

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.

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

import time

import numpy as np
import matplotlib.pyplot as plt
from nndl.conv_layers import *
from utils.data_utils import get_CIFAR10_data
from utils.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
from utils.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))))
```

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 [3]:

```

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)
dw_num = eval_numerical_gradient_array(lambda w: conv_forward_naive(x, w, b, conv_param)
db_num = eval_numerical_gradient_array(lambda b: conv_forward_naive(x, w, b, conv_param)

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: 8.327343047374856e-10
dw error: 2.0139230411245234e-09
db error: 2.6692627510510427e-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 [4]:

```
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 [5]:

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

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

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

# Your error should be 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:  3.2756229496900018e-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 utils. They are provided in `utils/fast_layers.py`.

The fast convolution implementation depends on a Cython extension ('pip install Cython' to your virtual environment); to compile it you need to run the following from the `utils` 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 [6]:

```
from utils.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: 5.400120s
Fast: 0.010899s
Speedup: 495.466269x
Difference: 2.2929071387675596e-11
```

```
Testing conv_backward_fast:
Naive: 9.441069s
Fast: 0.007226s
Speedup: 1306.586399x
dx difference: 7.625551401113985e-11
dw difference: 2.5109424891537813e-12
db difference: 3.432218184474099e-15
```

In [7]:

```
from utils.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: 0.574383s

fast: 0.004920s

speedup: 116.738722x

difference: 0.0

Testing pool_backward_fast:

Naive: 0.667588s

speedup: 55.241221x

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 [8]:

```

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), x, dx)
dw_num = eval_numerical_gradient_array(lambda w: conv_relu_pool_forward(x, w, b, conv_param), w, dw)
db_num = eval_numerical_gradient_array(lambda b: conv_relu_pool_forward(x, w, b, conv_param), b, db)

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: 8.137259622907255e-09
dw error: 7.179399022605902e-10
db error: 2.1927988059877098e-11

```

In [9]:

```

from nnlib.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), x, dx)
dw_num = eval_numerical_gradient_array(lambda w: conv_relu_forward(x, w, b, conv_param), w, dw)
db_num = eval_numerical_gradient_array(lambda b: conv_relu_forward(x, w, b, conv_param), b, db)

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: 2.7306685391136505e-08
dw error: 6.2689553780425465e-09
db error: 3.028745186715604e-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.

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 $(N \cdot H \cdot W, C)$ 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.

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

import time
import numpy as np
import matplotlib.pyplot as plt
from nndl.conv_layers import *
from utils.data_utils import get_CIFAR10_data
from utils.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
from utils.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))))
```

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 [2]:

```
# 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: [ 8.40782568 10.02877281 10.73764096]
  Stds:  [3.97430184 4.68968765 3.25480695]
After spatial batch normalization:
  Shape: (2, 3, 4, 5)
  Means: [-1.94289029e-17 -5.60662627e-16  4.57966998e-17]
  Stds:  [0.99999968 0.99999977 0.99999953]
After spatial batch normalization (nontrivial gamma, beta):
  Shape: (2, 3, 4, 5)
  Means: [6. 7. 8.]
  Stds:  [2.99999905 3.99999909 4.99999764]
```

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 [3]:

```
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:  4.161995488931059e-09
dgamma error:  6.005566401346872e-12
dbeta error:  6.306148566577526e-12
```

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.

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.

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

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

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

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

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

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

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

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

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 [4]: 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_name], verbose=False,
    e = rel_error(param_grad_num, grads[param_name])
    print('{} max relative error: {}'.format(param_name, rel_error(param_grad_num, grads

W1 max relative error: 0.002268986592757279
W2 max relative error: 0.0021176341224491218
W3 max relative error: 5.226928034487663e-06
b1 max relative error: 4.196696704863285e-06
b2 max relative error: 1.487895884269271e-07
b3 max relative error: 1.4628540331904963e-09
```

Overfit small dataset

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

```
In [5]: 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: 2.380656
(Epoch 0 / 10) train acc: 0.230000; val_acc: 0.159000
(Iteration 2 / 20) loss: 4.038913
(Epoch 1 / 10) train acc: 0.240000; val_acc: 0.126000
(Iteration 3 / 20) loss: 2.483374
(Iteration 4 / 20) loss: 2.097788
(Epoch 2 / 10) train acc: 0.230000; val_acc: 0.142000
(Iteration 5 / 20) loss: 2.446642
(Iteration 6 / 20) loss: 2.437417
(Epoch 3 / 10) train acc: 0.470000; val_acc: 0.184000
(Iteration 7 / 20) loss: 1.723952
(Iteration 8 / 20) loss: 1.821720
(Epoch 4 / 10) train acc: 0.450000; val_acc: 0.178000
(Iteration 9 / 20) loss: 1.672315
(Iteration 10 / 20) loss: 1.241660
(Epoch 5 / 10) train acc: 0.600000; val_acc: 0.189000
(Iteration 11 / 20) loss: 1.335591
(Iteration 12 / 20) loss: 1.040599
(Epoch 6 / 10) train acc: 0.650000; val_acc: 0.195000
(Iteration 13 / 20) loss: 1.166686
(Iteration 14 / 20) loss: 0.972074
(Epoch 7 / 10) train acc: 0.710000; val_acc: 0.200000
(Iteration 15 / 20) loss: 1.166310
(Iteration 16 / 20) loss: 0.747452
(Epoch 8 / 10) train acc: 0.790000; val_acc: 0.204000
(Iteration 17 / 20) loss: 0.725688
(Iteration 18 / 20) loss: 0.794187
(Epoch 9 / 10) train acc: 0.760000; val_acc: 0.180000
(Iteration 19 / 20) loss: 0.433434
(Iteration 20 / 20) loss: 0.691781
(Epoch 10 / 10) train acc: 0.850000; val_acc: 0.198000

```

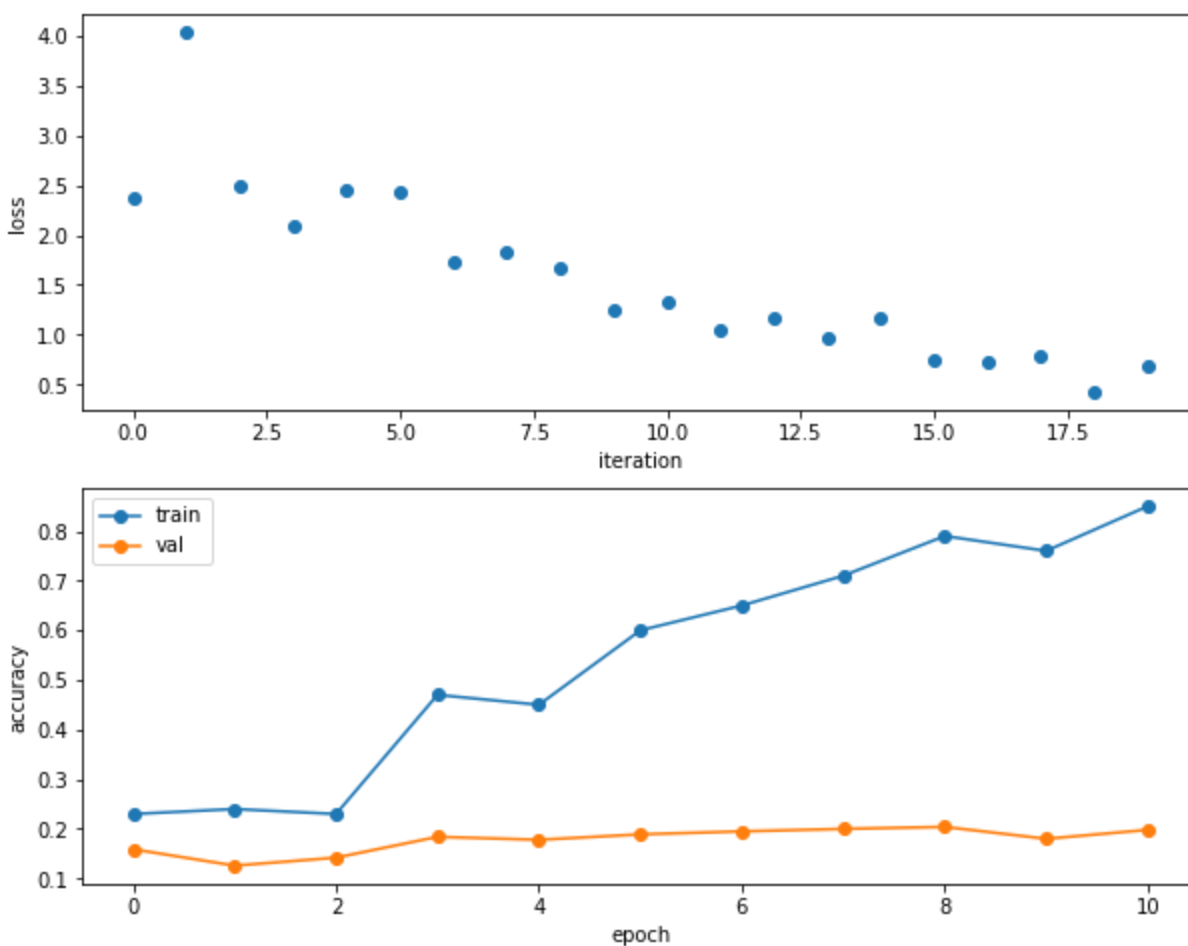
In [6]:

