# 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 [2]:
```python
## 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))))
```

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

# 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 [17]:
```python
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:  2.2121476575931688e-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 [31]:
```python
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,
dw_num = eval_numerical_gradient_array(lambda w: conv_forward_naive(x, w, b, conv_param)[0], w,
db_num = eval_numerical_gradient_array(lambda b: conv_forward_naive(x, w, b, conv_param)[0], b,

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))
```

```
(4, 3, 5, 5) (4, 2, 5, 5)
Testing conv_backward_naive function
dx error:  2.8496324342528193e-08
dw error:  1.145905092285234e-09
db error:  3.796437261992522e-12
```

## 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 [37]:
```python
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 [41]:
```python
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, d

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:   3.2756245539905995e-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 [42]:
```python
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: 14.657998s
Fast: 0.479885s
Speedup: 30.544792x
Difference:   4.3381416333149016e-12

Testing conv_backward_fast:
Naive: 37.213266s
Fast: 0.009992s
Speedup: 3724.438055x
dx difference:   1.9797083570797586e-11
dw difference:   8.705056450802472e-13
db difference:   0.0
```

In [43]:
```python
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: 0.388547s
fast: 0.003029s
speedup: 128.271232x
difference:  0.0

Testing pool_backward_fast:
Naive: 0.691170s
speedup: 63.002065x
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 [46]:
```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, po
dw_num = eval_numerical_gradient_array(lambda w: conv_relu_pool_forward(x, w, b, conv_param, po
db_num = eval_numerical_gradient_array(lambda b: conv_relu_pool_forward(x, w, b, conv_param, po

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:   4.5124399231583796e-08
dw error:   6.170788280844852e-10
db error:   2.2661151213160115e-11
```

In [47]:
```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)
```

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

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:  1.8323996630127003e-09
dw error:  1.0764646067955963e-09
db error:  4.147675486357139e-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*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.

```
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 [15]:
```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:  [ 9.24079377 11.41428531 10.43869742]
  Stds:  [3.42195716 3.50733199 3.75066663]
After spatial batch normalization:
  Shape:  (2, 3, 4, 5)
  Means:  [-4.44089210e-16  8.27116153e-16  5.66213743e-16]
  Stds:  [0.99999957 0.99999959 0.99999964]
After spatial batch normalization (nontrivial gamma, beta):
  Shape:  (2, 3, 4, 5)
  Means:  [6. 7. 8.]
  Stds:  [2.99999872 3.99999837 4.99999822]
```

# 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 [18]:
```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)
```

```python
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:  5.68794263178987e-09
dgamma error:  3.842018013168464e-12
dbeta error:  3.275106490366651e-12
```

In [ ]:

# 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]:
```python
# 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 [4]:
```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_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 [8]:
```python
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, h=1e-6
    e = rel_error(param_grad_num, grads[param_name])
    print('{} max relative error: {}'.format(param_name, rel_error(param_grad_num, grads[param_
```

```
W1 max relative error: 0.011582167206514055
W2 max relative error: 0.014215630010808635
W3 max relative error: 7.143934702505574e-05
b1 max relative error: 1.8072548730142395e-05
b2 max relative error: 1.5045429753527004e-06
b3 max relative error: 1.0474466450818472e-09
```

## Overfit small dataset

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

In [9]:
```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: 2.382580
```

