15 / 20) loss: 1.260498
(Iteration 16 / 20) loss: 1.126015
(Epoch 8 / 10) train acc: 0.560000; val_acc: 0.190000
(Iteration 17 / 20) loss: 1.230003
(Iteration 18 / 20) loss: 1.154813
(Epoch 9 / 10) train acc: 0.670000; val_acc: 0.211000
(Iteration 19 / 20) loss: 0.833167
(Iteration 20 / 20) loss: 0.915641
(Epoch 10 / 10) train acc: 0.740000; val_acc: 0.193000
```

In [12]:
```python
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 network

Now we train the 3 layer CNN on CIFAR-10 and assess its accuracy.

In [13]:
```python
model = ThreeLayerConvNet(weight_scale=0.001, hidden_dim=500, reg=0.00
1)

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()
```
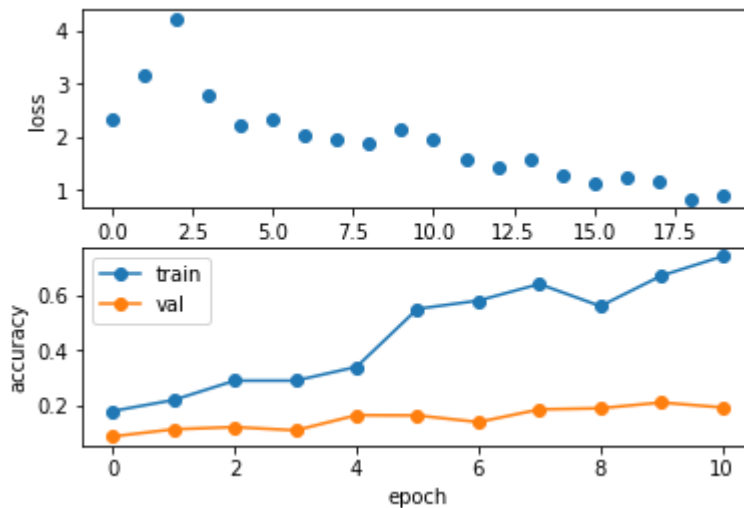
```
(Iteration 1 / 980) loss: 2.302710
(Epoch 0 / 1) train acc: 0.099000; val_acc: 0.107000
(Iteration 21 / 980) loss: 1.867889
(Iteration 41 / 980) loss: 1.823136
(Iteration 61 / 980) loss: 1.759033
(Iteration 81 / 980) loss: 1.863504
(Iteration 101 / 980) loss: 1.882208
(Iteration 121 / 980) loss: 1.974029
(Iteration 141 / 980) loss: 1.997073
(Iteration 161 / 980) loss: 1.563230
(Iteration 181 / 980) loss: 1.585403
(Iteration 201 / 980) loss: 1.759488
(Iteration 221 / 980) loss: 1.753828
(Iteration 241 / 980) loss: 2.131299
(Iteration 261 / 980) loss: 1.408404
(Iteration 281 / 980) loss: 1.405383
(Iteration 301 / 980) loss: 1.720499
(Iteration 321 / 980) loss: 1.567853
(Iteration 341 / 980) loss: 1.481917
(Iteration 361 / 980) loss: 1.491857
(Iteration 381 / 980) loss: 1.694706
(Iteration 401 / 980) loss: 1.911365
(Iteration 421 / 980) loss: 1.884573
(Iteration 441 / 980) loss: 1.511824
(Iteration 461 / 980) loss: 1.810901
(Iteration 481 / 980) loss: 1.654326
(Iteration 501 / 980) loss: 1.414904
(Iteration 521 / 980) loss: 1.754829
(Iteration 541 / 980) loss: 1.792938
(Iteration 561 / 980) loss: 1.961524
(Iteration 581 / 980) loss: 1.524744
(Iteration 601 / 980) loss: 1.466191
(Iteration 621 / 980) loss: 1.757177
(Iteration 641 / 980) loss: 1.750617
(Iteration 661 / 980) loss: 1.822015
(Iteration 681 / 980) loss: 1.608533
(Iteration 701 / 980) loss: 1.471488
(Iteration 721 / 980) loss: 1.617726
(Iteration 741 / 980) loss: 1.343086
(Iteration 761 / 980) loss: 1.537552
(Iteration 781 / 980) loss: 1.539328
(Iteration 801 / 980) loss: 1.431361
(Iteration 821 / 980) loss: 1.492462
(Iteration 841 / 980) loss: 1.674484
(Iteration 861 / 980) loss: 1.634723
(Iteration 881 / 980) loss: 1.608059
(Iteration 901 / 980) loss: 1.704482
(Iteration 921 / 980) loss: 1.465134
(Iteration 941 / 980) loss: 1.906660
(Iteration 961 / 980) loss: 1.723571
(Epoch 1 / 1) train acc: 0.448000; val_acc: 0.456000
```

# Get > 65% validation accuracy on CIFAR-10.

In the last part of the assignment, we'll now ask you to train a CNN to get better than 65% validation accuracy on CIFAR-10.

## Things you should try:

- Filter size: Above we used 7x7; but VGGNet and onwards showed stacks of 3x3 filters are good.
- Number of filters: Above we used 32 filters. Do more or fewer do better?
- Batch normalization: Try adding spatial batch normalization after convolution layers and vanilla batch normalization after affine layers. Do your networks train faster?
- Network architecture: Can a deeper CNN do better? Consider these architectures:
  - [conv-relu-pool]xN - conv - relu - [affine]xM - [softmax or SVM]
  - [conv-relu-pool]XN - [affine]XM - [softmax or SVM]
  - [conv-relu-conv-relu-pool]xN - [affine]xM - [softmax or SVM]