```

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 [7]: model = ThreeLayerConvNet(weight_scale=0.001, hidden_dim=500, reg=0.001)

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

solver.train()
```

```
(Iteration 1 / 980) loss: 2.304783
(Epoch 0 / 1) train acc: 0.102000; val_acc: 0.113000
(Iteration 21 / 980) loss: 2.074944
(Iteration 41 / 980) loss: 2.332429
(Iteration 61 / 980) loss: 1.825277
(Iteration 81 / 980) loss: 1.835692
(Iteration 101 / 980) loss: 1.916303
(Iteration 121 / 980) loss: 1.621939
(Iteration 141 / 980) loss: 1.465717
(Iteration 161 / 980) loss: 1.819977
(Iteration 181 / 980) loss: 1.743110
(Iteration 201 / 980) loss: 1.670612
(Iteration 221 / 980) loss: 1.865526
(Iteration 241 / 980) loss: 1.698929
(Iteration 261 / 980) loss: 1.720045
(Iteration 281 / 980) loss: 1.807330
(Iteration 301 / 980) loss: 1.909947
(Iteration 321 / 980) loss: 1.593170
(Iteration 341 / 980) loss: 1.849560
(Iteration 361 / 980) loss: 1.610693
(Iteration 381 / 980) loss: 1.802914
(Iteration 401 / 980) loss: 1.591736
(Iteration 421 / 980) loss: 1.561139
(Iteration 441 / 980) loss: 1.604666
(Iteration 461 / 980) loss: 1.673824
(Iteration 481 / 980) loss: 1.613188
(Iteration 501 / 980) loss: 1.643375
(Iteration 521 / 980) loss: 1.579830
(Iteration 541 / 980) loss: 1.431830
(Iteration 561 / 980) loss: 1.558496
(Iteration 581 / 980) loss: 1.815331
(Iteration 601 / 980) loss: 1.741974
(Iteration 621 / 980) loss: 1.871260
(Iteration 641 / 980) loss: 1.741922
(Iteration 661 / 980) loss: 1.948502
(Iteration 681 / 980) loss: 1.667782
(Iteration 701 / 980) loss: 1.505128
(Iteration 721 / 980) loss: 1.511095
(Iteration 741 / 980) loss: 2.047607
(Iteration 761 / 980) loss: 1.614136
(Iteration 781 / 980) loss: 1.481369
(Iteration 801 / 980) loss: 1.472241
(Iteration 821 / 980) loss: 1.406312
(Iteration 841 / 980) loss: 1.565292
(Iteration 861 / 980) loss: 1.417483
(Iteration 881 / 980) loss: 1.277716
(Iteration 901 / 980) loss: 1.976424
(Iteration 921 / 980) loss: 1.820521
(Iteration 941 / 980) loss: 1.568014
(Iteration 961 / 980) loss: 1.555147
(Epoch 1 / 1) train acc: 0.472000; val_acc: 0.480000
```

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 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 [9]:

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

model = ThreeLayerConvNet(weight_scale = 0.001, hidden_dim = 500, reg = 0.001,
                           filter_size = 3, num_filters = 64)

solver = Solver(model, data,
                 num_epochs = 10, batch_size = 500,
                 update_rule = 'adam',
                 optim_config = {
                     'learning_rate': 1e-3,
                 },
                 lr_decay = 0.9,
                 verbose = True, print_every = 10)
solver.train()

y_test_pred = np.argmax(model.loss(data['X_test']), axis = 1)
y_val_pred = np.argmax(model.loss(data['X_val']), axis = 1)
print("\nTest set accuracy: {}".format(np.mean(np.equal(y_test_pred, data['y_test']))))
print("Validation set accuracy: {}".format(np.mean(np.equal(y_val_pred, data['y_val']))))

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

