

# 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_gradients
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 [4]:

```

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

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

# Your errors should be around 1e-9'
print('Testing conv_backward_naive function')
print('dx error: ', rel_error(dx, dx_num))

```

```
print('dw error: ', rel_error(dw, dw_num))
print('db error: ', rel_error(db, db_num))
```

```
Testing conv_backward_naive function
dx error: 1.9360307413110262e-09
dw error: 4.6167829469755946e-10
db error: 1.3896841421511529e-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 [5]:

```
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 [6]:

```
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_
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: 3.2756375184481477e-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 [7]:

```
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: 4.342280s

Fast: 0.011171s

Speedup: 388.706477x

Difference: 1.8459009680469096e-11

Testing conv\_backward\_fast:

Naive: 6.967768s

Fast: 0.006382s

Speedup: 1091.786349x

dx difference: 1.6658138276324063e-11

dw difference: 2.346279647708712e-13

db difference: 1.4982429997675873e-15

In [8]:

```

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: 0.366097s
```

```
fast: 0.003871s
```

```
speedup: 94.575142x
```

```
difference: 0.0
```

```
Testing pool_backward_fast:
```

```
Naive: 0.934449s
```

```
speedup: 87.143432x
```

```
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 [10]:

```
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,
dw_num = eval_numerical_gradient_array(lambda w: conv_relu_pool_forward(x, w, b,
db_num = eval_numerical_gradient_array(lambda b: conv_relu_pool_forward(x, w, b,

print('Testing conv_relu_pool')
print('dx error: ', rel_error(dx_num, dx))
print('dw error: ', rel_error(dw_num, dw))
print('db error: ', rel_error(db_num, db))
```

```
Testing conv_relu_pool
```

```
dx error: 2.0650507323476647e-08
```

```
dw error: 1.3219786820449993e-09
```

```
db error: 1.3054464532345423e-10
```

In [11]:

```
from nndl.conv_layer_utils import conv_relu_forward, conv_relu_backward

x = np.random.randn(2, 3, 8, 8)
w = np.random.randn(3, 3, 3, 3)
b = np.random.randn(3,)
dout = np.random.randn(2, 3, 8, 8)
```

```
conv_param = {'stride': 1, 'pad': 1}

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

dx_num = eval_numerical_gradient_array(lambda x: conv_relu_forward(x, w, b, conv_param), 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:  2.6271548335129014e-09
dw error:  5.763066442438924e-10
db error:  1.0095419769079119e-10
```

## 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_gradien
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-ipytho
%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: [ 9.94964977 10.14779527 10.72345741]
Stds:  [3.29043416 4.29377971 4.30348    ]

```

After spatial batch normalization:

```

Shape: (2, 3, 4, 5)
Means: [ 1.18134669e-16 -7.91033905e-17  5.55111512e-17]

```

```

      Stds: [0.99999936 0.99999976 0.99999973]
After spatial batch normalization (nontrivial gamma, beta):
      Shape: (2, 3, 4, 5)
      Means: [6. 7. 8.]
      Stds: [2.99999809 3.99999905 4.99999864]

```

## 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:  3.4358596328264993e-09
dgamma error:  7.979001466048922e-10
dbeta error:  1.4458894766286844e-09

```

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]: # 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_g
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-ipytho
%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 `nnl/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 [3]:

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

```
W1 max relative error: 0.00010943822087774027
W2 max relative error: 0.005260871818128235
W3 max relative error: 0.00019657125816118783
b1 max relative error: 1.5143119472011596e-06
```

b2 max relative error: 1.729376044093205e-07  
 b3 max relative error: 1.4320398405898022e-09

## Overfit small dataset

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

In [7]:

```
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)