## Tips for training

For each network architecture that you try, you should tune the learning rate and regularization strength. When doing this there are a couple of 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.

```
In [16]:  # ==================================================================== #
          # YOUR CODE HERE:
          #   Implement a CNN to achieve greater than 65% validation accuracy
          #   on CIFAR-10.
          # ==================================================================== #

          model = ThreeLayerConvNet(weight_scale=0.001, hidden_dim=800, reg=0.01
          , num_filters=80, filter_size=3)

          solver = Solver(model, data,
                          num_epochs=15, batch_size=128,
                          update_rule='adam',
                          optim_config={
                              'learning_rate': 6e-4,
                          },
                          verbose=True, print_every=20)
          solver.train()

          # ==================================================================== #
          # END YOUR CODE HERE
          # ==================================================================== #
```

```
(Iteration 1 / 5730) loss: 2.302662
(Epoch 0 / 15) train acc: 0.100000; val_acc: 0.113000
(Iteration 21 / 5730) loss: 1.822297
(Iteration 41 / 5730) loss: 1.739436
(Iteration 61 / 5730) loss: 1.697439
(Iteration 81 / 5730) loss: 1.636978
(Iteration 101 / 5730) loss: 1.600068
(Iteration 121 / 5730) loss: 1.464973
(Iteration 141 / 5730) loss: 1.323235
(Iteration 161 / 5730) loss: 1.508464
(Iteration 181 / 5730) loss: 1.562103
(Iteration 201 / 5730) loss: 1.697535
(Iteration 221 / 5730) loss: 1.382375
(Iteration 241 / 5730) loss: 1.271008
(Iteration 261 / 5730) loss: 1.375135
(Iteration 281 / 5730) loss: 1.275878
(Iteration 301 / 5730) loss: 1.513031
(Iteration 321 / 5730) loss: 1.395120
(Iteration 341 / 5730) loss: 1.417776
(Iteration 361 / 5730) loss: 1.199762
(Iteration 381 / 5730) loss: 1.279455
(Epoch 1 / 15) train acc: 0.558000; val_acc: 0.565000
(Iteration 401 / 5730) loss: 1.220100
(Iteration 421 / 5730) loss: 1.225320
(Iteration 441 / 5730) loss: 1.278838
(Iteration 461 / 5730) loss: 1.276752
(Iteration 481 / 5730) loss: 1.307980
(Iteration 501 / 5730) loss: 1.218306
(Iteration 521 / 5730) loss: 1.165487
(Iteration 541 / 5730) loss: 1.083059
(Iteration 561 / 5730) loss: 1.273285
(Iteration 581 / 5730) loss: 1.119473
(Iteration 601 / 5730) loss: 1.335253
(Iteration 621 / 5730) loss: 1.070652
(Iteration 641 / 5730) loss: 1.318953
(Iteration 661 / 5730) loss: 1.097223
(Iteration 681 / 5730) loss: 1.294324
(Iteration 701 / 5730) loss: 1.118972
(Iteration 721 / 5730) loss: 1.143132
(Iteration 741 / 5730) loss: 1.135419
(Iteration 761 / 5730) loss: 0.984195
(Epoch 2 / 15) train acc: 0.615000; val_acc: 0.577000
(Iteration 781 / 5730) loss: 1.109513
(Iteration 801 / 5730) loss: 1.199963
(Iteration 821 / 5730) loss: 1.015277
(Iteration 841 / 5730) loss: 1.172329
(Iteration 861 / 5730) loss: 1.272491
(Iteration 881 / 5730) loss: 1.035057
(Iteration 901 / 5730) loss: 1.089504
(Iteration 921 / 5730) loss: 1.063471
(Iteration 941 / 5730) loss: 1.006050
(Iteration 961 / 5730) loss: 1.131010
(Iteration 981 / 5730) loss: 1.055226
(Iteration 1001 / 5730) loss: 1.220967
(Iteration 1021 / 5730) loss: 1.035626
(Iteration 1041 / 5730) loss: 1.187960
(Iteration 1061 / 5730) loss: 0.990191
```

```
(Iteration 1081 / 5730) loss: 0.880529
(Iteration 1101 / 5730) loss: 1.076975
(Iteration 1121 / 5730) loss: 1.132610
(Iteration 1141 / 5730) loss: 0.800732
(Epoch 3 / 15) train acc: 0.633000; val_acc: 0.609000
(Iteration 1161 / 5730) loss: 1.090642
(Iteration 1181 / 5730) loss: 1.331991
(Iteration 1201 / 5730) loss: 0.982377
(Iteration 1221 / 5730) loss: 1.047051
(Iteration 1241 / 5730) loss: 1.067707
(Iteration 1261 / 5730) loss: 0.959963
(Iteration 1281 / 5730) loss: 1.047880
(Iteration 1301 / 5730) loss: 0.945537
(Iteration 1321 / 5730) loss: 1.092515
(Iteration 1341 / 5730) loss: 0.921123
(Iteration 1361 / 5730) loss: 0.885517
(Iteration 1381 / 5730) loss: 0.960253
(Iteration 1401 / 5730) loss: 1.023843
(Iteration 1421 / 5730) loss: 1.011569
(Iteration 1441 / 5730) loss: 1.013666
(Iteration 1461 / 5730) loss: 1.334589
(Iteration 1481 / 5730) loss: 1.022286
(Iteration 1501 / 5730) loss: 0.872740
(Iteration 1521 / 5730) loss: 0.937490
(Epoch 4 / 15) train acc: 0.611000; val_acc: 0.595000
(Iteration 1541 / 5730) loss: 1.028597
(Iteration 1561 / 5730) loss: 0.876702
(Iteration 1581 / 5730) loss: 1.132168
(Iteration 1601 / 5730) loss: 0.929874
(Iteration 1621 / 5730) loss: 1.100173
(Iteration 1641 / 5730) loss: 1.092112
(Iteration 1661 / 5730) loss: 1.204253
(Iteration 1681 / 5730) loss: 0.998772
(Iteration 1701 / 5730) loss: 0.812532
(Iteration 1721 / 5730) loss: 1.070772
(Iteration 1741 / 5730) loss: 0.889282
(Iteration 1761 / 5730) loss: 0.962697
(Iteration 1781 / 5730) loss: 0.916171
(Iteration 1801 / 5730) loss: 0.896096
(Iteration 1821 / 5730) loss: 1.000129
(Iteration 1841 / 5730) loss: 0.843282
(Iteration 1861 / 5730) loss: 0.933782
(Iteration 1881 / 5730) loss: 0.877623
(Iteration 1901 / 5730) loss: 0.770611
(Epoch 5 / 15) train acc: 0.687000; val_acc: 0.632000
(Iteration 1921 / 5730) loss: 1.034500
(Iteration 1941 / 5730) loss: 1.072820
(Iteration 1961 / 5730) loss: 0.943722
(Iteration 1981 / 5730) loss: 0.903900
(Iteration 2001 / 5730) loss: 0.978961
(Iteration 2021 / 5730) loss: 0.944058
(Iteration 2041 / 5730) loss: 0.982706
(Iteration 2061 / 5730) loss: 0.945340
(Iteration 2081 / 5730) loss: 0.784325
(Iteration 2101 / 5730) loss: 0.766858
(Iteration 2121 / 5730) loss: 0.926472
(Iteration 2141 / 5730) loss: 0.905241
```

```
(Iteration 2161 / 5730) loss: 0.761714
(Iteration 2181 / 5730) loss: 0.817416
(Iteration 2201 / 5730) loss: 0.950651
(Iteration 2221 / 5730) loss: 0.817195
(Iteration 2241 / 5730) loss: 1.030922
(Iteration 2261 / 5730) loss: 1.037213
(Iteration 2281 / 5730) loss: 0.966074
(Epoch 6 / 15) train acc: 0.679000; val_acc: 0.628000
(Iteration 2301 / 5730) loss: 0.926695
(Iteration 2321 / 5730) loss: 1.200794
(Iteration 2341 / 5730) loss: 0.728167
(Iteration 2361 / 5730) loss: 0.849453
(Iteration 2381 / 5730) loss: 0.850361
(Iteration 2401 / 5730) loss: 0.951400
(Iteration 2421 / 5730) loss: 0.863175
(Iteration 2441 / 5730) loss: 0.938395
(Iteration 2461 / 5730) loss: 0.745604
(Iteration 2481 / 5730) loss: 0.849992
(Iteration 2501 / 5730) loss: 1.047825
(Iteration 2521 / 5730) loss: 0.984919
(Iteration 2541 / 5730) loss: 0.773999
(Iteration 2561 / 5730) loss: 0.876067
(Iteration 2581 / 5730) loss: 0.880775
(Iteration 2601 / 5730) loss: 0.753403
(Iteration 2621 / 5730) loss: 1.063312
(Iteration 2641 / 5730) loss: 0.764027
(Iteration 2661 / 5730) loss: 0.660783
(Epoch 7 / 15) train acc: 0.704000; val_acc: 0.657000
(Iteration 2681 / 5730) loss: 0.875050
(Iteration 2701 / 5730) loss: 0.731897
(Iteration 2721 / 5730) loss: 0.832229
(Iteration 2741 / 5730) loss: 0.827563
(Iteration 2761 / 5730) loss: 1.018010
(Iteration 2781 / 5730) loss: 0.782663
(Iteration 2801 / 5730) loss: 0.953400
(Iteration 2821 / 5730) loss: 0.898229
(Iteration 2841 / 5730) loss: 0.855345
(Iteration 2861 / 5730) loss: 1.111483