```
(Iteration 1 / 980) loss: 2.306681
(Epoch 0 / 10) train acc: 0.130000; val_acc: 0.143000
(Iteration 11 / 980) loss: 2.090823
(Iteration 21 / 980) loss: 1.785177
(Iteration 31 / 980) loss: 1.717308
(Iteration 41 / 980) loss: 1.634953
(Iteration 51 / 980) loss: 1.600705
(Iteration 61 / 980) loss: 1.453167
(Iteration 71 / 980) loss: 1.476606
(Iteration 81 / 980) loss: 1.370595
(Iteration 91 / 980) loss: 1.325285
(Epoch 1 / 10) train acc: 0.553000; val_acc: 0.544000
(Iteration 101 / 980) loss: 1.269836
(Iteration 111 / 980) loss: 1.338483
(Iteration 121 / 980) loss: 1.421310
(Iteration 131 / 980) loss: 1.294719
(Iteration 141 / 980) loss: 1.276342
(Iteration 151 / 980) loss: 1.296264
(Iteration 161 / 980) loss: 1.187775
(Iteration 171 / 980) loss: 1.259028
(Iteration 181 / 980) loss: 1.176344
(Iteration 191 / 980) loss: 1.122214
(Epoch 2 / 10) train acc: 0.611000; val_acc: 0.613000
(Iteration 201 / 980) loss: 1.262361
(Iteration 211 / 980) loss: 1.077568
(Iteration 221 / 980) loss: 1.266905
(Iteration 231 / 980) loss: 1.098652
(Iteration 241 / 980) loss: 1.106851
(Iteration 251 / 980) loss: 1.120836
(Iteration 261 / 980) loss: 1.132309
(Iteration 271 / 980) loss: 1.057133
(Iteration 281 / 980) loss: 1.165886
(Iteration 291 / 980) loss: 1.102857
(Epoch 3 / 10) train acc: 0.664000; val_acc: 0.620000
(Iteration 301 / 980) loss: 1.011122
(Iteration 311 / 980) loss: 1.019898
(Iteration 321 / 980) loss: 1.058973
(Iteration 331 / 980) loss: 1.053693
(Iteration 341 / 980) loss: 0.943551
(Iteration 351 / 980) loss: 0.985662
(Iteration 361 / 980) loss: 1.048992
(Iteration 371 / 980) loss: 1.061378
(Iteration 381 / 980) loss: 0.908078
(Iteration 391 / 980) loss: 0.883494
(Epoch 4 / 10) train acc: 0.721000; val_acc: 0.647000
(Iteration 401 / 980) loss: 0.966984
(Iteration 411 / 980) loss: 0.903183
(Iteration 421 / 980) loss: 0.768672
(Iteration 431 / 980) loss: 0.808488
(Iteration 441 / 980) loss: 0.918113
(Iteration 451 / 980) loss: 0.880160
(Iteration 461 / 980) loss: 0.819261
(Iteration 471 / 980) loss: 0.791451
(Iteration 481 / 980) loss: 0.758962
(Epoch 5 / 10) train acc: 0.770000; val_acc: 0.663000
(Iteration 491 / 980) loss: 0.742991
(Iteration 501 / 980) loss: 0.719659
(Iteration 511 / 980) loss: 0.759917
(Iteration 521 / 980) loss: 0.826521
(Iteration 531 / 980) loss: 0.856099
(Iteration 541 / 980) loss: 0.798033
(Iteration 551 / 980) loss: 0.742056
(Iteration 561 / 980) loss: 0.807927
(Iteration 571 / 980) loss: 0.779540
(Iteration 581 / 980) loss: 0.809450
(Epoch 6 / 10) train acc: 0.813000; val_acc: 0.670000
```



```
(Iteration 591 / 980) loss: 0.804336
(Iteration 601 / 980) loss: 0.796321
(Iteration 611 / 980) loss: 0.681944
(Iteration 621 / 980) loss: 0.739092
(Iteration 631 / 980) loss: 0.810331
(Iteration 641 / 980) loss: 0.730700
(Iteration 651 / 980) loss: 0.752555
(Iteration 661 / 980) loss: 0.757851
(Iteration 671 / 980) loss: 0.759011
(Iteration 681 / 980) loss: 0.676632
(Epoch 7 / 10) train acc: 0.803000; val_acc: 0.661000
(Iteration 691 / 980) loss: 0.597587
(Iteration 701 / 980) loss: 0.677737
(Iteration 711 / 980) loss: 0.667884
(Iteration 721 / 980) loss: 0.701386
(Iteration 731 / 980) loss: 0.641914
(Iteration 741 / 980) loss: 0.561046
(Iteration 751 / 980) loss: 0.585528
(Iteration 761 / 980) loss: 0.551552
(Iteration 771 / 980) loss: 0.715700
(Iteration 781 / 980) loss: 0.668288
(Epoch 8 / 10) train acc: 0.832000; val_acc: 0.657000
(Iteration 791 / 980) loss: 0.574248
(Iteration 801 / 980) loss: 0.574751
(Iteration 811 / 980) loss: 0.492692
(Iteration 821 / 980) loss: 0.559167
(Iteration 831 / 980) loss: 0.587724
(Iteration 841 / 980) loss: 0.555173
(Iteration 851 / 980) loss: 0.549317
(Iteration 861 / 980) loss: 0.555267
(Iteration 871 / 980) loss: 0.553400
(Iteration 881 / 980) loss: 0.520094
(Epoch 9 / 10) train acc: 0.866000; val_acc: 0.671000
(Iteration 891 / 980) loss: 0.565403
(Iteration 901 / 980) loss: 0.551831
(Iteration 911 / 980) loss: 0.498214
(Iteration 921 / 980) loss: 0.462239
(Iteration 931 / 980) loss: 0.474712
(Iteration 941 / 980) loss: 0.432006
(Iteration 951 / 980) loss: 0.470267
(Iteration 961 / 980) loss: 0.509695
(Iteration 971 / 980) loss: 0.490025
(Epoch 10 / 10) train acc: 0.884000; val_acc: 0.675000
```

Test set accuracy: 0.678

Validation set accuracy: 0.675

```

1 import numpy as np
2 from nnutils.layers import *
3 import pdb
4
5
6 def conv_forward_naive(x, w, b, conv_param):
7     """
8     A naive implementation of the forward pass for a convolutional layer.
9
10    The input consists of N data points, each with C channels, height H and
    width
11    W. We convolve each input with F different filters, where each filter spans
12    all C channels and has height HH and width WW.
13
14    Input:
15    - x: Input data of shape (N, C, H, W)
16    - w: Filter weights of shape (F, C, HH, WW)
17    - b: Biases, of shape (F,)
18    - conv_param: A dictionary with the following keys:
19        - 'stride': The number of pixels between adjacent receptive fields in the
20          horizontal and vertical directions.
21        - 'pad': The number of pixels that will be used to zero-pad the input.
22
23    Returns a tuple of:
24    - out: Output data, of shape (N, F, H', W') where H' and W' are given by
25        H' = 1 + (H + 2 * pad - HH) / stride
26        W' = 1 + (W + 2 * pad - WW) / stride
27    - cache: (x, w, b, conv_param)
28    """
29    out = None
30    pad = conv_param['pad']
31    stride = conv_param['stride']
32
33    # ===== #
34    # YOUR CODE HERE:
35    # Implement the forward pass of a convolutional neural network.
36    # Store the output as 'out'.
37    # Hint: to pad the array, you can use the function np.pad.
38    # ===== #
39
40    x_pad = np.pad(x, [(0, 0), (0, 0), (pad, pad), (pad, pad)],
41mode='constant')
42
43    N, C, H, W = x.shape
44    F, C, HH, WW = w.shape
45
46    H2 = int(1 + (H + 2 * pad - HH) / stride)
47    W2 = int(1 + (W + 2 * pad - WW) / stride)
48
49    out = np.zeros([N, F, H2, W2])
50
51    for n in np.arange(N):
52        for f in np.arange(F):
53            for row in np.arange(H2):
54                for col in np.arange(W2):
55                    out[n, f, row, col] = np.sum(w[f, :, :, :] * x_pad[n, :, row*stride
56: row*stride+HH, col*stride : col*stride+WW]) + b[f]
57
58    # ===== #

```