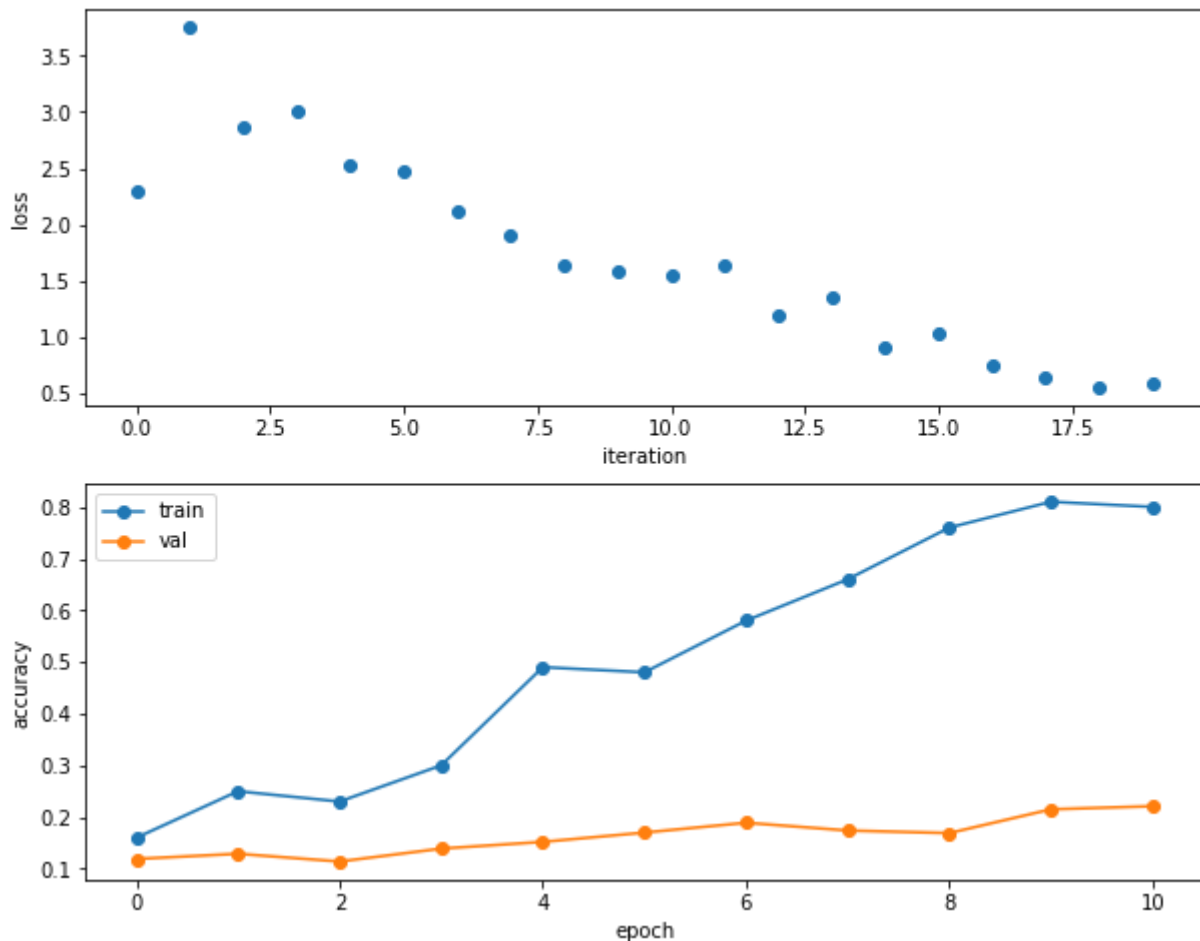
solver = Solver(model, small_data,
                 num_epochs=10, batch_size=50,
                 update_rule='adam',
                 optim_config={
                     'learning_rate': 1e-3,
                 },
                 verbose=True, print_every=1)

solver.train()

(Iteration 1 / 20) loss: 2.288154
(Epoch 0 / 10) train acc: 0.160000; val_acc: 0.119000
(Iteration 2 / 20) loss: 3.754403
(Epoch 1 / 10) train acc: 0.250000; val_acc: 0.129000
(Iteration 3 / 20) loss: 2.859510
(Iteration 4 / 20) loss: 3.013589
(Epoch 2 / 10) train acc: 0.230000; val_acc: 0.114000
(Iteration 5 / 20) loss: 2.531067
(Iteration 6 / 20) loss: 2.474353
(Epoch 3 / 10) train acc: 0.300000; val_acc: 0.139000
(Iteration 7 / 20) loss: 2.124515
(Iteration 8 / 20) loss: 1.911202
(Epoch 4 / 10) train acc: 0.490000; val_acc: 0.152000
(Iteration 9 / 20) loss: 1.644574
(Iteration 10 / 20) loss: 1.584084
(Epoch 5 / 10) train acc: 0.480000; val_acc: 0.170000
(Iteration 11 / 20) loss: 1.544988
(Iteration 12 / 20) loss: 1.644108
(Epoch 6 / 10) train acc: 0.580000; val_acc: 0.189000
(Iteration 13 / 20) loss: 1.198795
(Iteration 14 / 20) loss: 1.349414
(Epoch 7 / 10) train acc: 0.660000; val_acc: 0.174000
(Iteration 15 / 20) loss: 0.914195
(Iteration 16 / 20) loss: 1.029787
(Epoch 8 / 10) train acc: 0.760000; val_acc: 0.169000
(Iteration 17 / 20) loss: 0.755217
(Iteration 18 / 20) loss: 0.647375
(Epoch 9 / 10) train acc: 0.810000; val_acc: 0.215000
(Iteration 19 / 20) loss: 0.561457
(Iteration 20 / 20) loss: 0.593262
(Epoch 10 / 10) train acc: 0.800000; val_acc: 0.221000
```