```
(Epoch 0 / 10) train acc: 0.230000; val_acc: 0.127000
(Iteration 2 / 20) loss: 3.858268
(Epoch 1 / 10) train acc: 0.200000; val_acc: 0.114000
(Iteration 3 / 20) loss: 3.676738
(Iteration 4 / 20) loss: 2.294921
(Epoch 2 / 10) train acc: 0.250000; val_acc: 0.108000
(Iteration 5 / 20) loss: 2.364957
(Iteration 6 / 20) loss: 2.177080
(Epoch 3 / 10) train acc: 0.190000; val_acc: 0.128000
(Iteration 7 / 20) loss: 2.091676
(Iteration 8 / 20) loss: 1.936924
(Epoch 4 / 10) train acc: 0.350000; val_acc: 0.152000
(Iteration 9 / 20) loss: 1.948134
(Iteration 10 / 20) loss: 1.703696
(Epoch 5 / 10) train acc: 0.530000; val_acc: 0.175000
(Iteration 11 / 20) loss: 1.607247
(Iteration 12 / 20) loss: 1.332345
(Epoch 6 / 10) train acc: 0.600000; val_acc: 0.182000
(Iteration 13 / 20) loss: 1.358965
(Iteration 14 / 20) loss: 1.301250
(Epoch 7 / 10) train acc: 0.710000; val_acc: 0.217000
(Iteration 15 / 20) loss: 1.181836
(Iteration 16 / 20) loss: 1.537372
(Epoch 8 / 10) train acc: 0.730000; val_acc: 0.214000
(Iteration 17 / 20) loss: 0.771859
(Iteration 18 / 20) loss: 0.951309
(Epoch 9 / 10) train acc: 0.730000; val_acc: 0.198000
(Iteration 19 / 20) loss: 0.712975
(Iteration 20 / 20) loss: 0.561128
(Epoch 10 / 10) train acc: 0.820000; val_acc: 0.209000
```

In [10]:
```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 [11]:
```python
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()
```
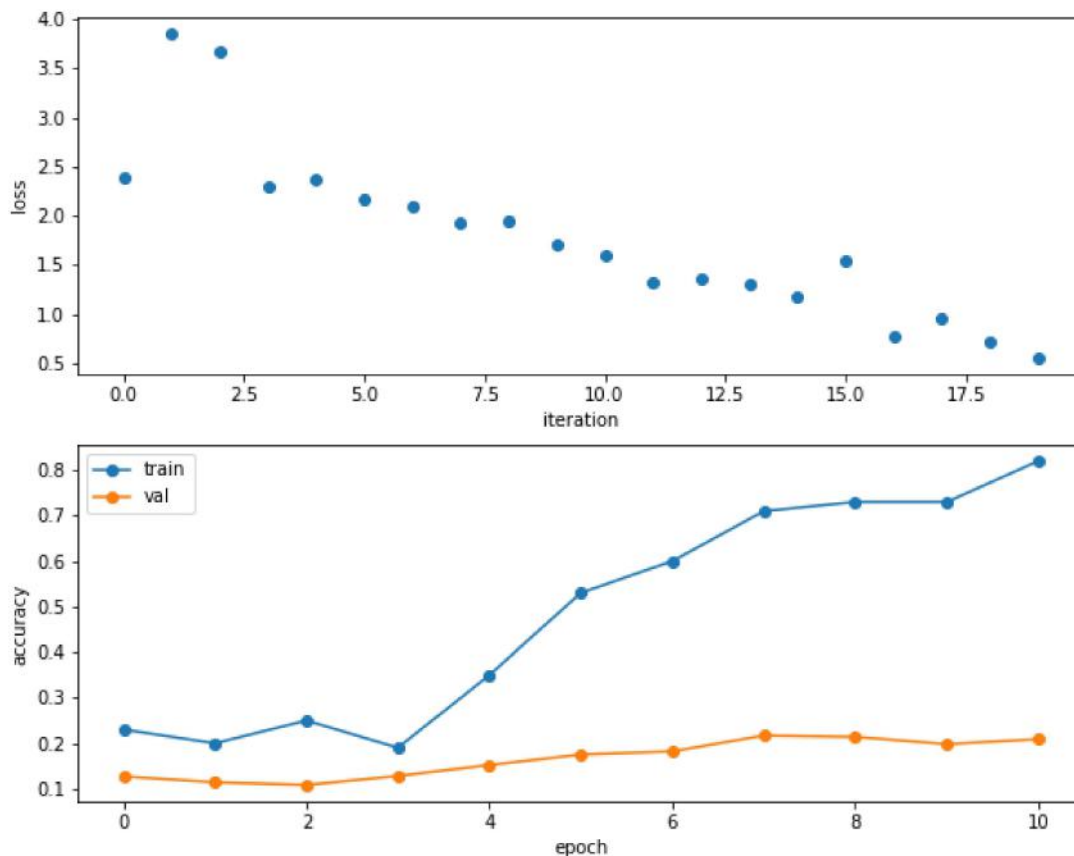
```
(Iteration 1 / 980) loss: 2.304531
(Epoch 0 / 1) train acc: 0.096000; val_acc: 0.098000
(Iteration 21 / 980) loss: 2.271377
(Iteration 41 / 980) loss: 2.002693
(Iteration 61 / 980) loss: 1.996322
(Iteration 81 / 980) loss: 1.863851
(Iteration 101 / 980) loss: 2.064730
(Iteration 121 / 980) loss: 1.992802
(Iteration 141 / 980) loss: 1.452346
(Iteration 161 / 980) loss: 1.869969
(Iteration 181 / 980) loss: 1.793654
(Iteration 201 / 980) loss: 2.008249
(Iteration 221 / 980) loss: 1.715329
(Iteration 241 / 980) loss: 1.515369
(Iteration 261 / 980) loss: 1.675155
(Iteration 281 / 980) loss: 1.806337
(Iteration 301 / 980) loss: 1.473572
(Iteration 321 / 980) loss: 1.447883
(Iteration 341 / 980) loss: 1.981338
(Iteration 361 / 980) loss: 1.718061
(Iteration 381 / 980) loss: 1.418680
(Iteration 401 / 980) loss: 1.740535
(Iteration 421 / 980) loss: 1.706089
```

```
(Iteration 441 / 980) loss: 1.512991
(Iteration 461 / 980) loss: 1.668131
(Iteration 481 / 980) loss: 1.629799
(Iteration 501 / 980) loss: 1.625585
(Iteration 521 / 980) loss: 1.501371
(Iteration 541 / 980) loss: 1.379852
(Iteration 561 / 980) loss: 1.410839
(Iteration 581 / 980) loss: 1.368521
(Iteration 601 / 980) loss: 1.558027
(Iteration 621 / 980) loss: 1.799553
(Iteration 641 / 980) loss: 1.564854
(Iteration 661 / 980) loss: 1.369897
(Iteration 681 / 980) loss: 1.381831
(Iteration 701 / 980) loss: 1.705887
(Iteration 721 / 980) loss: 1.601044
(Iteration 741 / 980) loss: 1.451873
(Iteration 761 / 980) loss: 1.671530
(Iteration 781 / 980) loss: 1.789960
(Iteration 801 / 980) loss: 1.541475
(Iteration 821 / 980) loss: 1.907693
(Iteration 841 / 980) loss: 1.492666
(Iteration 861 / 980) loss: 1.340623
(Iteration 881 / 980) loss: 1.759223
(Iteration 901 / 980) loss: 1.692885
(Iteration 921 / 980) loss: 1.639575
(Iteration 941 / 980) loss: 1.656631
(Iteration 961 / 980) loss: 1.858810
(Epoch 1 / 1) train acc: 0.462000; val_acc: 0.433000
```

# 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 aafter 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 [13]:    # =============================================================== #
            # YOUR CODE HERE:
            #   Implement a CNN to achieve greater than 65% validation accuracy
            #   on CIFAR-10.
```

```
# ================================================================ #

# tuning hyperparameters achieves >65% validation accuracy
model = ThreeLayerConvNet(weight_scale=0.001, hidden_dim=500, reg=0.001)

solver = Solver(model, data,
                num_epochs=10, batch_size=200,
                update_rule='adam',
                optim_config={
                  'learning_rate': 5e-4,
                },
                lr_decay=0.95,
                verbose=True, print_every=100)
solver.train()