(Iteration 2881 / 5730) loss: 0.789999
(Iteration 2901 / 5730) loss: 0.883426
(Iteration 2921 / 5730) loss: 0.902819
(Iteration 2941 / 5730) loss: 0.746762
(Iteration 2961 / 5730) loss: 0.995886
(Iteration 2981 / 5730) loss: 0.811427
(Iteration 3001 / 5730) loss: 0.834081
(Iteration 3021 / 5730) loss: 0.952963
(Iteration 3041 / 5730) loss: 0.944784
(Epoch 8 / 15) train acc: 0.711000; val_acc: 0.634000
(Iteration 3061 / 5730) loss: 0.837555
(Iteration 3081 / 5730) loss: 0.694587
(Iteration 3101 / 5730) loss: 0.821755
(Iteration 3121 / 5730) loss: 0.815190
(Iteration 3141 / 5730) loss: 0.768178
(Iteration 3161 / 5730) loss: 0.750169
(Iteration 3181 / 5730) loss: 0.842079
(Iteration 3201 / 5730) loss: 0.713645
(Iteration 3221 / 5730) loss: 0.853219
```

```
(Iteration 3241 / 5730) loss: 0.749841
(Iteration 3261 / 5730) loss: 0.806953
(Iteration 3281 / 5730) loss: 0.747577
(Iteration 3301 / 5730) loss: 0.759039
(Iteration 3321 / 5730) loss: 0.761904
(Iteration 3341 / 5730) loss: 0.769731
(Iteration 3361 / 5730) loss: 1.103440
(Iteration 3381 / 5730) loss: 0.806305
(Iteration 3401 / 5730) loss: 0.722995
(Iteration 3421 / 5730) loss: 0.910105
(Epoch 9 / 15) train acc: 0.741000; val_acc: 0.640000
(Iteration 3441 / 5730) loss: 0.789341
(Iteration 3461 / 5730) loss: 0.794981
(Iteration 3481 / 5730) loss: 0.808493
(Iteration 3501 / 5730) loss: 0.728320
(Iteration 3521 / 5730) loss: 0.813904
(Iteration 3541 / 5730) loss: 0.811307
(Iteration 3561 / 5730) loss: 0.813526
(Iteration 3581 / 5730) loss: 0.965592
(Iteration 3601 / 5730) loss: 1.043618
(Iteration 3621 / 5730) loss: 0.826015
(Iteration 3641 / 5730) loss: 0.874432
(Iteration 3661 / 5730) loss: 0.640500
(Iteration 3681 / 5730) loss: 0.719851
(Iteration 3701 / 5730) loss: 1.014668
(Iteration 3721 / 5730) loss: 0.896639
(Iteration 3741 / 5730) loss: 0.831060
(Iteration 3761 / 5730) loss: 0.738823
(Iteration 3781 / 5730) loss: 0.921722
(Iteration 3801 / 5730) loss: 0.878072
(Epoch 10 / 15) train acc: 0.692000; val_acc: 0.628000
(Iteration 3821 / 5730) loss: 0.881993
(Iteration 3841 / 5730) loss: 1.087112
(Iteration 3861 / 5730) loss: 0.659757
(Iteration 3881 / 5730) loss: 0.844886
(Iteration 3901 / 5730) loss: 0.838179
(Iteration 3921 / 5730) loss: 0.731724
(Iteration 3941 / 5730) loss: 0.667375
(Iteration 3961 / 5730) loss: 0.922343
(Iteration 3981 / 5730) loss: 0.805053
(Iteration 4001 / 5730) loss: 0.883172
(Iteration 4021 / 5730) loss: 0.736805
(Iteration 4041 / 5730) loss: 0.796258
(Iteration 4061 / 5730) loss: 0.598023
(Iteration 4081 / 5730) loss: 0.855600
(Iteration 4101 / 5730) loss: 0.630917
(Iteration 4121 / 5730) loss: 0.715877
(Iteration 4141 / 5730) loss: 0.509960
(Iteration 4161 / 5730) loss: 0.759190
(Iteration 4181 / 5730) loss: 0.613843
(Iteration 4201 / 5730) loss: 0.645178
(Epoch 11 / 15) train acc: 0.707000; val_acc: 0.601000
(Iteration 4221 / 5730) loss: 0.963549
(Iteration 4241 / 5730) loss: 0.704584
(Iteration 4261 / 5730) loss: 0.701464
(Iteration 4281 / 5730) loss: 0.718884
(Iteration 4301 / 5730) loss: 0.730558
```

```
                    (Iteration 4321 / 5730) loss: 0.632395
                    (Iteration 4341 / 5730) loss: 0.612482
                    (Iteration 4361 / 5730) loss: 0.668974
                    (Iteration 4381 / 5730) loss: 0.867646
                    (Iteration 4401 / 5730) loss: 0.687785
                    (Iteration 4421 / 5730) loss: 0.728129
                    (Iteration 4441 / 5730) loss: 0.708198
                    (Iteration 4461 / 5730) loss: 0.762829
                    (Iteration 4481 / 5730) loss: 0.719877
                    (Iteration 4501 / 5730) loss: 0.830460
                    (Iteration 4521 / 5730) loss: 0.923217
                    (Iteration 4541 / 5730) loss: 0.740081
                    (Iteration 4561 / 5730) loss: 0.800896
                    (Iteration 4581 / 5730) loss: 0.660932
                    (Epoch 12 / 15) train acc: 0.729000; val_acc: 0.655000
                    (Iteration 4601 / 5730) loss: 0.720061
                    (Iteration 4621 / 5730) loss: 0.676229
                    (Iteration 4641 / 5730) loss: 0.880973
                    (Iteration 4661 / 5730) loss: 0.705091
                    (Iteration 4681 / 5730) loss: 0.699740
                    (Iteration 4701 / 5730) loss: 0.609353
                    (Iteration 4721 / 5730) loss: 0.857134
                    (Iteration 4741 / 5730) loss: 0.729417
                    (Iteration 4761 / 5730) loss: 0.763072
                    (Iteration 4781 / 5730) loss: 0.818449
                    (Iteration 4801 / 5730) loss: 0.773662
                    (Iteration 4821 / 5730) loss: 0.895202
                    (Iteration 4841 / 5730) loss: 0.925682
                    (Iteration 4861 / 5730) loss: 0.781881
                    (Iteration 4881 / 5730) loss: 0.700775
                    (Iteration 4901 / 5730) loss: 0.856194
                    (Iteration 4921 / 5730) loss: 0.669190
                    (Iteration 4941 / 5730) loss: 0.847880
                    (Iteration 4961 / 5730) loss: 0.639630
                    (Epoch 13 / 15) train acc: 0.750000; val_acc: 0.639000
                    (Iteration 4981 / 5730) loss: 0.523683
                    (Iteration 5001 / 5730) loss: 0.764845
                    (Iteration 5021 / 5730) loss: 0.708426
                    (Iteration 5041 / 5730) loss: 0.773580
                    (Iteration 5061 / 5730) loss: 0.661314
                    (Iteration 5081 / 5730) loss: 0.582456
                    (Iteration 5101 / 5730) loss: 0.738108
                    (Iteration 5121 / 5730) loss: 0.901258
                    (Iteration 5141 / 5730) loss: 0.705411
                    (Iteration 5161 / 5730) loss: 0.678164
                    (Iteration 5181 / 5730) loss: 0.735594
                    (Iteration 5201 / 5730) loss: 0.920688
                    (Iteration 5221 / 5730) loss: 0.943875
                    (Iteration 5241 / 5730) loss: 0.883143
                    (Iteration 5261 / 5730) loss: 0.784406
                    (Iteration 5281 / 5730) loss: 0.787113
                    (Iteration 5301 / 5730) loss: 0.989408
                    (Iteration 5321 / 5730) loss: 0.717533
                    (Iteration 5341 / 5730) loss: 0.803206
                    (Epoch 14 / 15) train acc: 0.757000; val_acc: 0.647000
                    (Iteration 5361 / 5730) loss: 0.804627
                    (Iteration 5381 / 5730) loss: 0.674884
```

```
(Iteration 5401 / 5730) loss: 0.760615
(Iteration 5421 / 5730) loss: 0.815537
(Iteration 5441 / 5730) loss: 0.734104
(Iteration 5461 / 5730) loss: 0.764333
(Iteration 5481 / 5730) loss: 0.823910
(Iteration 5501 / 5730) loss: 0.800798
(Iteration 5521 / 5730) loss: 0.613245
(Iteration 5541 / 5730) loss: 0.620098
(Iteration 5561 / 5730) loss: 0.894341
(Iteration 5581 / 5730) loss: 0.603590
(Iteration 5601 / 5730) loss: 0.830513
(Iteration 5621 / 5730) loss: 0.872919
(Iteration 5641 / 5730) loss: 0.755896
(Iteration 5661 / 5730) loss: 0.766207
(Iteration 5681 / 5730) loss: 0.720361
(Iteration 5701 / 5730) loss: 0.923137
(Iteration 5721 / 5730) loss: 0.759570
(Epoch 15 / 15) train acc: 0.771000; val_acc: 0.669000
```