```
In [8]: 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 [9]: 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: 2.304728
(Epoch 0 / 1) train acc: 0.106000; val_acc: 0.105000
(Iteration 21 / 980) loss: 2.158532
(Iteration 41 / 980) loss: 1.891394
(Iteration 61 / 980) loss: 2.104487
(Iteration 81 / 980) loss: 1.804120
(Iteration 101 / 980) loss: 2.044281
(Iteration 121 / 980) loss: 1.773167
(Iteration 141 / 980) loss: 1.523350
(Iteration 161 / 980) loss: 1.910315
(Iteration 181 / 980) loss: 2.048969
(Iteration 201 / 980) loss: 1.567983
(Iteration 221 / 980) loss: 2.021541
(Iteration 241 / 980) loss: 1.987722
(Iteration 261 / 980) loss: 1.414708
(Iteration 281 / 980) loss: 1.758751
(Iteration 301 / 980) loss: 1.745545
(Iteration 321 / 980) loss: 1.639196
(Iteration 341 / 980) loss: 1.623240
(Iteration 361 / 980) loss: 1.675248
(Iteration 381 / 980) loss: 1.794584
(Iteration 401 / 980) loss: 1.630131
(Iteration 421 / 980) loss: 1.572791
(Iteration 441 / 980) loss: 1.529541
(Iteration 461 / 980) loss: 1.723137
(Iteration 481 / 980) loss: 1.758001
(Iteration 501 / 980) loss: 1.338156
(Iteration 521 / 980) loss: 1.449953
(Iteration 541 / 980) loss: 1.739851
(Iteration 561 / 980) loss: 1.648934
(Iteration 581 / 980) loss: 1.669570
(Iteration 601 / 980) loss: 1.705192
(Iteration 621 / 980) loss: 1.550650
(Iteration 641 / 980) loss: 1.791609
(Iteration 661 / 980) loss: 1.561385
(Iteration 681 / 980) loss: 1.849013
(Iteration 701 / 980) loss: 1.424552
(Iteration 721 / 980) loss: 1.509565
(Iteration 741 / 980) loss: 1.503256
(Iteration 761 / 980) loss: 1.438078
(Iteration 781 / 980) loss: 1.653617
(Iteration 801 / 980) loss: 1.593835
(Iteration 821 / 980) loss: 1.485103
(Iteration 841 / 980) loss: 1.449281
(Iteration 861 / 980) loss: 1.592750
(Iteration 881 / 980) loss: 1.493865
(Iteration 901 / 980) loss: 1.515934
(Iteration 921 / 980) loss: 1.560392
(Iteration 941 / 980) loss: 1.300246
(Iteration 961 / 980) loss: 1.524814
(Epoch 1 / 1) train acc: 0.457000; val_acc: 0.499000
```

## 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 [12]:

```
# ===== #
# 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)

solver = Solver(model, data,
                 num_epochs=7, batch_size=800,
                 update_rule='adam',
                 optim_config={
                     'learning_rate': 1e-3,
                 },
                 verbose=True, print_every=20)

solver.train()

y_val_max = np.argmax(model.loss(data['X_val']), axis=1)
y_test_max = np.argmax(model.loss(data['X_test']), axis=1)
print('Validation set accuracy: {}'.format(np.mean(y_val_max == data['y_val'])))
print('Test set accuracy: {}'.format(np.mean(y_test_max == data['y_test'])))