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

```
(Iteration 1 / 2450) loss: 2.304694
(Epoch 0 / 10) train acc: 0.085000; val_acc: 0.079000
(Iteration 101 / 2450) loss: 1.515280
(Iteration 201 / 2450) loss: 1.260570
(Epoch 1 / 10) train acc: 0.581000; val_acc: 0.555000
(Iteration 301 / 2450) loss: 1.188566
(Iteration 401 / 2450) loss: 1.167816
(Epoch 2 / 10) train acc: 0.644000; val_acc: 0.604000
(Iteration 501 / 2450) loss: 0.943109
(Iteration 601 / 2450) loss: 0.989317
(Iteration 701 / 2450) loss: 0.836393
(Epoch 3 / 10) train acc: 0.685000; val_acc: 0.614000
(Iteration 801 / 2450) loss: 0.987970
(Iteration 901 / 2450) loss: 1.054005
(Epoch 4 / 10) train acc: 0.727000; val_acc: 0.636000
(Iteration 1001 / 2450) loss: 0.963190
(Iteration 1101 / 2450) loss: 0.841740
(Iteration 1201 / 2450) loss: 0.807908
(Epoch 5 / 10) train acc: 0.765000; val_acc: 0.657000
(Iteration 1301 / 2450) loss: 0.740925
(Iteration 1401 / 2450) loss: 0.644606
(Epoch 6 / 10) train acc: 0.778000; val_acc: 0.661000
(Iteration 1501 / 2450) loss: 0.623828
(Iteration 1601 / 2450) loss: 0.649584
(Iteration 1701 / 2450) loss: 0.727654
(Epoch 7 / 10) train acc: 0.818000; val_acc: 0.646000
(Iteration 1801 / 2450) loss: 0.611627
(Iteration 1901 / 2450) loss: 0.639321
(Epoch 8 / 10) train acc: 0.829000; val_acc: 0.665000
(Iteration 2001 / 2450) loss: 0.581038
(Iteration 2101 / 2450) loss: 0.390481
(Iteration 2201 / 2450) loss: 0.517170
(Epoch 9 / 10) train acc: 0.878000; val_acc: 0.654000
(Iteration 2301 / 2450) loss: 0.501238
(Iteration 2401 / 2450) loss: 0.476331
(Epoch 10 / 10) train acc: 0.905000; val_acc: 0.671000
```