```

57 # END YOUR CODE HERE
58 # ===== #
59
60 cache = (x, w, b, conv_param)
61 return out, cache
62
63
64 def conv_backward_naive(dout, cache):
65     """
66     A naive implementation of the backward pass for a convolutional layer.
67
68     Inputs:
69     - dout: Upstream derivatives.
70     - cache: A tuple of (x, w, b, conv_param) as in conv_forward_naive
71
72     Returns a tuple of:
73     - dx: Gradient with respect to x
74     - dw: Gradient with respect to w
75     - db: Gradient with respect to b
76     """
77     dx, dw, db = None, None, None
78
79     N, F, out_height, out_width = dout.shape
80     x, w, b, conv_param = cache
81
82     stride, pad = [conv_param['stride'], conv_param['pad']]
83     xpad = np.pad(x, ((0,0), (0,0), (pad,pad), (pad,pad)), mode='constant')
84     num_filts, _, f_height, f_width = w.shape
85
86     # ===== #
87     # YOUR CODE HERE:
88     #   Implement the backward pass of a convolutional neural network.
89     #   Calculate the gradients: dx, dw, and db.
90     # ===== #
91
92     H = x.shape[2]
93     W = x.shape[3]
94
95     dx_pad = np.zeros_like(xpad)
96     dw = np.zeros_like(w)
97     db = np.zeros_like(b)
98
99     # dx
100    for n in np.arange(N):
101        for f in np.arange(F):
102            for row in np.arange(out_height):
103                for col in np.arange(out_width):
104                    dx_pad[n, :, row*stride : row*stride+f_height, col*stride :
col*stride+f_width] += dout[n, f, row, col] * w[f, :, :, :]
105                    dx = dx_pad[:, :, pad : pad+H, pad : pad+W]
106
107    # dw
108    for n in np.arange(N):
109        for f in np.arange(F):
110            for row in np.arange(out_height):
111                for col in np.arange(out_width):
112                    dw[f, :, :, :] += dout[n, f, row, col] * xpad[n, :, row*stride :
row*stride+f_height, col*stride : col*stride+f_width]
113
114    # db

```

```

115     for f in np.arange(F):
116         db[f] += np.sum(dout[:, f, :, :])
117
118     # ===== #
119     # END YOUR CODE HERE
120     # ===== #
121
122     return dx, dw, db
123
124
125 def max_pool_forward_naive(x, pool_param):
126     """
127     A naive implementation of the forward pass for a max pooling layer.
128
129     Inputs:
130     - x: Input data, of shape (N, C, H, W)
131     - pool_param: dictionary with the following keys:
132         - 'pool_height': The height of each pooling region
133         - 'pool_width': The width of each pooling region
134         - 'stride': The distance between adjacent pooling regions
135
136     Returns a tuple of:
137     - out: Output data
138     - cache: (x, pool_param)
139     """
140     out = None
141
142     # ===== #
143     # YOUR CODE HERE:
144     #   Implement the max pooling forward pass.
145     # ===== #
146
147     pool_height = pool_param['pool_height']
148     pool_width = pool_param['pool_width']
149     stride = pool_param['stride']
150
151     N, C, H, W = x.shape
152
153     H2 = int(1 + (H - pool_height) / stride)
154     W2 = int(1 + (W - pool_width) / stride)
155
156     out = np.zeros([N, C, H2, W2])
157
158     for n in np.arange(N):
159         for c in np.arange(C):
160             for row in np.arange(H2):
161                 for col in np.arange(W2):
162                     out[n, c, row, col] = np.max(x[n, c, row*stride :
163 row*stride+pool_height, col*stride : col*stride+pool_width])
164
165     # ===== #
166     # END YOUR CODE HERE
167     # ===== #
168     cache = (x, pool_param)
169     return out, cache
170
171 def max_pool_backward_naive(dout, cache):
172     """
173     A naive implementation of the backward pass for a max pooling layer.

```

```

174     Inputs:
175     - dout: Upstream derivatives
176     - cache: A tuple of (x, pool_param) as in the forward pass.
177
178     Returns:
179     - dx: Gradient with respect to x
180     """
181     dx = None
182     x, pool_param = cache
183     pool_height, pool_width, stride = pool_param['pool_height'],
pool_param['pool_width'], pool_param['stride']
184
185     # ===== #
186     # YOUR CODE HERE:
187     #   Implement the max pooling backward pass.
188     # ===== #
189
190     N, C, H, W = x.shape
191     out_height = dout.shape[2]
192     out_width = dout.shape[3]
193
194     dx = np.zeros_like(x)
195
196     for n in np.arange(N):
197         for c in np.arange(C):
198             for row in np.arange(out_height):
199                 for col in np.arange(out_width):
200                     max_idx = np.unravel_index(np.argmax(x[n, c, row*stride :
row*stride+pool_height, col*stride : col*stride+pool_width]), [pool_height,
pool_width])
201                     dx[n, c, row*stride+max_idx[0], col*stride+max_idx[1]] = dout[n, c,
row, col]
202
203     # ===== #
204     # END YOUR CODE HERE
205     # ===== #
206
207     return dx
208
209 def spatial_batchnorm_forward(x, gamma, beta, bn_param):
210     """
211     Computes the forward pass for spatial batch normalization.
212
213     Inputs:
214     - x: Input data of shape (N, C, H, W)
215     - gamma: Scale parameter, of shape (C,)
216     - beta: Shift parameter, of shape (C,)
217     - bn_param: Dictionary with the following keys:
218         - mode: 'train' or 'test'; required
219         - eps: Constant for numeric stability
220         - momentum: Constant for running mean / variance. momentum=0 means that
221           old information is discarded completely at every time step, while
222           momentum=1 means that new information is never incorporated. The
223           default of momentum=0.9 should work well in most situations.
224         - running_mean: Array of shape (D,) giving running mean of features
225         - running_var: Array of shape (D,) giving running variance of features
226
227     Returns a tuple of:
228     - out: Output data, of shape (N, C, H, W)
229     - cache: Values needed for the backward pass

```

```

230     """
231     out, cache = None, None
232
233     # ===== #
234     # YOUR CODE HERE:
235     #   Implement the spatial batchnorm forward pass.
236     #
237     #   You may find it useful to use the batchnorm forward pass you
238     #   implemented in HW #4.
239     # ===== #
240
241     N, C, H, W = x.shape
242
243     x_reshape = np.reshape(np.transpose(x, (0, 2, 3, 1)), (N*H*W, C))
244     out_2D, cache = batchnorm_forward(x_reshape, gamma, beta, bn_param)
245
246     out = np.transpose(np.reshape(out_2D, (N, H, W, C)), (0, 3, 1, 2))
247
248     # ===== #
249     # END YOUR CODE HERE
250     # ===== #
251
252     return out, cache
253
254
255 def spatial_batchnorm_backward(dout, cache):
256     """
257     Computes the backward pass for spatial batch normalization.
258
259     Inputs:
260     - dout: Upstream derivatives, of shape (N, C, H, W)
261     - cache: Values from the forward pass
262
263     Returns a tuple of:
264     - dx: Gradient with respect to inputs, of shape (N, C, H, W)
265     - dgamma: Gradient with respect to scale parameter, of shape (C,)
266     - dbeta: Gradient with respect to shift parameter, of shape (C,)
267     """
268     dx, dgamma, dbeta = None, None, None
269
270     # ===== #
271     # YOUR CODE HERE:
272     #   Implement the spatial batchnorm backward pass.
273     #
274     #   You may find it useful to use the batchnorm forward pass you
275     #   implemented in HW #4.
276     # ===== #
277
278     N, C, H, W = dout.shape
279
280     dout_reshape = np.reshape(np.transpose(dout, (0, 2, 3, 1)), (N*H*W, C))
281     dx_2D, dgamma, dbeta = batchnorm_backward(dout_reshape, cache)
282
283     dx = np.transpose(np.reshape(dx_2D, (N, H, W, C)), (0, 3, 1, 2))
284
285     # ===== #
286     # END YOUR CODE HERE
287     # ===== #
288
289     return dx, dgamma, dbeta

```