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



```
(Iteration 1 / 427) loss: 2.304590
(Epoch 0 / 7) train acc: 0.116000; val_acc: 0.119000
(Iteration 21 / 427) loss: 1.857259
(Iteration 41 / 427) loss: 1.625057
(Iteration 61 / 427) loss: 1.480301
(Epoch 1 / 7) train acc: 0.482000; val_acc: 0.489000
(Iteration 81 / 427) loss: 1.381730
(Iteration 101 / 427) loss: 1.324259
(Iteration 121 / 427) loss: 1.245069
(Epoch 2 / 7) train acc: 0.565000; val_acc: 0.568000
(Iteration 141 / 427) loss: 1.258634
(Iteration 161 / 427) loss: 1.305722
(Iteration 181 / 427) loss: 1.084395
(Epoch 3 / 7) train acc: 0.641000; val_acc: 0.594000
(Iteration 201 / 427) loss: 1.118002
(Iteration 221 / 427) loss: 1.105004
(Iteration 241 / 427) loss: 1.123589
(Epoch 4 / 7) train acc: 0.642000; val_acc: 0.594000
(Iteration 261 / 427) loss: 1.004464
(Iteration 281 / 427) loss: 0.981249
(Iteration 301 / 427) loss: 1.017331
(Epoch 5 / 7) train acc: 0.710000; val_acc: 0.619000
(Iteration 321 / 427) loss: 0.951913
(Iteration 341 / 427) loss: 0.853187
(Iteration 361 / 427) loss: 0.877747
(Epoch 6 / 7) train acc: 0.708000; val_acc: 0.629000
(Iteration 381 / 427) loss: 0.870140
(Iteration 401 / 427) loss: 0.828386
(Iteration 421 / 427) loss: 0.870131
(Epoch 7 / 7) train acc: 0.729000; val_acc: 0.657000
Validation set accuracy: 0.657
Test set accuracy: 0.65
```

In [ ]:

```

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.
        # ===== #

        C, H, W = input_dim

        size_W1 = (num_filters, C, filter_size, filter_size)
        size_b1 = num_filters

        Con_output = (num_filters, C, H, W)
        size_W2 = (hidden_dim, (H//2)*(W//2)*num_filters)
        size_b2 = hidden_dim

        size_W3 = (num_classes, hidden_dim)
        size_b3 = num_classes

```

```

self.params['W1'] = np.random.normal(loc=0.0, scale=weight_scale, size = size_W1)
self.params['b1'] = np.zeros(size_b1)
self.params['W2'] = np.random.normal(loc=0.0, scale=weight_scale, size = size_W2).T
self.params['b2'] = np.zeros(size_b2)
self.params['W3'] = np.random.normal(loc=0.0, scale=weight_scale, size = size_W3).T
self.params['b3'] = np.zeros(size_b3)

# ===== #
# 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".
    # ===== #

    out_pool_fast, cache_pool_fast = conv_relu_pool_forward(X, W1, b1, conv_param, pool_param)
    out_affine_ReLU_fast, cache_affine_ReLU_fast = affine_relu_forward(out_pool_fast, W2, b2)
    scores, cache_affine_fast = affine_forward(out_affine_ReLU_fast, 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, dz = softmax_loss(scores, y)
    loss += 0.5 * self.reg * (np.sum(W1 * W1) + np.sum(W2 * W2) + np.sum(W3 * W3))

    dh_affine, dw3, db3 = affine_backward(dz, cache_affine_fast)

```

```
dh_affine_relu, dw2, db2 = affine_relu_backward(dh_affine, cache_affine_ReLU_fast)

dh_Conv, dw1, db1 = conv_relu_pool_backward(dh_affine_relu, cache_pool_fast)

grads['W1'] = dw1 + self.reg * W1
grads['b1'] = db1
grads['W2'] = dw2 + self.reg * W2
grads['b2'] = db2
grads['W3'] = dw3 + self.reg * W3
grads['b3'] = db3

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

return loss, grads
```

pass

```
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
```

```

import numpy as np
from nn.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 WW.

    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
       $H' = 1 + (H + 2 * \text{pad} - \text{HH}) / \text{stride}$ 
       $W' = 1 + (W + 2 * \text{pad} - \text{WW}) / \text{stride}$ 
    - 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, _, HH, WW = w.shape

    H_conv = 1 + (H + 2 * pad - HH) // stride
    W_conv = 1 + (W + 2 * pad - WW) // stride

    out_shape = (N, F, H_conv, W_conv)
    out = np.zeros(out_shape)

    npad = ((0,0), (0,0), (pad, pad), (pad, pad))
    x = np.pad(x, npad, mode='constant')

    for i in np.arange(N):
        for j in np.arange(F):
            for h1 in np.arange(H_conv):
                for w1 in np.arange(W_conv):
                    start_h = h1*stride
                    start_w = w1*stride
                    x_selected = x[i, :, start_h:(start_h + HH), start_w:(start_w + WW)]
                    w_selected = w[j, :, :, :]
                    b_selected = b[j]
                    out[i, j, h1, w1] = np.sum(x_selected * w_selected) + b_selected

    # ===== #
    # 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
    """