# 4   Code for all sections

optim.py

```
1    import numpy as np
2
3    """
4    This code was originally written for CS 231n at Stanford University
5    (cs231n.stanford.edu).  It has been modified in various areas for use in the
6    ECE 239AS class at UCLA.  This includes the descriptions of what code to
7    implement as well as some slight potential changes in variable names to be
8    consistent with class nomenclature.  We thank Justin Johnson & Serena Yeung for
9    permission to use this code.  To see the original version, please visit
10   cs231n.stanford.edu.
11   """
12
13   """
14   This file implements various first-order update rules that are commonly used for
15   training neural networks. Each update rule accepts current weights and the
16   gradient of the loss with respect to those weights and produces the next set of
17   weights. Each update rule has the same interface:
18
19   def update(w, dw, config=None):
20
21   Inputs:
22     - w: A numpy array giving the current weights.
23     - dw: A numpy array of the same shape as w giving the gradient of the
24       loss with respect to w.
25     - config: A dictionary containing hyperparameter values such as learning rate,
26       momentum, etc. If the update rule requires caching values over many
27       iterations, then config will also hold these cached values.
28
29   Returns:
30     - next_w: The next point after the update.
31     - config: The config dictionary to be passed to the next iteration of the
32       update rule.
33
34   NOTE: For most update rules, the default learning rate will probably not perform
35   well; however the default values of the other hyperparameters should work well
36   for a variety of different problems.
37
38   For efficiency, update rules may perform in-place updates, mutating w and
39   setting next_w equal to w.
40   """
41
42
43   def sgd(w, dw, config=None):
44     """
45     Performs vanilla stochastic gradient descent.
46
47     config format:
48     - learning_rate: Scalar learning rate.
49     """
50     if config is None: config = {}
51     config.setdefault('learning_rate', 1e-2)
52
53     w -= config['learning_rate'] * dw
54     return w, config
55
56
57   def sgd_momentum(w, dw, config=None):
58     """
59     Performs stochastic gradient descent with momentum.
60
61     config format:
62     - learning_rate: Scalar learning rate.
63     - momentum: Scalar between 0 and 1 giving the momentum value.
64       Setting momentum = 0 reduces to sgd.
65     - velocity: A numpy array of the same shape as w and dw used to store a moving
66       average of the gradients.
67     """
68     if config is None: config = {}
69     config.setdefault('learning_rate', 1e-2)
70     config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
```

```
71      v = config.get('velocity', np.zeros_like(w))        # gets velocity, else sets it to zero.
72
73      # ================================================================ #
74      # YOUR CODE HERE:
75      #    Implement the momentum update formula.  Return the updated weights
76      #    as next_w, and the updated velocity as v.
77      # ================================================================ #
78      v = config['momentum'] * v - config['learning_rate'] * dw
79      next_w = w + v
80
81      # ================================================================ #
82      # END YOUR CODE HERE
83      # ================================================================ #
84
85      config['velocity'] = v
86
87      return next_w, config
88
89  def sgd_nesterov_momentum(w, dw, config=None):
90      """
91      Performs stochastic gradient descent with Nesterov momentum.
92
93      config format:
94      - learning_rate: Scalar learning rate.
95      - momentum: Scalar between 0 and 1 giving the momentum value.
96        Setting momentum = 0 reduces to sgd.
97      - velocity: A numpy array of the same shape as w and dw used to store a moving
98        average of the gradients.
99      """
100     if config is None: config = {}
101     config.setdefault('learning_rate', 1e-2)
102     config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
103     v = config.get('velocity', np.zeros_like(w))        # gets velocity, else sets it to zero.
104
105     # ================================================================ #
106     # YOUR CODE HERE:
107     #    Implement the momentum update formula.  Return the updated weights
108     #    as next_w, and the updated velocity as v.
109     # ================================================================ #
110     v_old = v
111     v = config['momentum'] * v - config['learning_rate'] * dw
112     next_w = w + v + config['momentum'] * (v - v_old)
113
114     # ================================================================ #
115     # END YOUR CODE HERE
116     # ================================================================ #
117
118     config['velocity'] = v
119
120     return next_w, config
121
122  def rmsprop(w, dw, config=None):
123      """
124      Uses the RMSProp update rule, which uses a moving average of squared gradient
125      values to set adaptive per-parameter learning rates.
126
127      config format:
128      - learning_rate: Scalar learning rate.
129      - decay_rate: Scalar between 0 and 1 giving the decay rate for the squared
130        gradient cache.
131      - epsilon: Small scalar used for smoothing to avoid dividing by zero.
132      - beta: Moving average of second moments of gradients.
133      """
134     if config is None: config = {}
135     config.setdefault('learning_rate', 1e-2)
136     config.setdefault('decay_rate', 0.99)
137     config.setdefault('epsilon', 1e-8)
138     config.setdefault('a', np.zeros_like(w))
139
140     next_w = None
141
142     # ================================================================ #
143     # YOUR CODE HERE:
```

```python
144    #    Implement RMSProp.  Store the next value of w as next_w.  You need
145    #    to also store in config['a'] the moving average of the second
146    #    moment gradients, so they can be used for future gradients. Concretely,
147    #    config['a'] corresponds to "a" in the lecture notes.
148    # ================================================================ #
149    config['a'] = config['decay_rate'] * config['a'] + (1 - config['decay_rate']) * np.multiply(dw, dw)
150    next_w = w - np.multiply(config['learning_rate'] / (np.sqrt(config['a']) + config['epsilon']), dw)
151
152    # ================================================================ #
153    # END YOUR CODE HERE
154    # ================================================================ #
155
156    return next_w, config
157
158
159 def adam(w, dw, config=None):
160    """
161    Uses the Adam update rule, which incorporates moving averages of both the
162    gradient and its square and a bias correction term.
163
164    config format:
165    - learning_rate: Scalar learning rate.
166    - beta1: Decay rate for moving average of first moment of gradient.
167    - beta2: Decay rate for moving average of second moment of gradient.
168    - epsilon: Small scalar used for smoothing to avoid dividing by zero.
169    - m: Moving average of gradient.
170    - v: Moving average of squared gradient.
171    - t: Iteration number.
172    """
173    if config is None: config = {}
174    config.setdefault('learning_rate', 1e-3)
175    config.setdefault('beta1', 0.9)
176    config.setdefault('beta2', 0.999)
177    config.setdefault('epsilon', 1e-8)
178    config.setdefault('v', np.zeros_like(w))
179    config.setdefault('a', np.zeros_like(w))
180    config.setdefault('t', 0)
181
182    next_w = None
183
184    # ================================================================ #
185    # YOUR CODE HERE:
186    #    Implement Adam.  Store the next value of w as next_w.  You need
187    #    to also store in config['a'] the moving average of the second
188    #    moment gradients, and in config['v'] the moving average of the
189    #    first moments.  Finally, store in config['t'] the increasing time.
190    # ================================================================ #
191
192    config['t'] += 1
193    config['v'] = config['beta1'] * config['v'] + (1 - config['beta1']) * dw
194    config['a'] = config['beta2'] * config['a'] + (1 - config['beta2']) * np.multiply(dw, dw)
195    v_corr = config['v'] / (1 - (config['beta1'] ** config['t']))
196    a_corr = config['a'] / (1 - (config['beta2'] ** config['t']))
197    next_w = w - np.multiply((config['learning_rate'] / (np.sqrt(a_corr) + config['epsilon'])), v_corr)
198    # ================================================================ #
199    # END YOUR CODE HERE
200    # ================================================================ #
201
202    return next_w, config
203
204
205
206
207
208
```