In [ ]:

```python
import numpy as np

from nndl.layers import *
from nndl.conv_layers import *
from utils.fast_layers import *
from nndl.layer_utils import *
from nndl.conv_layer_utils import *

import pdb

class ThreeLayerConvNet(object):
  """
  A three-layer convolutional network with the following architecture:

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

  The network operates on minibatches of data that have shape (N, C, H, W)
  consisting of N images, each with height H and width W and with C input
  channels.
  """

  def __init__(self, input_dim=(3, 32, 32), num_filters=32, filter_size=7,
               hidden_dim=100, num_classes=10, weight_scale=1e-3, reg=0.0,
               dtype=np.float32, use_batchnorm=False):
    """
    Initialize a new network.

    Inputs:
    - input_dim: Tuple (C, H, W) giving size of input data
    - num_filters: Number of filters to use in the convolutional layer
    - filter_size: Size of filters to use in the convolutional layer
    - hidden_dim: Number of units to use in the fully-connected hidden layer
    - num_classes: Number of scores to produce from the final affine layer.
    - weight_scale: Scalar giving standard deviation for random initialization
      of weights.
    - reg: Scalar giving L2 regularization strength
    - dtype: numpy datatype to use for computation.
    """
    self.use_batchnorm = use_batchnorm
    self.params = {}
    self.reg = reg
    self.dtype = dtype


    # ================================================================ #
    # YOUR CODE HERE:
    #   Initialize the weights and biases of a three layer CNN. To initialize:
    #     - the biases should be initialized to zeros.
    #     - the weights should be initialized to a matrix with entries
    #         drawn from a Gaussian distribution with zero mean and
    #         standard deviation given by weight_scale.
    # ================================================================ #

    self.params['W1'] = np.random.normal(loc=0.0, scale=weight_scale,
      size=(num_filters,input_dim[0],filter_size,filter_size))
    self.params['b1'] = np.zeros((num_filters))
    self.params['W2'] = np.random.normal(loc=0.0, scale=weight_scale,
      size=(int(num_filters*input_dim[1]*input_dim[2]/4),hidden_dim))
    self.params['b2'] = np.zeros((hidden_dim))
    self.params['W3'] = np.random.normal(loc=0.0, scale=weight_scale,
      size=(hidden_dim,num_classes))
    self.params['b3'] = np.zeros((num_classes))
    # ================================================================ #
    # END YOUR CODE HERE
    # ================================================================ #

    for k, v in self.params.items():
      self.params[k] = v.astype(dtype)


  def loss(self, X, y=None):
    """
```

```
          Evaluate loss and gradient for the three-layer convolutional network.

          Input / output: Same API as TwoLayerNet in fc_net.py.
          """
          W1, b1 = self.params['W1'], self.params['b1']
          W2, b2 = self.params['W2'], self.params['b2']
          W3, b3 = self.params['W3'], self.params['b3']

          # pass conv_param to the forward pass for the convolutional layer
          filter_size = W1.shape[2]
          conv_param = {'stride': 1, 'pad': (filter_size - 1) / 2}

          # pass pool_param to the forward pass for the max-pooling layer
          pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}

          scores = None

          # ================================================================ #
          # YOUR CODE HERE:
          #    Implement the forward pass of the three layer CNN.  Store the output
          #    scores as the variable "scores".
          # ================================================================ #
          a,c1 = conv_relu_pool_forward(X, W1, b1, conv_param, pool_param)
          a,c2 = affine_relu_forward(a, W2, b2)
          scores,c3 = affine_forward(a,W3,b3)
          # ================================================================ #
          # END YOUR CODE HERE
          # ================================================================ #

          if y is None:
            return scores

          loss, grads = 0, {}
          # ================================================================ #
          # YOUR CODE HERE:
          #    Implement the backward pass of the three layer CNN.  Store the grads
          #    in the grads dictionary, exactly as before (i.e., the gradient of
          #    self.params[k] will be grads[k]).  Store the loss as "loss", and
          #    don't forget to add regularization on ALL weight matrices.
          # ================================================================ #

          loss, ds = softmax_loss(scores,y)
          da,grads['W3'],grads['b3'] = affine_backward(ds,c3)
          da,grads['W2'],grads['b2'] = affine_relu_backward(da,c2)
          da,grads['W1'],grads['b1'] = conv_relu_pool_backward(da,c1)
          for w in ['W1','W2','W3']:
              grads[w]+= self.reg*self.params[w]
              loss += self.reg*0.5*np.sum(np.square(self.params[w]))
          # ================================================================ #
          # END YOUR CODE HERE
          # ================================================================ #

          return loss, grads


      pass
```

```python
import numpy as np
from nndl.layers import *
import pdb


def conv_forward_naive(x, w, b, conv_param):
  """
  A naive implementation of the forward pass for a convolutional layer.

  The input consists of N data points, each with C channels, height H and width
  W. We convolve each input with F different filters, where each filter spans
  all C channels and has height HH and width HH.

  Input:
  - x: Input data of shape (N, C, H, W)
  - w: Filter weights of shape (F, C, HH, WW)
  - b: Biases, of shape (F,)
  - conv_param: A dictionary with the following keys:
    - 'stride': The number of pixels between adjacent receptive fields in the
      horizontal and vertical directions.
    - 'pad': The number of pixels that will be used to zero-pad the input.

  Returns a tuple of:
  - out: Output data, of shape (N, F, H', W') where H' and W' are given by
    H' = 1 + (H + 2 * pad - HH) / stride
    W' = 1 + (W + 2 * pad - WW) / stride
  - cache: (x, w, b, conv_param)
  """
  out = None
  pad = conv_param['pad']
  stride = conv_param['stride']

  # ================================================================ #
  # YOUR CODE HERE:
  #   Implement the forward pass of a convolutional neural network.
  #   Store the output as 'out'.
  #   Hint: to pad the array, you can use the function np.pad.
  # ================================================================ #
  N,C,H,W = x.shape
  F, C, HH, WW = w.shape
  npad = ((0, 0), (0, 0), (pad, pad), (pad, pad))
  xpad = np.pad(x,npad)

  H_ = int(1 + (H + 2 * pad - HH) / stride)
  W_ = int(1 + (W + 2 * pad - WW) / stride)
  out = np.zeros((N,F,H_,W_))
  for i in range(N):
    for j in range(F):
        for h in range(H_):
            for width in range(W_):
                cur=0
                for i2 in range(C):
                    for j2 in range(HH):
                        for k2 in range(WW):
                            cur+=xpad[i,i2,h*stride+j2,width*stride+k2]*w[j,i2,j2,k2]
                out[i,j,h,width] = cur + b[j]
  # ================================================================ #
  # END YOUR CODE HERE
  # ================================================================ #

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


def conv_backward_naive(dout, cache):
  """
  A naive implementation of the backward pass for a convolutional layer.