```

```

Returns a tuple of:
- dx: Gradient with respect to x
- dw: Gradient with respect to w
- db: Gradient with respect to b
"""
dx, dw, db = None, None, None

N, F, out_height, out_width = dout.shape
x, w, b, conv_param = cache

stride, pad = [conv_param['stride'], conv_param['pad']]
xpad = np.pad(x, ((0,0), (0,0), (pad,pad), (pad,pad)), mode='constant')
num_filts, _, f_height, f_width = w.shape

# ===== #
# YOUR CODE HERE:
# Implement the backward pass of a convolutional neural network.
# Calculate the gradients: dx, dw, and db.
# ===== #

dx = np.zeros(x.shape)
dw = np.zeros(w.shape)
db = np.zeros(b.shape)

N, C, H, W = x.shape

HConv = 1 + (H - f_height) // stride
WConv = 1 + (W - f_width) // stride

for i in np.arange(N):
    for j in np.arange(num_filts):
        if (i == 0):
            db[j] += np.sum(dout[:, j, :, :])

        for h1 in np.arange(HConv):
            start_h = h1 * stride
            for w1 in np.arange(WConv):
                start_w = w1 * stride

                upstream_selected = dout[i, j, h1, w1]
                x_selected = x[i, :, start_h : (start_h + f_height), start_w : (start_w + f_width)]
                w_selected = w[j, :, :, :]

                dx[i, :, start_h : (start_h + f_height), start_w : (start_w + f_width)] += w[j, :, :, :] * dout
                dw[j, :, :, :] += x[i, :, start_h : (start_h + f_height), start_w : (start_w + f_width)] * dout

dx = dx[:, :, pad : -pad, pad : -pad]

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

return dx, dw, db

def max_pool_forward_naive(x, pool_param):
    """
    A naive implementation of the forward pass for a max pooling layer.

    Inputs:
    - x: Input data, of shape (N, C, H, W)
    - pool_param: dictionary with the following keys:
    - 'pool_height': The height of each pooling region
    - 'pool_width': The width of each pooling region
    - 'stride': The distance between adjacent pooling regions

    Returns a tuple of:
    - out: Output data
    - cache: (x, pool_param)
    """
    out = None

    # ===== #
    # YOUR CODE HERE:
    # Implement the max pooling forward pass.
    # ===== #

```



```

f_height = pool_param['pool_height']
f_width = pool_param['pool_width']
stride = pool_param['stride']

N, C, H, W = x.shape

HMaxPool = 1 + (H - f_height) // stride
WMaxPool = 1 + (W - f_width) // stride

out_shape = (N, C, HMaxPool, WMaxPool)
out = np.zeros(out_shape)

for i in np.arange(N):
    for j in np.arange(C):
        for h1 in np.arange(HMaxPool):
            for w1 in np.arange(WMaxPool):
                start_h = h1*stride
                start_w = w1*stride
                x_selected = x[i,j,start_h:(start_h + f_height),start_w:(start_w + f_width)]
                out[i, j, h1, w1] = np.max(x_selected)

# ===== #
# END YOUR CODE HERE
# ===== #
cache = (x, pool_param)
return out, cache

def max_pool_backward_naive(dout, cache):
    """
    A naive implementation of the backward pass for a max pooling layer.

    Inputs:
    - dout: Upstream derivatives
    - cache: A tuple of (x, pool_param) as in the forward pass.

    Returns:
    - dx: Gradient with respect to x
    """
    dx = None
    x, pool_param = cache
    pool_height, pool_width, stride = pool_param['pool_height'], pool_param['pool_width'], pool_param['stride']

    # ===== #
    # YOUR CODE HERE:
    # Implement the max pooling backward pass.
    # ===== #

    N, C, H, W = x.shape

    # height and width of output
    HMaxPool = 1 + (H - pool_height) // stride
    WMaxPool = 1 + (W - pool_width) // stride

    dx = np.zeros(x.shape)

    # for loops
    for i in np.arange(N):
        for j in np.arange(C):
            for h1 in np.arange(HMaxPool):
                for w1 in np.arange(WMaxPool):
                    start_h = h1*stride
                    start_w = w1*stride

                    x_selected = x[i,j,start_h:(start_h + pool_height),start_w:(start_w + pool_width)]
                    upstream_selected = dout[i,j,h1,w1]
                    local_gradient = (x_selected==np.max(x_selected))

                    dx[i,j,start_h:(start_h+pool_height),start_w:(start_w+pool_width)] += local_gradient * upst

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

    return dx

def spatial_batchnorm_forward(x, gamma, beta, bn_param):