```
1 from nndl.layers import *
2 from utils.fast_layers import *
3
4
5 def conv_relu_forward(x, w, b, conv_param):
6     """
7     A convenience layer that performs a convolution followed by a ReLU.
8
9     Inputs:
10    - x: Input to the convolutional layer
11    - w, b, conv_param: Weights and parameters for the convolutional layer
12
13    Returns a tuple of:
14    - out: Output from the ReLU
15    - cache: Object to give to the backward pass
16    """
17    a, conv_cache = conv_forward_fast(x, w, b, conv_param)
18    out, relu_cache = relu_forward(a)
19    cache = (conv_cache, relu_cache)
20    return out, cache
21
22
23 def conv_relu_backward(dout, cache):
24     """
25     Backward pass for the conv-relu convenience layer.
26     """
27    conv_cache, relu_cache = cache
28    da = relu_backward(dout, relu_cache)
29    dx, dw, db = conv_backward_fast(da, conv_cache)
30    return dx, dw, db
31
32
33 def conv_relu_pool_forward(x, w, b, conv_param, pool_param):
34     """
35     Convenience layer that performs a convolution, a ReLU, and a pool.
36
37     Inputs:
38    - x: Input to the convolutional layer
39    - w, b, conv_param: Weights and parameters for the convolutional layer
40    - pool_param: Parameters for the pooling layer
41
42    Returns a tuple of:
43    - out: Output from the pooling layer
44    - cache: Object to give to the backward pass
45    """
46    a, conv_cache = conv_forward_fast(x, w, b, conv_param)
47    s, relu_cache = relu_forward(a)
48    out, pool_cache = max_pool_forward_fast(s, pool_param)
49    cache = (conv_cache, relu_cache, pool_cache)
50    return out, cache
51
52
53 def conv_relu_pool_backward(dout, cache):
54     """
55     Backward pass for the conv-relu-pool convenience layer
56     """
57    conv_cache, relu_cache, pool_cache = cache
58    ds = max_pool_backward_fast(dout, pool_cache)
59    da = relu_backward(ds, relu_cache)
```

```
60 | dx, dw, db = conv_backward_fast(da, conv_cache)
61 | return dx, dw, db
```



```

1 import numpy as np
2 import pdb
3
4 def affine_forward(x, w, b):
5     """
6     Computes the forward pass for an affine (fully-connected) layer.
7
8     The input x has shape (N, d_1, ..., d_k) and contains a minibatch of N
9     examples, where each example x[i] has shape (d_1, ..., d_k). We will
10    reshape each input into a vector of dimension D = d_1 * ... * d_k, and
11    then transform it to an output vector of dimension M.
12
13    Inputs:
14    - x: A numpy array containing input data, of shape (N, d_1, ..., d_k)
15    - w: A numpy array of weights, of shape (D, M)
16    - b: A numpy array of biases, of shape (M,)
17
18    Returns a tuple of:
19    - out: output, of shape (N, M)
20    - cache: (x, w, b)
21    """
22
23    # ===== #
24    # YOUR CODE HERE:
25    # Calculate the output of the forward pass. Notice the dimensions
26    # of w are D x M, which is the transpose of what we did in earlier
27    # assignments.
28    # ===== #
29
30    x_reshape = x.reshape((x.shape[0], w.shape[0])) #N x D
31    out = np.dot(x_reshape, w) + b.reshape((1, b.shape[0])) #N x M
32
33    # ===== #
34    # END YOUR CODE HERE
35    # ===== #
36
37    cache = (x, w, b)
38    return out, cache
39
40
41 def affine_backward(dout, cache):
42     """
43     Computes the backward pass for an affine layer.
44
45     Inputs:
46     - dout: Upstream derivative, of shape (N, M)
47     - cache: Tuple of:
48       - x: Input data, of shape (N, d_1, ... d_k)
49       - w: Weights, of shape (D, M)
50
51     Returns a tuple of:
52     - dx: Gradient with respect to x, of shape (N, d1, ..., d_k)
53     - dw: Gradient with respect to w, of shape (D, M)
54     - db: Gradient with respect to b, of shape (M,)
55     """
56     x, w, b = cache
57     dx, dw, db = None, None, None
58
59    # ===== #

```

```

60 # YOUR CODE HERE:
61 #   Calculate the gradients for the backward pass.
62 # ===== #
63
64 x_reshape = np.reshape(x, (x.shape[0], w.shape[0])) #N x D
65 dx_reshape = np.dot(dout, w.T)
66
67 dx = np.reshape(dx_reshape, x.shape) #N x D
68 dw = np.dot(x_reshape.T, dout) #D x M
69 db = np.dot(dout.T, np.ones(x.shape[0])) #M x 1
70
71 # ===== #
72 # END YOUR CODE HERE
73 # ===== #
74
75 return dx, dw, db
76
77 def relu_forward(x):
78     """
79     Computes the forward pass for a layer of rectified linear units (ReLUs).
80
81     Input:
82     - x: Inputs, of any shape
83
84     Returns a tuple of:
85     - out: Output, of the same shape as x
86     - cache: x
87     """
88     # ===== #
89     # YOUR CODE HERE:
90     #   Implement the ReLU forward pass.
91     # ===== #
92
93     out = np.maximum(x, 0)
94
95     # ===== #
96     # END YOUR CODE HERE
97     # ===== #
98
99     cache = x
100    return out, cache
101
102
103 def relu_backward(dout, cache):
104     """
105     Computes the backward pass for a layer of rectified linear units (ReLUs).
106
107     Input:
108     - dout: Upstream derivatives, of any shape
109     - cache: Input x, of same shape as dout
110
111     Returns:
112     - dx: Gradient with respect to x
113     """
114     x = cache
115
116     # ===== #
117     # YOUR CODE HERE:
118     #   Implement the ReLU backward pass
119     # ===== #

```