## layers.py

```
1   import numpy as np
2   import pdb
3
4   """
5   This code was originally written for CS 231n at Stanford University
6   (cs231n.stanford.edu).  It has been modified in various areas for use in the
7   ECE 239AS class at UCLA.  This includes the descriptions of what code to
8   implement as well as some slight potential changes in variable names to be
9   consistent with class nomenclature.  We thank Justin Johnson & Serena Yeung for
10  permission to use this code.  To see the original version, please visit
11  cs231n.stanford.edu.
12  """
13
14  def affine_forward(x, w, b):
15    """
16    Computes the forward pass for an affine (fully-connected) layer.
17
18    The input x has shape (N, d_1, ..., d_k) and contains a minibatch of N
19    examples, where each example x[i] has shape (d_1, ..., d_k). We will
20    reshape each input into a vector of dimension D = d_1 * ... * d_k, and
21    then transform it to an output vector of dimension M.
22
23    Inputs:
24    - x: A numpy array containing input data, of shape (N, d_1, ..., d_k)
25    - w: A numpy array of weights, of shape (D, M)
26    - b: A numpy array of biases, of shape (M,)
27
28    Returns a tuple of:
29    - out: output, of shape (N, M)
30    - cache: (x, w, b)
31    """
32
33    # ================================================================ #
34    # YOUR CODE HERE:
35    #   Calculate the output of the forward pass.  Notice the dimensions
36    #   of w are D x M, which is the transpose of what we did in earlier
37    #   assignments.
38    # ================================================================ #
39
40    out = x.reshape(x.shape[0], -1).dot(w) + b
41    # ================================================================ #
42    # END YOUR CODE HERE
43    # ================================================================ #
44
45    cache = (x, w, b)
46    return out, cache
47
48
49  def affine_backward(dout, cache):
50    """
51    Computes the backward pass for an affine layer.
52
53    Inputs:
54    - dout: Upstream derivative, of shape (N, M)
55    - cache: Tuple of:
56      - x: Input data, of shape (N, d_1, ... d_k)
57      - w: Weights, of shape (D, M)
58
59    Returns a tuple of:
60    - dx: Gradient with respect to x, of shape (N, d1, ..., d_k)
61    - dw: Gradient with respect to w, of shape (D, M)
62    - db: Gradient with respect to b, of shape (M,)
63    """
64    x, w, b = cache
65    dx, dw, db = None, None, None
66
67    # ================================================================ #
68    # YOUR CODE HERE:
69    #   Calculate the gradients for the backward pass.
70    # ================================================================ #
71
72    db = np.sum(dout, axis=0)
73    dx = np.array(dout).dot(w.T).reshape(x.shape)
74    dw = x.reshape(x.shape[0], -1).T.dot(dout)
```

```
 75      # ================================================================ #
 76      # END YOUR CODE HERE
 77      # ================================================================ #
 78
 79      return dx, dw, db
 80
 81  def relu_forward(x):
 82      """
 83      Computes the forward pass for a layer of rectified linear units (ReLUs).
 84
 85      Input:
 86      - x: Inputs, of any shape
 87
 88      Returns a tuple of:
 89      - out: Output, of the same shape as x
 90      - cache: x
 91      """
 92      # ================================================================ #
 93      # YOUR CODE HERE:
 94      #    Implement the ReLU forward pass.
 95      # ================================================================ #
 96      out = x * (x > 0)
 97      # ================================================================ #
 98      # END YOUR CODE HERE
 99      # ================================================================ #
100
101      cache = x
102      return out, cache
103
104
105  def relu_backward(dout, cache):
106      """
107      Computes the backward pass for a layer of rectified linear units (ReLUs).
108
109      Input:
110      - dout: Upstream derivatives, of any shape
111      - cache: Input x, of same shape as dout
112
113      Returns:
114      - dx: Gradient with respect to x
115      """
116      x = cache
117
118      # ================================================================ #
119      # YOUR CODE HERE:
120      #    Implement the ReLU backward pass
121      # ================================================================ #
122      dx = dout * (x > 0)
123
124      # ================================================================ #
125      # END YOUR CODE HERE
126      # ================================================================ #
127
128      return dx
129
130  def batchnorm_forward(x, gamma, beta, bn_param):
131      """
132      Forward pass for batch normalization.
133
134      During training the sample mean and (uncorrected) sample variance are
135      computed from minibatch statistics and used to normalize the incoming data.
136      During training we also keep an exponentially decaying running mean of the mean
137      and variance of each feature, and these averages are used to normalize data
138      at test-time.
139
140      At each timestep we update the running averages for mean and variance using
141      an exponential decay based on the momentum parameter:
142
143      running_mean = momentum * running_mean + (1 - momentum) * sample_mean
144      running_var = momentum * running_var + (1 - momentum) * sample_var
145
146      Note that the batch normalization paper suggests a different test-time
147      behavior: they compute sample mean and variance for each feature using a
148      large number of training images rather than using a running average. For
149      this implementation we have chosen to use running averages instead since
150      they do not require an additional estimation step; the torch7 implementation
151      of batch normalization also uses running averages.
```

```
152
153      Input:
154      - x: Data of shape (N, D)
155      - gamma: Scale parameter of shape (D,)
156      - beta: Shift paremeter of shape (D,)
157      - bn_param: Dictionary with the following keys:
158        - mode: 'train' or 'test'; required
159        - eps: Constant for numeric stability
160        - momentum: Constant for running mean / variance.
161        - running_mean: Array of shape (D,) giving running mean of features
162        - running_var Array of shape (D,) giving running variance of features
163
164      Returns a tuple of:
165      - out: of shape (N, D)
166      - cache: A tuple of values needed in the backward pass
167      """
168      mode = bn_param['mode']
169      eps = bn_param.get('eps', 1e-5)
170      momentum = bn_param.get('momentum', 0.9)
171
172      N, D = x.shape
173      running_mean = bn_param.get('running_mean', np.zeros(D, dtype=x.dtype))
174      running_var = bn_param.get('running_var', np.zeros(D, dtype=x.dtype))
175
176      out, cache = None, None
177      if mode == 'train':
178
179        # ================================================================ #
180        # YOUR CODE HERE:
181        #   A few steps here:
182        #     (1) Calculate the running mean and variance of the minibatch.
183        #     (2) Normalize the activations with the running mean and variance.
184        #     (3) Scale and shift the normalized activations.  Store this
185        #         as the variable 'out'
186        #     (4) Store any variables you may need for the backward pass in
187        #         the 'cache' variable.
188        # ================================================================ #
189        running_mean = np.mean(x, axis=0)
190        running_var = np.var(x, axis=0)
191        x_hat = (x - running_mean) / np.sqrt(eps + running_var)
192        out = gamma * x_hat + beta
193        cache = (x_hat, x, running_mean, running_var, eps, gamma)
194
195        # ================================================================ #
196        # END YOUR CODE HERE
197        # ================================================================ #
198
199      elif mode == 'test':
200
201        # ================================================================ #
202        # YOUR CODE HERE:
203        #   Calculate the testing time normalized activation.  Normalize using
204        #   the running mean and variance, and then scale and shift appropriately.
205        #   Store the output as 'out'.
206        # ================================================================ #
207        x_hat = (x - running_mean) / np.sqrt(eps + running_var)
208        out = gamma * x_hat + beta
209
210        # ================================================================ #
211        # END YOUR CODE HERE
212        # ================================================================ #
213
214      else:
215        raise ValueError('Invalid forward batchnorm mode "%s"' % mode)
216
217      # Store the updated running means back into bn_param
218      bn_param['running_mean'] = running_mean
219      bn_param['running_var'] = running_var
220
221      return out, cache
222
223  def batchnorm_backward(dout, cache):
224      """
225      Backward pass for batch normalization.
226
227      For this implementation, you should write out a computation graph for
228      batch normalization on paper and propagate gradients backward through
```

```
229        intermediate nodes.
230
231        Inputs:
232        - dout: Upstream derivatives, of shape (N, D)
233        - cache: Variable of intermediates from batchnorm_forward.
234
235        Returns a tuple of:
236        - dx: Gradient with respect to inputs x, of shape (N, D)
237        - dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
238        - dbeta: Gradient with respect to shift parameter beta, of shape (D,)
239        """
240        dx, dgamma, dbeta = None, None, None
241
242        # ================================================================ #
243        # YOUR CODE HERE:
244        #    Implement the batchnorm backward pass, calculating dx, dgamma, and dbeta.
245        # ================================================================ #
246        x_hat, x, mean, var, eps, gamma = cache
247        N, D = x.shape
248        dbeta = np.sum(dout, axis=0)
249        dgamma = np.sum(dout * x_hat, axis=0)
250        dl_dxhat = dout * gamma
251        dl_dvar = (-1 / 2) * np.sum((1 / ((var + eps) ** (3 / 2))) * (x - mean) * dl_dxhat, axis=0)
252        dl_dmean = (-1 / (np.sqrt(var + eps))) * np.sum(dl_dxhat, axis=0)
253        dx = np.array((1 / np.sqrt(var + eps)) * dl_dxhat + (2 * (x - mean) / N) * dl_dvar + (1 / N) * dl_dmean)
254
255        # ================================================================ #
256        # END YOUR CODE HERE
257        # ================================================================ #
258
259        return dx, dgamma, dbeta
260
261    def dropout_forward(x, dropout_param):
262        """
263        Performs the forward pass for (inverted) dropout.
264
265        Inputs:
266        - x: Input data, of any shape
267        - dropout_param: A dictionary with the following keys:
268          - p: Dropout parameter. We drop each neuron output with probability p.
269          - mode: 'test' or 'train'. If the mode is train, then perform dropout;
270            if the mode is test, then just return the input.
271          - seed: Seed for the random number generator. Passing seed makes this
272            function deterministic, which is needed for gradient checking but not in
273            real networks.
274
275        Outputs:
276        - out: Array of the same shape as x.
277        - cache: A tuple (dropout_param, mask). In training mode, mask is the dropout
278          mask that was used to multiply the input; in test mode, mask is None.
279        """
280        p, mode = dropout_param['p'], dropout_param['mode']
281        if 'seed' in dropout_param:
282            np.random.seed(dropout_param['seed'])
283
284        mask = None
285        out = None
286
287        if mode == 'train':
288            # ================================================================ #
289            # YOUR CODE HERE:
290            #    Implement the inverted dropout forward pass during training time.
291            #    Store the masked and scaled activations in out, and store the
292            #    dropout mask as the variable mask.
293            # ================================================================ #
294            mask = np.random.rand(*x.shape) < p
295            out = x * mask / p
296
297            # ================================================================ #
298            # END YOUR CODE HERE
299            # ================================================================ #
300
301        elif mode == 'test':
302
303            # ================================================================ #
304            # YOUR CODE HERE:
305            #    Implement the inverted dropout forward pass during test time.
```

```python
306        # ================================================================ #
307        out = x
308
309        # ================================================================ #
310        # END YOUR CODE HERE
311        # ================================================================ #
312
313    cache = (dropout_param, mask)
314    out = out.astype(x.dtype, copy=False)
315
316    return out, cache
317
318  def dropout_backward(dout, cache):
319    """
320    Perform the backward pass for (inverted) dropout.
321
322    Inputs:
323    - dout: Upstream derivatives, of any shape
324    - cache: (dropout_param, mask) from dropout_forward.
325    """
326    dropout_param, mask = cache
327    mode = dropout_param['mode']
328
329    dx = None
330    if mode == 'train':
331      # ================================================================ #
332      # YOUR CODE HERE:
333      #   Implement the inverted dropout backward pass during training time.
334      # ================================================================ #
335      dx = dout * mask / dropout_param['p']
336
337      # ================================================================ #
338      # END YOUR CODE HERE
339      # ================================================================ #
340    elif mode == 'test':
341      # ================================================================ #
342      # YOUR CODE HERE:
343      #   Implement the inverted dropout backward pass during test time.
344      # ================================================================ #
345      dx = dout
346      # ================================================================ #
347      # END YOUR CODE HERE
348      # ================================================================ #
349    return dx
350
351  def svm_loss(x, y):
352    """
353    Computes the loss and gradient using for multiclass SVM classification.
354
355    Inputs:
356    - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
357      for the ith input.
358    - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
359      0 <= y[i] < C
360
361    Returns a tuple of:
362    - loss: Scalar giving the loss
363    - dx: Gradient of the loss with respect to x
364    """
365    N = x.shape[0]
366    correct_class_scores = x[np.arange(N), y]
367    margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
368    margins[np.arange(N), y] = 0
369    loss = np.sum(margins) / N
370    num_pos = np.sum(margins > 0, axis=1)
371    dx = np.zeros_like(x)
372    dx[margins > 0] = 1
373    dx[np.arange(N), y] -= num_pos
374    dx /= N
375    return loss, dx
376
377
378  def softmax_loss(x, y):
379    """
380    Computes the loss and gradient for softmax classification.
381
382    Inputs:
```

```
383      - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
384        for the ith input.
385      - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
386        0 <= y[i] < C
387
388      Returns a tuple of:
389      - loss: Scalar giving the loss
390      - dx: Gradient of the loss with respect to x
391      """
392
393      probs = np.exp(x - np.max(x, axis=1, keepdims=True))
394      probs /= np.sum(probs, axis=1, keepdims=True)
395      N = x.shape[0]
396      loss = -np.sum(np.log(probs[np.arange(N), y])) / N
397      dx = probs.copy()
398      dx[np.arange(N), y] -= 1
399      dx /= N
400      return loss, dx
401
```