  Inputs:
  - dout: Upstream derivatives.
  - cache: A tuple of (x, w, b, conv_param) as in conv_forward_naive
```

```python
 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        db=np.sum(dout,axis=(0,2,3))
 93        dxpad=np.zeros(xpad.shape)
 94        dw=np.zeros(w.shape)
 95
 96
 97        N,C,H,W = x.shape
 98        F, C, HH, WW = w.shape
 99        H_  = int(1 + (H + 2 * pad - HH) / stride)
100        W_  = int(1 + (W + 2 * pad - WW) / stride)
101        for i in range(N):
102          for j in range(F):
103              for h in range(H_):
104                  for width in range(W_):
105                      cur=0
106                      for i2 in range(C):
107                          for j2 in range(HH):
108                              for k2 in range(WW):
109                                  #cur+=xpad[i,i2,h*stride+j2,width*stride+k2]*w[j,i2,j2,k2]
110                                  dxpad[i,i2,h*stride+j2,width*stride+k2] +=
                                      w[j,i2,j2,k2]*dout[i,j,h,width]
111                                  dw[j,i2,j2,k2] +=
                                      xpad[i,i2,h*stride+j2,width*stride+k2]*dout[i,j,h,width]
112                      #out[i,j,h,width] = cur + b[j]
113        dx = dxpad[:,:,pad:-pad, pad:-pad]
114        # ================================================================ #
115        # END YOUR CODE HERE
116        # ================================================================ #
117
118        return dx, dw, db
119
120
121    def max_pool_forward_naive(x, pool_param):
122        """
123        A naive implementation of the forward pass for a max pooling layer.
124
125        Inputs:
126        - x: Input data, of shape (N, C, H, W)
127        - pool_param: dictionary with the following keys:
128          - 'pool_height': The height of each pooling region
129          - 'pool_width': The width of each pooling region
130          - 'stride': The distance between adjacent pooling regions
131
132        Returns a tuple of:
133        - out: Output data
134        - cache: (x, pool_param)
135        """
136        out = None
137
138        # ================================================================ #
139        # YOUR CODE HERE:
140        #    Implement the max pooling forward pass.
141        # ================================================================ #
142        N, C, H, W = x.shape
```

```
143        ph = pool_param['pool_height']
144        pw = pool_param['pool_width']
145        s = pool_param['stride']
146        H_ = int(1 + (H-ph)/s)
147        W_ = int(1 + (W-pw)/s)
148        out = np.zeros((N,C,H_,W_))
149        for i in range(N):
150          for j in range(C):
151              for h in range(H_):
152                  for w in range(W_):
153                      out[i,j,h,w]=np.max(x[i,j,h*s:h*s+ph,w*s:w*s+pw])
154        # ================================================================ #
155        # END YOUR CODE HERE
156        # ================================================================ #
157        cache = (x, pool_param)
158        return out, cache
159
160    def max_pool_backward_naive(dout, cache):
161        """
162        A naive implementation of the backward pass for a max pooling layer.
163
164        Inputs:
165        - dout: Upstream derivatives
166        - cache: A tuple of (x, pool_param) as in the forward pass.
167
168        Returns:
169        - dx: Gradient with respect to x
170        """
171        dx = None
172        x, pool_param = cache
173        pool_height, pool_width, stride = pool_param['pool_height'],
               pool_param['pool_width'], pool_param['stride']
174
175        # ================================================================ #
176        # YOUR CODE HERE:
177        #    Implement the max pooling backward pass.
178        # ================================================================ #
179        dx = np.zeros(x.shape)
180        N, C, H, W = x.shape
181        ph = pool_param['pool_height']
182        pw = pool_param['pool_width']
183        s = pool_param['stride']
184        H_ = int(1 + (H-ph)/s)
185        W_ = int(1 + (W-pw)/s)
186        for i in range(N):
187          for j in range(C):
188              for h in range(H_):
189                  for w in range(W_):
190                      cur=np.max(x[i,j,h*s:h*s+ph,w*s:w*s+pw])
191                      for i2 in range(ph):
192                          for j2 in range(pw):
193                              if x[i,j,h*s+i2,w*s+j2]==cur:
194                                  dx[i,j,h*s+i2,w*s+j2]=dout[i,j,h,w]
195
196        # ================================================================ #
197        # END YOUR CODE HERE
198        # ================================================================ #
199
200        return dx
201
202    def spatial_batchnorm_forward(x, gamma, beta, bn_param):
203        """
204        Computes the forward pass for spatial batch normalization.
205
206        Inputs:
207        - x: Input data of shape (N, C, H, W)
208        - gamma: Scale parameter, of shape (C,)
209        - beta: Shift parameter, of shape (C,)
210        - bn_param: Dictionary with the following keys:
211          - mode: 'train' or 'test'; required
212          - eps: Constant for numeric stability
213          - momentum: Constant for running mean / variance. momentum=0 means that
```

```python
214              old information is discarded completely at every time step, while
215              momentum=1 means that new information is never incorporated. The
216              default of momentum=0.9 should work well in most situations.
217          - running_mean: Array of shape (D,) giving running mean of features
218          - running_var Array of shape (D,) giving running variance of features
219
220      Returns a tuple of:
221      - out: Output data, of shape (N, C, H, W)
222      - cache: Values needed for the backward pass
223      """
224      out, cache = None, None
225
226      # ================================================================ #
227      # YOUR CODE HERE:
228      #    Implement the spatial batchnorm forward pass.
229      #
230      #    You may find it useful to use the batchnorm forward pass you
231      #    implemented in HW #4.
232      # ================================================================ #
233      N, C, H, W = x.shape
234      xb = np.transpose(x, (0,2,3,1))
235      xb = xb.reshape((-1,C))
236      out,cache = batchnorm_forward(xb, gamma, beta, bn_param)
237      out = out.reshape((N,H,W,C))
238      out = np.transpose(out, (0,3,1,2))
239      # ================================================================ #
240      # END YOUR CODE HERE
241      # ================================================================ #
242
243      return out, cache
244
245
246  def spatial_batchnorm_backward(dout, cache):
247      """
248      Computes the backward pass for spatial batch normalization.
249
250      Inputs:
251      - dout: Upstream derivatives, of shape (N, C, H, W)
252      - cache: Values from the forward pass
253
254      Returns a tuple of:
255      - dx: Gradient with respect to inputs, of shape (N, C, H, W)
256      - dgamma: Gradient with respect to scale parameter, of shape (C,)
257      - dbeta: Gradient with respect to shift parameter, of shape (C,)
258      """
259      dx, dgamma, dbeta = None, None, None
260
261      # ================================================================ #
262      # YOUR CODE HERE:
263      #    Implement the spatial batchnorm backward pass.
264      #
265      #    You may find it useful to use the batchnorm forward pass you