```

120
121     dx = dout * (x > 0)
122
123     # ===== #
124     # END YOUR CODE HERE
125     # ===== #
126
127     return dx
128
129 def batchnorm_forward(x, gamma, beta, bn_param):
130     """
131     Forward pass for batch normalization.
132
133     During training the sample mean and (uncorrected) sample variance are
134     computed from minibatch statistics and used to normalize the incoming data.
135     During training we also keep an exponentially decaying running mean of the
136     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
151     implementation
152     of batch normalization also uses running averages.
153
154     Input:
155     - x: Data of shape (N, D)
156     - gamma: Scale parameter of shape (D,)
157     - beta: Shift parameter of shape (D,)
158     - bn_param: Dictionary with the following keys:
159       - mode: 'train' or 'test'; required
160       - eps: Constant for numeric stability
161       - momentum: Constant for running mean / variance.
162       - running_mean: Array of shape (D,) giving running mean of features
163       - running_var: Array of shape (D,) giving running variance of features
164
165     Returns a tuple of:
166     - out: of shape (N, D)
167     - cache: A tuple of values needed in the backward pass
168     """
169     mode = bn_param['mode']
170     eps = bn_param.get('eps', 1e-5)
171     momentum = bn_param.get('momentum', 0.9)
172
173     N, D = x.shape
174     running_mean = bn_param.get('running_mean', np.zeros(D, dtype=x.dtype))
175     running_var = bn_param.get('running_var', np.zeros(D, dtype=x.dtype))
176
177     out, cache = None, None
178     if mode == 'train':

```

```

178 # ===== #
179 # YOUR CODE HERE:
180 #   A few steps here:
181 #   (1) Calculate the running mean and variance of the minibatch.
182 #   (2) Normalize the activations with the running mean and variance.
183 #   (3) Scale and shift the normalized activations. Store this
184 #       as the variable 'out'
185 #   (4) Store any variables you may need for the backward pass in
186 #       the 'cache' variable.
187 # ===== #
188
189 mean = np.mean(x, axis = 0)
190 var = np.var(x, axis = 0)
191 normalize_x = (x - mean) / np.sqrt(var + eps)
192
193 running_mean = momentum * running_mean + (1 - momentum) * mean
194 running_var = momentum * running_var + (1 - momentum) * var
195
196 out = gamma * normalize_x + beta
197
198 cache = {'normalize_x': normalize_x,
199         'x_minus_mean': (x - mean),
200         'sqrt_var_eps': np.sqrt(var + eps),
201         'gamma': gamma
202         }
203
204 # ===== #
205 # END YOUR CODE HERE
206 # ===== #
207
208 elif mode == 'test':
209
210     # ===== #
211     # YOUR CODE HERE:
212     #   Calculate the testing time normalized activation. Normalize using
213     #   the running mean and variance, and then scale and shift
214     #   appropriately.
215     #   Store the output as 'out'.
216     # ===== #
217
218     normalize_x = (x - running_mean) / np.sqrt(running_var + eps)
219     out = gamma * normalize_x + beta
220
221     # ===== #
222     # END YOUR CODE HERE
223     # ===== #
224
225 else:
226     raise ValueError('Invalid forward batchnorm mode "%s"' % mode)
227
228 # Store the updated running means back into bn_param
229 bn_param['running_mean'] = running_mean
230 bn_param['running_var'] = running_var
231
232 return out, cache
233
234 def batchnorm_backward(dout, cache):
235     """
236     Backward pass for batch normalization.

```

```

237 For this implementation, you should write out a computation graph for
238 batch normalization on paper and propagate gradients backward through
239 intermediate nodes.
240
241 Inputs:
242 - dout: Upstream derivatives, of shape (N, D)
243 - cache: Variable of intermediates from batchnorm_forward.
244
245 Returns a tuple of:
246 - dx: Gradient with respect to inputs x, of shape (N, D)
247 - dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
248 - dbeta: Gradient with respect to shift parameter beta, of shape (D,)
249 """
250 dx, dgamma, dbeta = None, None, None
251
252 # ===== #
253 # YOUR CODE HERE:
254 # Implement the batchnorm backward pass, calculating dx, dgamma, and
255 # dbeta.
256 # ===== #
257
258 normalize_x = cache.get('normalize_x')
259 x_minus_mean = cache.get('x_minus_mean')
260 sqrt_var_eps = cache.get('sqrt_var_eps')
261 gamma = cache.get('gamma')
262 N = dout.shape[0]
263
264 dbeta = np.sum(dout, axis = 0)
265 dgamma = np.sum(dout * normalize_x, axis = 0)
266
267 dnormalize_x = dout * gamma
268
269 db = x_minus_mean * dnormalize_x # b = 1 / sqrt_var_eps
270 dc = (-1 / (sqrt_var_eps * sqrt_var_eps)) * db # c = sqrt_var_eps
271 de = (1 / (2 * sqrt_var_eps)) * dc # e = sqrt_var_eps * sqrt_var_eps
272 dvar = np.sum(de, axis = 0)
273
274 da = dnormalize_x / sqrt_var_eps # a = x - mu
275 dm_u = -np.sum(da, axis = 0) - 2 * np.sum(x_minus_mean, axis = 0) * dvar / N
276
277 dx = da + 2 * x_minus_mean * dvar / N + dm_u / N
278
279 # ===== #
280 # END YOUR CODE HERE
281 # ===== #
282
283 return dx, dgamma, dbeta
284
285 def dropout_forward(x, dropout_param):
286     """
287     Performs the forward pass for (inverted) dropout.
288
289     Inputs:
290     - x: Input data, of any shape
291     - dropout_param: A dictionary with the following keys:
292       - p: Dropout parameter. We drop each neuron output with probability p.
293       - mode: 'test' or 'train'. If the mode is train, then perform dropout;
294         if the mode is test, then just return the input.
295       - seed: Seed for the random number generator. Passing seed makes this

```

```

295         function deterministic, which is needed for gradient checking but not
in
296         real networks.
297
298     Outputs:
299     - out: Array of the same shape as x.
300     - cache: A tuple (dropout_param, mask). In training mode, mask is the
dropout
301     mask that was used to multiply the input; in test mode, mask is None.
302     """
303     p, mode = dropout_param['p'], dropout_param['mode']
304     if 'seed' in dropout_param:
305         np.random.seed(dropout_param['seed'])
306
307     mask = None
308     out = None
309
310     if mode == 'train':
311         # ===== #
312         # YOUR CODE HERE:
313         # Implement the inverted dropout forward pass during training time.
314         # Store the masked and scaled activations in out, and store the
315         # dropout mask as the variable mask.
316         # ===== #
317
318         mask = (np.random.rand(*x.shape) < (1 - p)) / (1 - p)
319         out = x * mask
320
321         # ===== #
322         # END YOUR CODE HERE
323         # ===== #
324
325     elif mode == 'test':
326
327         # ===== #
328         # YOUR CODE HERE:
329         # Implement the inverted dropout forward pass during test time.
330         # ===== #
331
332         out = x
333
334         # ===== #
335         # END YOUR CODE HERE
336         # ===== #
337
338     cache = (dropout_param, mask)
339     out = out.astype(x.dtype, copy=False)
340
341     return out, cache
342
343 def dropout_backward(dout, cache):
344     """
345     Perform the backward pass for (inverted) dropout.
346
347     Inputs:
348     - dout: Upstream derivatives, of any shape
349     - cache: (dropout_param, mask) from dropout_forward.
350     """
351     dropout_param, mask = cache
352     mode = dropout_param['mode']

```