# layer_utils.py

```python
1    from .layers import *
2
3    """
4    This code was originally written for CS 231n at Stanford University
5    (cs231n.stanford.edu).  It has been modified in various areas for use in the
6    ECE 239AS class at UCLA.  This includes the descriptions of what code to
7    implement as well as some slight potential changes in variable names to be
8    consistent with class nomenclature.  We thank Justin Johnson & Serena Yeung for
9    permission to use this code.  To see the original version, please visit
10   cs231n.stanford.edu.
11   """
12
13   def affine_relu_forward(x, w, b):
14     """
15     Convenience layer that performs an affine transform followed by a ReLU
16
17     Inputs:
18     - x: Input to the affine layer
19     - w, b: Weights for the affine layer
20
21     Returns a tuple of:
22     - out: Output from the ReLU
23     - cache: Object to give to the backward pass
24     """
25     a, fc_cache = affine_forward(x, w, b)
26     out, relu_cache = relu_forward(a)
27     cache = (fc_cache, relu_cache)
28     return out, cache
29
30
31   def affine_relu_backward(dout, cache):
32     """
33     Backward pass for the affine-relu convenience layer
34     """
35     fc_cache, relu_cache = cache
36     da = relu_backward(dout, relu_cache)
37     dx, dw, db = affine_backward(da, fc_cache)
38     return dx, dw, db
39
40   def affine_batchnorm_relu_forward(x, w, b, gamma, beta, bn_param):
41     """
42     Convenience layer that performs an affine transform followed by batch normalization and then ReLU
43
44     Inputs:
45     - x: Input to the affine layer
46     - w, b: Weights for the affine layer
47     - gamma, beta: the gamma and beta parameters associated with this layer.
48     - bn_param: a set of parameters corresponding to the layer. Relevant mainly for testing
49
50     Returns a tuple of:
51     - out: Output from the ReLU
52     - cache: Object to give to the backward pass
53     """
54     a, fc_cache = affine_forward(x, w, b)
55     a, bn_cache = batchnorm_forward(a, gamma, beta, bn_param)
56     out, relu_cache = relu_forward(a)
57     cache = (fc_cache, bn_cache, relu_cache)
58     return out, cache
59
60   def affine_batchnorm_relu_backward(dout, cache):
61     """
62     Backward pass for the affine-batchnorm-relu convenience layer
63     """
64     fc_cache, bn_cache, relu_cache = cache
65     da = relu_backward(dout, relu_cache)
66     da, dgamma, dbeta = batchnorm_backward(da, bn_cache)
67     dx, dw, db = affine_backward(da, fc_cache)
68     return dx, dw, db, dgamma, dbeta
```