```

```

"""
Computes the forward pass for spatial batch normalization.

Inputs:
- x: Input data of shape (N, C, H, W)
- gamma: Scale parameter, of shape (C,)
- beta: Shift parameter, of shape (C,)
- bn_param: Dictionary with the following keys:
- mode: 'train' or 'test'; required
- eps: Constant for numeric stability
- momentum: Constant for running mean / variance. momentum=0 means that
  old information is discarded completely at every time step, while
  momentum=1 means that new information is never incorporated. The
  default of momentum=0.9 should work well in most situations.
- running_mean: Array of shape (D,) giving running mean of features
- running_var: Array of shape (D,) giving running variance of features

Returns a tuple of:
- out: Output data, of shape (N, C, H, W)
- cache: Values needed for the backward pass
"""
out, cache = None, None

# ===== #
# YOUR CODE HERE:
#   Implement the spatial batchnorm forward pass.
#
#   You may find it useful to use the batchnorm forward pass you
#   implemented in HW #4.
# ===== #

N, C, H, W = x.shape
x = x.reshape((N,H,W,C))
x = x.reshape((N*H*W,C))

out, cache = batchnorm_forward(x, gamma, beta, bn_param)
out = out.T
out = out.reshape(C,N,H,W)
out = out.swapaxes(0,1)

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

return out, cache

def spatial_batchnorm_backward(dout, cache):
    """
    Computes the backward pass for spatial batch normalization.

    Inputs:
    - dout: Upstream derivatives, of shape (N, C, H, W)
    - cache: Values from the forward pass

    Returns a tuple of:
    - dx: Gradient with respect to inputs, of shape (N, C, H, W)
    - dgamma: Gradient with respect to scale parameter, of shape (C,)
    - dbeta: Gradient with respect to shift parameter, of shape (C,)
    """
    dx, dgamma, dbeta = None, None, None

    # ===== #
    # YOUR CODE HERE:
    #   Implement the spatial batchnorm backward pass.
    #
    #   You may find it useful to use the batchnorm forward pass you
    #   implemented in HW #4.
    # ===== #

    N, C, H, W = dout.shape

    dout = dout.swapaxes(0,1)
    dout = dout.reshape(C,N*H*W)
    dout = dout.T

    dx, dgamma, dbeta = batchnorm_backward(dout, cache)

```

```
dx = dx.reshape((N,C,H,W))
dgamma = dgamma.reshape((C,))
dbeta = dbeta.reshape((C,))

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

return dx, dgamma, dbeta
```

```

from nndl.layers import *
from utils.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
from nndl.layer_utils import affine_relu_forward, affine_relu_backward
from nndl.fc_net import FullyConnectedNet

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

def affine_forward_test():
    # Test the affine_forward function

    num_inputs = 2
    input_shape = (4, 5, 6)
    output_dim = 3

    input_size = num_inputs * np.prod(input_shape)
    weight_size = output_dim * np.prod(input_shape)

    x = np.linspace(-0.1, 0.5, num=input_size).reshape(num_inputs, *input_shape)
    w = np.linspace(-0.2, 0.3, num=weight_size).reshape(np.prod(input_shape), output_dim)
    b = np.linspace(-0.3, 0.1, num=output_dim)

    out, _ = affine_forward(x, w, b)
    correct_out = np.array([[ 1.49834967,  1.70660132,  1.91485297],
                           [ 3.25553199,  3.5141327,  3.77273342]])

    # Compare your output with ours. The error should be around 1e-9.
    print('If affine_forward function is working, difference should be less than 1e-9:')
    print('difference: {}'.format(rel_error(out, correct_out)))

def affine_backward_test():
    # Test the affine_backward function

    x = np.random.randn(10, 2, 3)
    w = np.random.randn(6, 5)
    b = np.random.randn(5)
    dout = np.random.randn(10, 5)

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

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

    # The error should be around 1e-10
    print('If affine_backward is working, error should be less than 1e-9:')
    print('dx error: {}'.format(rel_error(dx_num, dx)))
    print('dw error: {}'.format(rel_error(dw_num, dw)))
    print('db error: {}'.format(rel_error(db_num, db)))

def relu_forward_test():
    # Test the relu_forward function

    x = np.linspace(-0.5, 0.5, num=12).reshape(3, 4)

    out, _ = relu_forward(x)
    correct_out = np.array([[ 0.,          0.,          0.,          0.],
                           [ 0.,          0.,          0.04545455,  0.13636364],
                           [ 0.22727273,  0.31818182,  0.40909091,  0.5,          ]])

    # Compare your output with ours. The error should be around 1e-8
    print('If relu_forward function is working, difference should be around 1e-8:')
    print('difference: {}'.format(rel_error(out, correct_out)))

def relu_backward_test():
    x = np.random.randn(10, 10)
    dout = np.random.randn(*x.shape)

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