```

353
354 dx = None
355 if mode == 'train':
356     # ===== #
357     # YOUR CODE HERE:
358     # Implement the inverted dropout backward pass during training time.
359     # ===== #
360
361     dx = dout * mask
362
363     # ===== #
364     # END YOUR CODE HERE
365     # ===== #
366 elif mode == 'test':
367     # ===== #
368     # YOUR CODE HERE:
369     # Implement the inverted dropout backward pass during test time.
370     # ===== #
371
372     dx = dout
373
374     # ===== #
375     # END YOUR CODE HERE
376     # ===== #
377 return dx
378
379 def svm_loss(x, y):
380     """
381     Computes the loss and gradient using for multiclass SVM classification.
382
383     Inputs:
384     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth
385 class
386     for the ith input.
387     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
388 0 <= y[i] < C
389
390 Returns a tuple of:
391 - loss: Scalar giving the loss
392 - dx: Gradient of the loss with respect to x
393 """
394 N = x.shape[0]
395 correct_class_scores = x[np.arange(N), y]
396 margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
397 margins[np.arange(N), y] = 0
398 loss = np.sum(margins) / N
399 num_pos = np.sum(margins > 0, axis=1)
400 dx = np.zeros_like(x)
401 dx[margins > 0] = 1
402 dx[np.arange(N), y] -= num_pos
403 dx /= N
404 return loss, dx
405
406 def softmax_loss(x, y):
407     """
408     Computes the loss and gradient for softmax classification.
409
410     Inputs:

```

```
411     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth
class
412     for the ith input.
413     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
414         0 <= y[i] < C
415
416     Returns a tuple of:
417     - loss: Scalar giving the loss
418     - dx: Gradient of the loss with respect to x
419     """
420
421     probs = np.exp(x - np.max(x, axis=1, keepdims=True))
422     probs /= np.sum(probs, axis=1, keepdims=True)
423     N = x.shape[0]
424     loss = -np.sum(np.log(probs[np.arange(N), y])) / N
425     dx = probs.copy()
426     dx[np.arange(N), y] -= 1
427     dx /= N
428     return loss, dx
429
```



```
1 from .layers import *
2
3 def affine_relu_forward(x, w, b):
4     """
5     Convenience layer that performs an affine transform followed by a ReLU
6
7     Inputs:
8     - x: Input to the affine layer
9     - w, b: Weights for the affine layer
10
11     Returns a tuple of:
12     - out: Output from the ReLU
13     - cache: Object to give to the backward pass
14     """
15     a, fc_cache = affine_forward(x, w, b)
16     out, relu_cache = relu_forward(a)
17     cache = (fc_cache, relu_cache)
18     return out, cache
19
20
21 def affine_relu_backward(dout, cache):
22     """
23     Backward pass for the affine-relu convenience layer
24     """
25     fc_cache, relu_cache = cache
26     da = relu_backward(dout, relu_cache)
27     dx, dw, db = affine_backward(da, fc_cache)
28     return dx, dw, db
```

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

```

55     - velocity: A numpy array of the same shape as w and dw used to store a
moving
56     average of the gradients.
57     """
58     if config is None: config = {}
59     config.setdefault('learning_rate', 1e-2)
60     config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
61     v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets
it to zero.
62
63     # ===== #
64     # YOUR CODE HERE:
65     #     Implement the momentum update formula. Return the updated weights
66     #     as next_w, and the updated velocity as v.
67     # ===== #
68
69     alpha = config['momentum']
70     eps = config['learning_rate']
71
72     v = alpha * v - eps * dw
73     w += v
74
75     next_w = w
76
77     # ===== #
78     # END YOUR CODE HERE
79     # ===== #
80
81     config['velocity'] = v
82
83     return next_w, config
84
85 def sgd_nesterov_momentum(w, dw, config=None):
86     """
87     Performs stochastic gradient descent with Nesterov momentum.
88
89     config format:
90     - learning_rate: Scalar learning rate.
91     - momentum: Scalar between 0 and 1 giving the momentum value.
92       Setting momentum = 0 reduces to sgd.
93     - velocity: A numpy array of the same shape as w and dw used to store a
moving
94     average of the gradients.
95     """
96     if config is None: config = {}
97     config.setdefault('learning_rate', 1e-2)
98     config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
99     v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets
it to zero.
100
101     # ===== #
102     # YOUR CODE HERE:
103     #     Implement the momentum update formula. Return the updated weights
104     #     as next_w, and the updated velocity as v.
105     # ===== #
106
107     alpha = config['momentum']
108     eps = config['learning_rate']
109
110     v_old = v

```

```

111     v = alpha * v - eps * dw
112     w += (v + alpha * (v - v_old))
113
114     next_w = w
115
116     # ===== #
117     # END YOUR CODE HERE
118     # ===== #
119
120     config['velocity'] = v
121
122     return next_w, config
123
124 def rmsprop(w, dw, config=None):
125     """
126     Uses the RMSProp update rule, which uses a moving average of squared
127     gradient
128     values to set adaptive per-parameter learning rates.
129
130     config format:
131     - learning_rate: Scalar learning rate.
132     - decay_rate: Scalar between 0 and 1 giving the decay rate for the squared
133       gradient cache.
134     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
135     - beta: Moving average of second moments of gradients.
136     """
137     if config is None: config = {}
138     config.setdefault('learning_rate', 1e-2)
139     config.setdefault('decay_rate', 0.99)
140     config.setdefault('epsilon', 1e-8)
141     config.setdefault('a', np.zeros_like(w))
142
143     next_w = None
144
145     # ===== #
146     # YOUR CODE HERE:
147     # Implement RMSProp. Store the next value of w as next_w. You need
148     # to also store in config['a'] the moving average of the second
149     # moment gradients, so they can be used for future gradients. Concretely,
150     # config['a'] corresponds to "a" in the lecture notes.
151     # ===== #
152
153     a = config['a']
154     beta = config['decay_rate']
155     eps = config['learning_rate']
156     nu = config['epsilon']
157
158     a = beta * a + (1 - beta) * dw * dw
159     w -= (eps * dw) / (np.sqrt(a) + nu)
160
161     config['a'] = a
162     next_w = w
163
164     # ===== #
165     # END YOUR CODE HERE
166     # ===== #
167
168     return next_w, config
169

```