conv_layers.py

```
1    import numpy as np
2    from nndl.layers import *
3    import pdb
4
5    """
6    This code was originally written for CS 231n at Stanford University
7    (cs231n.stanford.edu).  It has been modified in various areas for use in the
8    ECE 239AS class at UCLA.  This includes the descriptions of what code to
9    implement as well as some slight potential changes in variable names to be
10   consistent with class nomenclature.  We thank Justin Johnson & Serena Yeung for
11   permission to use this code.  To see the original version, please visit
12   cs231n.stanford.edu.
13   """
14
15   def conv_forward_naive(x, w, b, conv_param):
16     """
17     A naive implementation of the forward pass for a convolutional layer.
18
19     The input consists of N data points, each with C channels, height H and width
20     W. We convolve each input with F different filters, where each filter spans
21     all C channels and has height HH and width HH.
22
23     Input:
24     - x: Input data of shape (N, C, H, W)
25     - w: Filter weights of shape (F, C, HH, WW)
26     - b: Biases, of shape (F,)
27     - conv_param: A dictionary with the following keys:
28       - 'stride': The number of pixels between adjacent receptive fields in the
29         horizontal and vertical directions.
30       - 'pad': The number of pixels that will be used to zero-pad the input.
31
32     Returns a tuple of:
33     - out: Output data, of shape (N, F, H', W') where H' and W' are given by
34       H' = 1 + (H + 2 * pad - HH) / stride
35       W' = 1 + (W + 2 * pad - WW) / stride
36     - cache: (x, w, b, conv_param)
37     """
38     out = None
39     pad = conv_param['pad']
40     stride = conv_param['stride']
41
42     # ================================================================ #
43     # YOUR CODE HERE:
44     #   Implement the forward pass of a convolutional neural network.
45     #   Store the output as 'out'.
46     #   Hint: to pad the array, you can use the function np.pad.
47     # ================================================================ #
48     padded_x = np.pad(x, [(0, 0), (0, 0), (pad, pad), (pad, pad)], mode='constant')
49     out = np.zeros(shape=(x.shape[0], w.shape[0], int(1 + (x.shape[2] + 2 * pad - w.shape[2]) / stride),
50                           int(1 + (x.shape[3] + 2 * pad - w.shape[3]) / stride)))
51     for example in range(x.shape[0]):
52       for f in range(out.shape[1]):
53         for i in range(out.shape[2]):
54           for j in range(out.shape[3]):
55             out[example, f, i, j] = b[f] + np.sum(w[f] * padded_x[example, :, i * stride:i * stride + w.shape[2],
56                                                   j * stride:j * stride + w.shape[3]])
57     # ================================================================ #
58     # END YOUR CODE HERE
59     # ================================================================ #
60
61     cache = (x, w, b, conv_param)
62     return out, cache
63
64
65   def conv_backward_naive(dout, cache):
66     """
67     A naive implementation of the backward pass for a convolutional layer.
68
69     Inputs:
70     - dout: Upstream derivatives.
71     - cache: A tuple of (x, w, b, conv_param) as in conv_forward_naive
72
73     Returns a tuple of:
74     - dx: Gradient with respect to x
75     - dw: Gradient with respect to w
76     - db: Gradient with respect to b
77     """
78     dx, dw, db = None, None, None
79
80     N, F, out_height, out_width = dout.shape
81     x, w, b, conv_param = cache
82
83     stride, pad = [conv_param['stride'], conv_param['pad']]
84     xpad = np.pad(x, ((0,0), (0,0), (pad,pad), (pad,pad)), mode='constant')
85     num_filts, _, f_height, f_width = w.shape
86
87     # ================================================================ #
```

```
 88        # YOUR CODE HERE:
 89        #   Implement the backward pass of a convolutional neural network.
 90        #   Calculate the gradients: dx, dw, and db.
 91        # ================================================================ #
 92        dw = np.zeros_like(w)
 93        dx = np.zeros_like(x)
 94        db = np.zeros_like(b)
 95        dxpad = np.zeros_like(xpad)
 96
 97        for f in range(dout.shape[1]):   # F
 98          for example in range(dout.shape[0]):   # N
 99            for h_tag in range(dout.shape[2]):   # H'
100              for w_tag in range(dout.shape[3]):   # W'
101                offset_h = stride * h_tag
102                offset_w = stride * w_tag
103                dw[f] += dout[example, f, h_tag, w_tag] * xpad[example, :, offset_h:offset_h+w.shape[2],
104                                                 offset_w:offset_w+w.shape[3]]
105                dxpad[example, :, offset_h:offset_h+w.shape[2],
106                          offset_w:offset_w+w.shape[3]] += dout[example, f, h_tag, w_tag] * w[f]
107        db = np.sum(dout, axis=(0, 2, 3))
108        dx = dxpad[:, :, pad:-pad, pad:-pad]   # The padded parameters are not relevant.
109        # ================================================================ #
110        # END YOUR CODE HERE
111        # ================================================================ #
112
113        return dx, dw, db
114
115
116    def max_pool_forward_naive(x, pool_param):
117        """
118        A naive implementation of the forward pass for a max pooling layer.
119
120        Inputs:
121        - x: Input data, of shape (N, C, H, W)
122        - pool_param: dictionary with the following keys:
123          - 'pool_height': The height of each pooling region
124          - 'pool_width': The width of each pooling region
125          - 'stride': The distance between adjacent pooling regions
126
127        Returns a tuple of:
128        - out: Output data
129        - cache: (x, pool_param)
130        """
131        out = None
132
133        # ================================================================ #
134        # YOUR CODE HERE:
135        #   Implement the max pooling forward pass.
136        # ================================================================ #
137        pool_height = pool_param['pool_height']
138        pool_width = pool_param['pool_width']
139        stride = pool_param['stride']
140
141        out_height = int((x.shape[2] - pool_height) / stride) + 1
142        out_width = int((x.shape[3] - pool_width) / stride) + 1
143        out = np.zeros(shape=(x.shape[0], x.shape[1], out_height, out_width))
144
145        for example in range(out.shape[0]):
146          for c in range(out.shape[1]):
147            for h in range(out.shape[2]):
148              for w in range(out.shape[3]):
149                out[example, c, h, w] = \
150                  np.amax(x[example, c, h * stride:h * stride + pool_height, w * stride:w * stride + pool_width])
151        # ================================================================ #
152        # END YOUR CODE HERE
153        # ================================================================ #
154        cache = (x, pool_param)
155        return out, cache
156
157    def max_pool_backward_naive(dout, cache):
158        """
159        A naive implementation of the backward pass for a max pooling layer.
160
161        Inputs:
162        - dout: Upstream derivatives
163        - cache: A tuple of (x, pool_param) as in the forward pass.
164
165        Returns:
166        - dx: Gradient with respect to x
167        """
168        dx = None
169        x, pool_param = cache
170        pool_height, pool_width, stride = pool_param['pool_height'], pool_param['pool_width'], pool_param['stride']
171
172        # ================================================================ #
173        # YOUR CODE HERE:
174        #   Implement the max pooling backward pass.
175        # ================================================================ #
176        dx = np.zeros_like(x)
177        for f in range(dout.shape[1]):   # F
178          for example in range(dout.shape[0]):   # N
```

```python
179              for h_tag in range(dout.shape[2]):  # H'
180                for w_tag in range(dout.shape[3]):  # W'
181                    pool_window = x[example, f, h_tag * stride:h_tag * stride + pool_height, w_tag * stride:w_tag * stride + pool_width]
182                    h_max_index, w_max_index = np.unravel_index(np.argmax(pool_window, axis=None), pool_window.shape)
183                    dx[example, f, h_max_index + h_tag*stride, w_max_index + w_tag*stride] = dout[example, f, h_tag, w_tag]
184
185      # ================================================================ #
186      # END YOUR CODE HERE
187      # ================================================================ #
188
189      return dx
190
191  def spatial_batchnorm_forward(x, gamma, beta, bn_param):
192      """
193      Computes the forward pass for spatial batch normalization.
194
195      Inputs:
196      - x: Input data of shape (N, C, H, W)
197      - gamma: Scale parameter, of shape (C,)
198      - beta: Shift parameter, of shape (C,)
199      - bn_param: Dictionary with the following keys:
200        - mode: 'train' or 'test'; required
201        - eps: Constant for numeric stability
202        - momentum: Constant for running mean / variance. momentum=0 means that
203          old information is discarded completely at every time step, while
204          momentum=1 means that new information is never incorporated. The
205          default of momentum=0.9 should work well in most situations.
206        - running_mean: Array of shape (D,) giving running mean of features
207        - running_var Array of shape (D,) giving running variance of features
208
209      Returns a tuple of:
210      - out: Output data, of shape (N, C, H, W)
211      - cache: Values needed for the backward pass
212      """
213      out, cache = None, None
214
215      # ================================================================ #
216      # YOUR CODE HERE:
217      #    Implement the spatial batchnorm forward pass.
218      #
219      #    You may find it useful to use the batchnorm forward pass you
220      #    implemented in HW #4.
221      # ================================================================ #
222      x_hat = x.swapaxes(0, 1).reshape(x.shape[1], -1).T
223      out, cache = np.array(batchnorm_forward(x_hat, gamma, beta, bn_param))
224      out = out.reshape(x.shape[0], x.shape[2], x.shape[3], -1).swapaxes(1, 3).swapaxes(2, 3)
225
226      # ================================================================ #
227      # END YOUR CODE HERE
228      # ================================================================ #
229
230      return out, cache
231
232
233  def spatial_batchnorm_backward(dout, cache):
234      """
235      Computes the backward pass for spatial batch normalization.
236
237      Inputs:
238      - dout: Upstream derivatives, of shape (N, C, H, W)
239      - cache: Values from the forward pass
240
241      Returns a tuple of:
242      - dx: Gradient with respect to inputs, of shape (N, C, H, W)
243      - dgamma: Gradient with respect to scale parameter, of shape (C,)
244      - dbeta: Gradient with respect to shift parameter, of shape (C,)
245      """
246      dx, dgamma, dbeta = None, None, None
247
248      # ================================================================ #
249      # YOUR CODE HERE:
250      #    Implement the spatial batchnorm backward pass.
251      #
252      #    You may find it useful to use the batchnorm forward pass you
253      #    implemented in HW #4.
254      # ================================================================ #
255
256      dout_hat = dout.swapaxes(0, 1).reshape(dout.shape[1], -1).T
257      dx, dgamma, dbeta = np.array(batchnorm_backward(dout_hat, cache))
258      dx = dx.reshape(dout.shape[0], dout.shape[2], dout.shape[3], -1).swapaxes(1, 3).swapaxes(2, 3)
259
260      # ================================================================ #
261      # END YOUR CODE HERE
262      # ================================================================ #
263
264      return dx, dgamma, dbeta
```