266      #    implemented in HW #4.
267      # ================================================================ #
268      N, C, H, W = dout.shape
269      dout = np.transpose(dout, (0,2,3,1))
270      dout = dout.reshape(-1,C)
271      #x,x_n,mean,var,eps,gamma = cache
272      #x = x.reshape((N,H,W,C))
273      #x = np.transpose(x, (0,3,1,2))
274      #x_n = x_n.reshape((N,H,W,C))
275      #x_n = np.transpose(x_n, (0,3,1,2))
276      #cache = x,x_n,mean,var,eps,gamma
277      dx, dgamma, dbeta = batchnorm_backward(dout,cache)
278      dx = dx.reshape((N,H,W,C))
279      dx = np.transpose(dx, (0,3,1,2))
280      # ================================================================ #
281      # END YOUR CODE HERE
282      # ================================================================ #
283
284      return dx, dgamma, dbeta
```

```python
from nndl.layers import *
from utils.fast_layers import *


def conv_relu_forward(x, w, b, conv_param):
    """
    A convenience layer that performs a convolution followed by a ReLU.

    Inputs:
    - x: Input to the convolutional layer
    - w, b, conv_param: Weights and parameters for the convolutional layer

    Returns a tuple of:
    - out: Output from the ReLU
    - cache: Object to give to the backward pass
    """
    a, conv_cache = conv_forward_fast(x, w, b, conv_param)
    out, relu_cache = relu_forward(a)
    cache = (conv_cache, relu_cache)
    return out, cache


def conv_relu_backward(dout, cache):
    """
    Backward pass for the conv-relu convenience layer.
    """
    conv_cache, relu_cache = cache
    da = relu_backward(dout, relu_cache)
    dx, dw, db = conv_backward_fast(da, conv_cache)
    return dx, dw, db


def conv_relu_pool_forward(x, w, b, conv_param, pool_param):
    """
    Convenience layer that performs a convolution, a ReLU, and a pool.

    Inputs:
    - x: Input to the convolutional layer
    - w, b, conv_param: Weights and parameters for the convolutional layer
    - pool_param: Parameters for the pooling layer

    Returns a tuple of:
    - out: Output from the pooling layer
    - cache: Object to give to the backward pass
    """
    a, conv_cache = conv_forward_fast(x, w, b, conv_param)
    s, relu_cache = relu_forward(a)
    out, pool_cache = max_pool_forward_fast(s, pool_param)
    cache = (conv_cache, relu_cache, pool_cache)
    return out, cache


def conv_relu_pool_backward(dout, cache):
    """
    Backward pass for the conv-relu-pool convenience layer
    """
    conv_cache, relu_cache, pool_cache = cache
    ds = max_pool_backward_fast(dout, pool_cache)
    da = relu_backward(ds, relu_cache)
    dx, dw, db = conv_backward_fast(da, conv_cache)
    return dx, dw, db
```

```python
import numpy as np
import pdb

def affine_forward(x, w, b):
    """
    Computes the forward pass for an affine (fully-connected) layer.

    The input x has shape (N, d_1, ..., d_k) and contains a minibatch of N
    examples, where each example x[i] has shape (d_1, ..., d_k). We will
    reshape each input into a vector of dimension D = d_1 * ... * d_k, and
    then transform it to an output vector of dimension M.

    Inputs:
    - x: A numpy array containing input data, of shape (N, d_1, ..., d_k)
    - w: A numpy array of weights, of shape (D, M)
    - b: A numpy array of biases, of shape (M,)

    Returns a tuple of:
    - out: output, of shape (N, M)
    - cache: (x, w, b)
    """

    # ================================================================ #
    # YOUR CODE HERE:
    #   Calculate the output of the forward pass.  Notice the dimensions
    #   of w are D x M, which is the transpose of what we did in earlier
    #   assignments.
    # ================================================================ #

    out = x.reshape(x.shape[0],-1).dot(w)+b
    # ================================================================ #
    # END YOUR CODE HERE
    # ================================================================ #

    cache = (x, w, b)
    return out, cache


def affine_backward(dout, cache):
    """
    Computes the backward pass for an affine layer.

    Inputs:
    - dout: Upstream derivative, of shape (N, M)
    - cache: Tuple of:
      - x: Input data, of shape (N, d_1, ... d_k)
      - w: Weights, of shape (D, M)

    Returns a tuple of:
    - dx: Gradient with respect to x, of shape (N, d1, ..., d_k)
    - dw: Gradient with respect to w, of shape (D, M)
    - db: Gradient with respect to b, of shape (M,)
    """
    x, w, b = cache
    dx, dw, db = None, None, None

    # ================================================================ #
    # YOUR CODE HERE:
    #   Calculate the gradients for the backward pass.
    # ================================================================ #

    dx = dout.dot(w.T).reshape(x.shape)
    dw = x.reshape(x.shape[0],-1).T.dot(dout).reshape(w.shape)
    db = np.sum(dout,axis=0)
    # ================================================================ #
    # END YOUR CODE HERE
    # ================================================================ #

    return dx, dw, db

def relu_forward(x):
    """
```

```python
    Computes the forward pass for a layer of rectified linear units (ReLUs).

    Input:
    - x: Inputs, of any shape

    Returns a tuple of:
    - out: Output, of the same shape as x
    - cache: x
    """
    # ================================================================ #
    # YOUR CODE HERE:
    #    Implement the ReLU forward pass.
    # ================================================================ #
    out = np.maximum(x,0)
    # ================================================================ #
    # END YOUR CODE HERE
    # ================================================================ #

    cache = x
    return out, cache


def relu_backward(dout, cache):
    """
    Computes the backward pass for a layer of rectified linear units (ReLUs).

    Input:
    - dout: Upstream derivatives, of any shape
    - cache: Input x, of same shape as dout

    Returns:
    - dx: Gradient with respect to x
    """
    x = cache

    # ================================================================ #
    # YOUR CODE HERE:
    #    Implement the ReLU backward pass
    # ================================================================ #

    dx = dout
    dx[x<0]=0
    # ================================================================ #
    # END YOUR CODE HERE
    # ================================================================ #

    return dx

def batchnorm_forward(x, gamma, beta, bn_param):
    """
    Forward pass for batch normalization.

    During training the sample mean and (uncorrected) sample variance are
    computed from minibatch statistics and used to normalize the incoming data.
    During training we also keep an exponentially decaying running mean of the mean
    and variance of each feature, and these averages are used to normalize data
    at test-time.

    At each timestep we update the running averages for mean and variance using
    an exponential decay based on the momentum parameter:

    running_mean = momentum * running_mean + (1 - momentum) * sample_mean
    running_var = momentum * running_var + (1 - momentum) * sample_var

    Note that the batch normalization paper suggests a different test-time
    behavior: they compute sample mean and variance for each feature using a
    large number of training images rather than using a running average. For
    this implementation we have chosen to use running averages instead since
    they do not require an additional estimation step; the torch7 implementation
    of batch normalization also uses running averages.

    Input:
```