```

170 def adam(w, dw, config=None):
171     """
172     Uses the Adam update rule, which incorporates moving averages of both the
173     gradient and its square and a bias correction term.
174
175     config format:
176     - learning_rate: Scalar learning rate.
177     - beta1: Decay rate for moving average of first moment of gradient.
178     - beta2: Decay rate for moving average of second moment of gradient.
179     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
180     - m: Moving average of gradient.
181     - v: Moving average of squared gradient.
182     - t: Iteration number.
183     """
184     if config is None: config = {}
185     config.setdefault('learning_rate', 1e-3)
186     config.setdefault('beta1', 0.9)
187     config.setdefault('beta2', 0.999)
188     config.setdefault('epsilon', 1e-8)
189     config.setdefault('v', np.zeros_like(w))
190     config.setdefault('a', np.zeros_like(w))
191     config.setdefault('t', 0)
192
193     next_w = None
194
195     # ===== #
196     # YOUR CODE HERE:
197     # Implement Adam. Store the next value of w as next_w. You need
198     # to also store in config['a'] the moving average of the second
199     # moment gradients, and in config['v'] the moving average of the
200     # first moments. Finally, store in config['t'] the increasing time.
201     # ===== #
202
203     t = config['t']
204     v = config['v']
205     a = config['a']
206     eps = config['learning_rate']
207     nu = config['epsilon']
208     beta1 = config['beta1']
209     beta2 = config['beta2']
210
211     t += 1
212     v = beta1 * v + (1 - beta1) * dw
213     a = beta2 * a + (1 - beta2) * dw * dw
214     v_u = v / (1 - beta1**t)
215     a_u = a / (1 - beta2**t)
216     w -= (eps * v_u) / (np.sqrt(a_u) + nu)
217
218     config['t'] = t
219     config['v'] = v
220     config['a'] = a
221     next_w = w
222
223     # ===== #
224     # END YOUR CODE HERE
225     # ===== #
226
227     return next_w, config
228
229

```

230
231
232
233

```

1 import numpy as np
2
3 from nndl.layers import *
4 from nndl.conv_layers import *
5 from utils.fast_layers import *
6 from nndl.layer_utils import *
7 from nndl.conv_layer_utils import *
8
9 import pdb
10
11 class ThreeLayerConvNet(object):
12     """
13     A three-layer convolutional network with the following architecture:
14
15     conv - relu - 2x2 max pool - affine - relu - affine - softmax
16
17     The network operates on minibatches of data that have shape (N, C, H, W)
18     consisting of N images, each with height H and width W and with C input
19     channels.
20     """
21
22     def __init__(self, input_dim=(3, 32, 32), num_filters=32, filter_size=7,
23                 hidden_dim=100, num_classes=10, weight_scale=1e-3, reg=0.0,
24                 dtype=np.float32, use_batchnorm=False):
25         """
26         Initialize a new network.
27
28         Inputs:
29         - input_dim: Tuple (C, H, W) giving size of input data
30         - num_filters: Number of filters to use in the convolutional layer
31         - filter_size: Size of filters to use in the convolutional layer
32         - hidden_dim: Number of units to use in the fully-connected hidden layer
33         - num_classes: Number of scores to produce from the final affine layer.
34         - weight_scale: Scalar giving standard deviation for random
35           initialization
36           of weights.
37         - reg: Scalar giving L2 regularization strength
38         - dtype: numpy datatype to use for computation.
39         """
40         self.use_batchnorm = use_batchnorm
41         self.params = {}
42         self.reg = reg
43         self.dtype = dtype
44
45         # ===== #
46         # YOUR CODE HERE:
47         # Initialize the weights and biases of a three layer CNN. To
48         initialize:
49         # - the biases should be initialized to zeros.
50         # - the weights should be initialized to a matrix with entries
51         #   drawn from a Gaussian distribution with zero mean and
52         #   standard deviation given by weight_scale.
53         # ===== #
54
55         # 1st Layer
56         C, H, W = input_dim
57
58         stride = 1

```

```

58     pad = (filter_size - 1) / 2
59
60     out_conv_H = int(1 + (H + 2 * pad - filter_size) / stride)
61     out_conv_W = int(1 + (W + 2 * pad - filter_size) / stride)
62
63     self.params['W1'] = np.random.normal(0, weight_scale, [num_filters, C,
filter_size, filter_size])
64     self.params['b1'] = np.zeros([num_filters])
65
66     # 2nd Layer
67     pool_stride = 2
68     pool_filter_size = 2
69
70     out_pool_H = int(1 + (out_conv_H - pool_filter_size) / pool_stride)
71     out_pool_W = int(1 + (out_conv_W - pool_filter_size) / pool_stride)
72
73     self.params['W2'] = np.random.normal(0, weight_scale,
[num_filters*out_pool_H*out_pool_W, hidden_dim])
74     self.params['b2'] = np.zeros([hidden_dim])
75
76     # 3rd Layer
77     self.params['W3'] = np.random.normal(0, weight_scale, [hidden_dim,
num_classes])
78     self.params['b3'] = np.zeros([num_classes])
79
80     # ===== #
81     # END YOUR CODE HERE
82     # ===== #
83
84     for k, v in self.params.items():
85         self.params[k] = v.astype(dtype)
86
87
88 def loss(self, X, y=None):
89     """
90     Evaluate loss and gradient for the three-layer convolutional network.
91
92     Input / output: Same API as TwoLayerNet in fc_net.py.
93     """
94     W1, b1 = self.params['W1'], self.params['b1']
95     W2, b2 = self.params['W2'], self.params['b2']
96     W3, b3 = self.params['W3'], self.params['b3']
97
98     # pass conv_param to the forward pass for the convolutional layer
99     filter_size = W1.shape[2]
100    conv_param = {'stride': 1, 'pad': (filter_size - 1) / 2}
101
102    # pass pool_param to the forward pass for the max-pooling layer
103    pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}
104
105    scores = None
106
107    # ===== #
108    # YOUR CODE HERE:
109    # Implement the forward pass of the three layer CNN. Store the output
110    # scores as the variable "scores".
111    # ===== #
112
113    conv_outs, conv_caches = conv_relu_pool_forward(X, W1, b1, conv_param,
pool_param)

```



```

114     fc_outs, fc_caches = affine_relu_forward(conv_outs, W2, b2)
115     scores, caches = affine_forward(fc_outs, W3, b3)
116
117     # ===== #
118     # END YOUR CODE HERE
119     # ===== #
120
121     if y is None:
122         return scores
123
124     loss, grads = 0, {}
125     # ===== #
126     # YOUR CODE HERE:
127     # Implement the backward pass of the three layer CNN. Store the grads
128     # in the grads dictionary, exactly as before (i.e., the gradient of
129     # self.params[k] will be grads[k]). Store the loss as "loss", and
130     # don't forget to add regularization on ALL weight matrices.
131     # ===== #
132
133     loss, dx = softmax_loss(scores, y)
134     reg_loss_sum = 0
135     reg_loss_sum += (np.linalg.norm(W1)**2 + np.linalg.norm(W2)**2 +
136 np.linalg.norm(W3)**2)
137
138     loss += 0.5 * self.reg * reg_loss_sum
139
140     dx3, dW3, db3 = affine_backward(dx, caches)
141     dx2, dW2, db2 = affine_relu_backward(dx3, fc_caches)
142     dx1, dW1, db1 = conv_relu_pool_backward(dx2, conv_caches)
143
144     grads['W1'] = dW1 + self.reg * W1
145     grads['b1'] = db1
146     grads['W2'] = dW2 + self.reg * W2
147     grads['b2'] = db2
148     grads['W3'] = dW3 + self.reg * W3
149     grads['b3'] = db3
150
151     # ===== #
152     # END YOUR CODE HERE
153     # ===== #
154
155     return loss, grads
156
157 pass
158

```