# conv_layer_utils.py

```python
1   from nndl.layers import *
2   from cs231n.fast_layers import *
3
4   """
5   This code was originally written for CS 231n at Stanford University
6   (cs231n.stanford.edu).  It has been modified in various areas for use in the
7   ECE 239AS class at UCLA.  This includes the descriptions of what code to
8   implement as well as some slight potential changes in variable names to be
9   consistent with class nomenclature.  We thank Justin Johnson & Serena Yeung for
10  permission to use this code.  To see the original version, please visit
11  cs231n.stanford.edu.
12  """
13
14
15  def conv_relu_forward(x, w, b, conv_param):
16    """
17    A convenience layer that performs a convolution followed by a ReLU.
18
19    Inputs:
20    - x: Input to the convolutional layer
21    - w, b, conv_param: Weights and parameters for the convolutional layer
22
23    Returns a tuple of:
24    - out: Output from the ReLU
25    - cache: Object to give to the backward pass
26    """
27    a, conv_cache = conv_forward_fast(x, w, b, conv_param)
28    out, relu_cache = relu_forward(a)
29    cache = (conv_cache, relu_cache)
30    return out, cache
31
32
33  def conv_relu_backward(dout, cache):
34    """
35    Backward pass for the conv-relu convenience layer.
36    """
37    conv_cache, relu_cache = cache
38    da = relu_backward(dout, relu_cache)
39    dx, dw, db = conv_backward_fast(da, conv_cache)
40    return dx, dw, db
41
42
43  def conv_relu_pool_forward(x, w, b, conv_param, pool_param):
44    """
45    Convenience layer that performs a convolution, a ReLU, and a pool.
46
47    Inputs:
48    - x: Input to the convolutional layer
49    - w, b, conv_param: Weights and parameters for the convolutional layer
50    - pool_param: Parameters for the pooling layer
51
52    Returns a tuple of:
53    - out: Output from the pooling layer
54    - cache: Object to give to the backward pass
55    """
56    a, conv_cache = conv_forward_fast(x, w, b, conv_param)
57    s, relu_cache = relu_forward(a)
58    out, pool_cache = max_pool_forward_fast(s, pool_param)
59    cache = (conv_cache, relu_cache, pool_cache)
60    return out, cache
61
62
```

```python
63    def conv_relu_pool_backward(dout, cache):
64      """
65      Backward pass for the conv-relu-pool convenience layer
66      """
67      conv_cache, relu_cache, pool_cache = cache
68      ds = max_pool_backward_fast(dout, pool_cache)
69      da = relu_backward(ds, relu_cache)
70      dx, dw, db = conv_backward_fast(da, conv_cache)
71      return dx, dw, db
```

cnn.py

```python
1    import numpy as np
2
3    from nndl.layers import *
4    from nndl.conv_layers import *
5    from cs231n.fast_layers import *
6    from nndl.layer_utils import *
7    from nndl.conv_layer_utils import *
8
9    import pdb
10
11   """
12   This code was originally written for CS 231n at Stanford University
13   (cs231n.stanford.edu).  It has been modified in various areas for use in the
14   ECE 239AS class at UCLA.  This includes the descriptions of what code to
15   implement as well as some slight potential changes in variable names to be
16   consistent with class nomenclature.  We thank Justin Johnson & Serena Yeung for
17   permission to use this code.  To see the original version, please visit
18   cs231n.stanford.edu.
19   """
20
21   class ThreeLayerConvNet(object):
22     """
23     A three-layer convolutional network with the following architecture:
24
25     conv - relu - 2x2 max pool - affine - relu - affine - softmax
26
27     The network operates on minibatches of data that have shape (N, C, H, W)
28     consisting of N images, each with height H and width W and with C input
29     channels.
30     """
31
32     def __init__(self, input_dim=(3, 32, 32), num_filters=32, filter_size=7,
33                  hidden_dim=100, num_classes=10, weight_scale=1e-3, reg=0.0,
34                  dtype=np.float32, use_batchnorm=False):
35       """
36       Initialize a new network.
37
38       Inputs:
39       - input_dim: Tuple (C, H, W) giving size of input data
40       - num_filters: Number of filters to use in the convolutional layer
41       - filter_size: Size of filters to use in the convolutional layer
42       - hidden_dim: Number of units to use in the fully-connected hidden layer
43       - num_classes: Number of scores to produce from the final affine layer.
44       - weight_scale: Scalar giving standard deviation for random initialization
45         of weights.
46       - reg: Scalar giving L2 regularization strength
47       - dtype: numpy datatype to use for computation.
48       """
49       self.use_batchnorm = use_batchnorm
50       self.params = {}
51       self.reg = reg
52       self.dtype = dtype
53
54
55       # ================================================================ #
56       # YOUR CODE HERE:
57       #   Initialize the weights and biases of a three layer CNN. To initialize:
58       #     - the biases should be initialized to zeros.
59       #     - the weights should be initialized to a matrix with entries
60       #         drawn from a Gaussian distribution with zero mean and
61       #         standard deviation given by weight_scale.
62       # ================================================================ #
63
64       self.params['b1'] = np.zeros((num_filters))
65       self.params['b2'] = np.zeros((hidden_dim))
66       self.params['b3'] = np.zeros((num_classes))
67       self.params['W1'] = np.random.normal(0, weight_scale, size=(num_filters, input_dim[0], filter_size, filter_size))
68       self.params['W2'] = np.random.normal(0, weight_scale, size=(int(((input_dim[1] - 2) / 2 + 1) *
69                                                                   ((input_dim[2] - 2) / 2 + 1) * num_filters),
70                                                                   hidden_dim))
71       self.params['W3'] = np.random.normal(0, weight_scale, size=(hidden_dim, num_classes))
72       # ================================================================ #
73       # END YOUR CODE HERE
74       # ================================================================ #
75
76       for k, v in self.params.items():
77         self.params[k] = v.astype(dtype)
78
79
80     def loss(self, X, y=None):
81       """
```

```
 82        Evaluate loss and gradient for the three-layer convolutional network.
 83
 84        Input / output: Same API as TwoLayerNet in fc_net.py.
 85        """
 86        W1, b1 = self.params['W1'], self.params['b1']
 87        W2, b2 = self.params['W2'], self.params['b2']
 88        W3, b3 = self.params['W3'], self.params['b3']
 89
 90        # pass conv_param to the forward pass for the convolutional layer
 91        filter_size = W1.shape[2]
 92        conv_param = {'stride': 1, 'pad': (filter_size - 1) / 2}
 93
 94        # pass pool_param to the forward pass for the max-pooling layer
 95        pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}
 96
 97        scores = None
 98
 99        # ================================================================= #
100        # YOUR CODE HERE:
101        #   Implement the forward pass of the three layer CNN.  Store the output
102        #   scores as the variable "scores".
103        # ================================================================= #
104        conv_out, conv_cache = conv_relu_pool_forward(X, W1, b1, conv_param, pool_param)
105        affine1_out, affine1_cache = affine_relu_forward(conv_out, W2, b2)
106        scores, affine2_cache = affine_forward(affine1_out, W3, b3)
107        # ================================================================= #
108        # END YOUR CODE HERE
109        # ================================================================= #
110
111        if y is None:
112            return scores
113
114        loss, grads = 0, {}
115        # ================================================================= #
116        # YOUR CODE HERE:
117        #   Implement the backward pass of the three layer CNN.  Store the grads
118        #   in the grads dictionary, exactly as before (i.e., the gradient of
119        #   self.params[k] will be grads[k]).  Store the loss as "loss", and
120        #   don't forget to add regularization on ALL weight matrices.
121        # ================================================================= #
122        loss, dout = softmax_loss(scores, y)
123        loss += 0.5 * self.reg * ((np.linalg.norm(self.params['W1'] ** 2)) + (np.linalg.norm(self.params['W2'] ** 2)) +
124                                  (np.linalg.norm(self.params['W3'] ** 2)))
125        dout, grads['W3'], grads['b3'] = affine_backward(dout, affine2_cache)
126        dout, grads['W2'], grads['b2'] = affine_relu_backward(dout, affine1_cache)
127        dout, grads['W1'], grads['b1'] = conv_relu_pool_backward(dout, conv_cache)
128
129        grads['W3'] += self.reg * self.params['W3']
130        grads['W2'] += self.reg * self.params['W2']
131        grads['W1'] += self.reg * self.params['W1']
132        # ================================================================= #
133        # END YOUR CODE HERE
134        # ================================================================= #
135
136        return loss, grads
137
```