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

```python
    """
    Backward pass for batch normalization.

    For this implementation, you should write out a computation graph for
    batch normalization on paper and propagate gradients backward through
    intermediate nodes.

    Inputs:
    - dout: Upstream derivatives, of shape (N, D)
    - cache: Variable of intermediates from batchnorm_forward.

    Returns a tuple of:
    - dx: Gradient with respect to inputs x, of shape (N, D)
    - dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
    - dbeta: Gradient with respect to shift parameter beta, of shape (D,)
    """
    dx, dgamma, dbeta = None, None, None

    # ================================================================ #
    # YOUR CODE HERE:
    #   Implement the batchnorm backward pass, calculating dx, dgamma, and dbeta.
    # ================================================================ #

    x,x_n,mean,var,eps,gamma = cache
    dbeta = np.sum(dout,axis=0)
    dgamma = np.sum(dout*x_n,axis=0)
    dxh = dout*gamma
    std = np.sqrt(var+eps)
    du = -np.sum((1/std)*dxh,axis=0)
    dv = -np.sum((1/(2*((std)**3)))*(x-mean)*dxh,axis=0)
    dx = dxh*(1/std) + du/x.shape[0] + (2*(x-mean)/x.shape[0])*dv
    # ================================================================ #
    # END YOUR CODE HERE
    # ================================================================ #

    return dx, dgamma, dbeta

def dropout_forward(x, dropout_param):
    """
    Performs the forward pass for (inverted) dropout.

    Inputs:
    - x: Input data, of any shape
    - dropout_param: A dictionary with the following keys:
      - p: Dropout parameter. We drop each neuron output with probability p.
      - mode: 'test' or 'train'. If the mode is train, then perform dropout;
        if the mode is test, then just return the input.
      - seed: Seed for the random number generator. Passing seed makes this
        function deterministic, which is needed for gradient checking but not in
        real networks.

    Outputs:
    - out: Array of the same shape as x.
    - cache: A tuple (dropout_param, mask). In training mode, mask is the dropout
      mask that was used to multiply the input; in test mode, mask is None.
    """
    p, mode = dropout_param['p'], dropout_param['mode']
    if 'seed' in dropout_param:
      np.random.seed(dropout_param['seed'])

    mask = None
    out = None

    if mode == 'train':
      # ================================================================ #
      # YOUR CODE HERE:
      #   Implement the inverted dropout forward pass during training time.
      #   Store the masked and scaled activations in out, and store the
      #   dropout mask as the variable mask.
      # ================================================================ #

      mask = (np.random.rand(*x.shape) < (1-p) )/(1-p)
```

```python
289        out = x*mask
290        # =============================================================== #
291        # END YOUR CODE HERE
292        # =============================================================== #

294    elif mode == 'test':

296        # =============================================================== #
297        # YOUR CODE HERE:
298        #    Implement the inverted dropout forward pass during test time.
299        # =============================================================== #
300        out=x

302        # =============================================================== #
303        # END YOUR CODE HERE
304        # =============================================================== #

306    cache = (dropout_param, mask)
307    out = out.astype(x.dtype, copy=False)

309    return out, cache

311  def dropout_backward(dout, cache):
312      """
313      Perform the backward pass for (inverted) dropout.

315      Inputs:
316      - dout: Upstream derivatives, of any shape
317      - cache: (dropout_param, mask) from dropout_forward.
318      """
319      dropout_param, mask = cache
320      mode = dropout_param['mode']

322      dx = None
323      if mode == 'train':
324        # =============================================================== #
325        # YOUR CODE HERE:
326        #    Implement the inverted dropout backward pass during training time.
327        # =============================================================== #
328        (dropout_param, mask) = cache
329        dx = dout*mask

331        # =============================================================== #
332        # END YOUR CODE HERE
333        # =============================================================== #
334      elif mode == 'test':
335        # =============================================================== #
336        # YOUR CODE HERE:
337        #    Implement the inverted dropout backward pass during test time.
338        # =============================================================== #
339        dx=dout
340        # =============================================================== #
341        # END YOUR CODE HERE
342        # =============================================================== #
343      return dx

345  def svm_loss(x, y):
346      """
347      Computes the loss and gradient using for multiclass SVM classification.

349      Inputs:
350      - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
351        for the ith input.
352      - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
353        0 <= y[i] < C

355      Returns a tuple of:
356      - loss: Scalar giving the loss
357      - dx: Gradient of the loss with respect to x
358      """
359      N = x.shape[0]
360      correct_class_scores = x[np.arange(N), y]
```

```python
    margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
    margins[np.arange(N), y] = 0
    loss = np.sum(margins) / N
    num_pos = np.sum(margins > 0, axis=1)
    dx = np.zeros_like(x)
    dx[margins > 0] = 1
    dx[np.arange(N), y] -= num_pos
    dx /= N
    return loss, dx


def softmax_loss(x, y):
    """
    Computes the loss and gradient for softmax classification.

    Inputs:
    - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
      for the ith input.
    - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
      0 <= y[i] < C

    Returns a tuple of:
    - loss: Scalar giving the loss
    - dx: Gradient of the loss with respect to x
    """

    probs = np.exp(x - np.max(x, axis=1, keepdims=True))
    probs /= np.sum(probs, axis=1, keepdims=True)
    N = x.shape[0]
    loss = -np.sum(np.log(probs[np.arange(N), y])) / N
    dx = probs.copy()
    dx[np.arange(N), y] -= 1
    dx /= N
    return loss, dx
```

```python
import numpy as np

"""
This file implements various first-order update rules that are commonly used for
training neural networks. Each update rule accepts current weights and the
gradient of the loss with respect to those weights and produces the next set of
weights. Each update rule has the same interface:

def update(w, dw, config=None):

Inputs:
  - w: A numpy array giving the current weights.
  - dw: A numpy array of the same shape as w giving the gradient of the
    loss with respect to w.
  - config: A dictionary containing hyperparameter values such as learning rate,
    momentum, etc. If the update rule requires caching values over many
    iterations, then config will also hold these cached values.

Returns:
  - next_w: The next point after the update.
  - config: The config dictionary to be passed to the next iteration of the
    update rule.

NOTE: For most update rules, the default learning rate will probably not perform
well; however the default values of the other hyperparameters should work well
for a variety of different problems.

For efficiency, update rules may perform in-place updates, mutating w and
setting next_w equal to w.
"""


def sgd(w, dw, config=None):
    """
    Performs vanilla stochastic gradient descent.

    config format:
    - learning_rate: Scalar learning rate.
    """
    if config is None: config = {}
    config.setdefault('learning_rate', 1e-2)

    w -= config['learning_rate'] * dw
    return w, config


def sgd_momentum(w, dw, config=None):
    """
    Performs stochastic gradient descent with momentum.

    config format:
    - learning_rate: Scalar learning rate.
    - momentum: Scalar between 0 and 1 giving the momentum value.
      Setting momentum = 0 reduces to sgd.
    - velocity: A numpy array of the same shape as w and dw used to store a moving
      average of the gradients.
    """
    if config is None: config = {}
    config.setdefault('learning_rate', 1e-2)
    config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
    v = config.get('velocity', np.zeros_like(w))   # gets velocity, else sets it to
    zero.

    # ================================================================ #
    # YOUR CODE HERE:
    #   Implement the momentum update formula.  Return the updated weights
    #   as next_w, and the updated velocity as v.
    # ================================================================ #
    v = config['momentum'] * v - config['learning_rate'] * dw
    next_w = w+v
    # ================================================================ #
    # END YOUR CODE HERE
```

```python
     # ================================================================ #

    config['velocity'] = v

    return next_w, config

  def sgd_nesterov_momentum(w, dw, config=None):
    """
    Performs stochastic gradient descent with Nesterov momentum.

    config format:
    - learning_rate: Scalar learning rate.
    - momentum: Scalar between 0 and 1 giving the momentum value.
      Setting momentum = 0 reduces to sgd.
    - velocity: A numpy array of the same shape as w and dw used to store a moving
      average of the gradients.
    """
    if config is None: config = {}
    config.setdefault('learning_rate', 1e-2)
    config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
    v = config.get('velocity', np.zeros_like(w))   # gets velocity, else sets it to
    zero.

    # ================================================================ #
    # YOUR CODE HERE:
    #    Implement the momentum update formula.  Return the updated weights
    #    as next_w, and the updated velocity as v.
    # ================================================================ #
    v_new = config['momentum'] * v - config['learning_rate'] * dw
    next_w = w+v_new + config['momentum'] *(v_new-v)
    v = v_new

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

    config['velocity'] = v

    return next_w, config

  def rmsprop(w, dw, config=None):
    """
    Uses the RMSProp update rule, which uses a moving average of squared gradient
    values to set adaptive per-parameter learning rates.

    config format:
    - learning_rate: Scalar learning rate.
    - decay_rate: Scalar between 0 and 1 giving the decay rate for the squared
      gradient cache.
    - epsilon: Small scalar used for smoothing to avoid dividing by zero.
    - beta: Moving average of second moments of gradients.
    """
    if config is None: config = {}
    config.setdefault('learning_rate', 1e-2)
    config.setdefault('decay_rate', 0.99)
    config.setdefault('epsilon', 1e-8)
    config.setdefault('a', np.zeros_like(w))

    next_w = None

    # ================================================================ #
    # YOUR CODE HERE:
    #    Implement RMSProp.  Store the next value of w as next_w.  You need
    #    to also store in config['a'] the moving average of the second
    #    moment gradients, so they can be used for future gradients. Concretely,
    #    config['a'] corresponds to "a" in the lecture notes.
    # ================================================================ #
    config['a'] = config['decay_rate']*config['a'] + (1-config['decay_rate'])*dw*dw
    next_w = w - (config['learning_rate']*dw)/(np.sqrt(config['a'])+config['epsilon'])

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

```python
    # ================================================================= #

    return next_w, config


def adam(w, dw, config=None):
    """
    Uses the Adam update rule, which incorporates moving averages of both the
    gradient and its square and a bias correction term.

    config format:
    - learning_rate: Scalar learning rate.
    - beta1: Decay rate for moving average of first moment of gradient.
    - beta2: Decay rate for moving average of second moment of gradient.
    - epsilon: Small scalar used for smoothing to avoid dividing by zero.
    - m: Moving average of gradient.
    - v: Moving average of squared gradient.
    - t: Iteration number.
    """
    if config is None: config = {}
    config.setdefault('learning_rate', 1e-3)
    config.setdefault('beta1', 0.9)
    config.setdefault('beta2', 0.999)
    config.setdefault('epsilon', 1e-8)
    config.setdefault('v', np.zeros_like(w))
    config.setdefault('a', np.zeros_like(w))
    config.setdefault('t', 0)

    next_w = None

    # ================================================================= #
    # YOUR CODE HERE:
    #    Implement Adam.  Store the next value of w as next_w.  You need
    #    to also store in config['a'] the moving average of the second
    #    moment gradients, and in config['v'] the moving average of the
    #    first moments.  Finally, store in config['t'] the increasing time.
    # ================================================================= #
    config['t']+=1
    config['v'] = config['beta1']*config['v'] + (1-config['beta1'])*dw
    config['a'] = config['beta2']*config['a'] + (1-config['beta2'])*dw*dw
    v = config['v']/(1-config['beta1']**config['t'])
    a = config['a']/(1-config['beta2']**config['t'])
    next_w = w -(config['learning_rate']*v)/(np.sqrt(a)+config['epsilon'])
    # ================================================================= #
    # END YOUR CODE HERE
    # ================================================================= #

    return next_w, config
```