```

```

_, cache = relu_forward(x)
dx = relu_backward(dout, cache)

# The error should be around 1e-12
print('If relu_forward function is working, error should be less than 1e-9:')
print('dx error: {}'.format(rel_error(dx_num, dx)))

def affine_relu_test():

    x = np.random.randn(2, 3, 4)
    w = np.random.randn(12, 10)
    b = np.random.randn(10)
    dout = np.random.randn(2, 10)

    out, cache = affine_relu_forward(x, w, b)
    dx, dw, db = affine_relu_backward(dout, cache)

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

    print('If affine_relu_forward and affine_relu_backward are working, error should be less than 1e-9:')
    print('dx error: {}'.format(rel_error(dx_num, dx)))
    print('dw error: {}'.format(rel_error(dw_num, dw)))
    print('db error: {}'.format(rel_error(db_num, db)))

def fc_net_test():
    N, D, H1, H2, C = 2, 15, 20, 30, 10
    X = np.random.randn(N, D)
    y = np.random.randint(C, size=(N,))

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

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

        for name in sorted(grads):
            f = lambda _: model.loss(X, y)[0]
            grad_num = eval_numerical_gradient(f, model.params[name], verbose=False, h=1e-5)
            print('{} relative error: {}'.format(name, rel_error(grad_num, grads[name])))

```

```
from .layers import *
```

```
def affine_relu_forward(x, w, b):  
    """
```

Convenience layer that performs an affine transform followed by a ReLU

Inputs:

- x: Input to the affine layer
- w, b: Weights for the affine layer

Returns a tuple of:

- out: Output from the ReLU
- cache: Object to give to the backward pass

```
    """
```

```
    a, fc_cache = affine_forward(x, w, b)  
    out, relu_cache = relu_forward(a)  
    cache = (fc_cache, relu_cache)  
    return out, cache
```

```
def affine_relu_backward(dout, cache):  
    """
```

Backward pass for the affine-relu convenience layer

```
    """
```

```
    fc_cache, relu_cache = cache  
    da = relu_backward(dout, relu_cache)  
    dx, dw, db = affine_backward(da, fc_cache)  
    return dx, dw, db
```

```
def affine_batchnorm_relu_forward(x, w, b, gamma, beta, bn_param):
```

```
    aff_out, aff_cache = affine_forward(x, w, b)  
    batch_out, batch_cache = batchnorm_forward(aff_out, gamma, beta, bn_param)  
    out, relu_cache = relu_forward(batch_out)  
    cache = (aff_cache, relu_cache, batch_cache)  
    return out, cache
```

```
def affine_batchnorm_relu_backward(dout, cache):
```

```
    aff_cache, relu_cache, batch_cache = cache  
    dbatch = relu_backward(dout, relu_cache)  
    daffine, dgamma, dbeta = batchnorm_backward(dbatch, batch_cache)  
    dx, dw, db = affine_backward(daffine, aff_cache)  
    return dx, dw, db, dgamma, dbeta
```

```

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.
    # ===== #

    x_transformed = x.reshape(x.shape[0], -1)
    out = x_transformed @ 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.
    # ===== #

    x_transformed = x.reshape(x.shape[0], -1)
    dx = dout @ w.T
    dx = dx.reshape(x.shape)
    dw = x_transformed.T @ dout
    db = np.sum(dout, axis=0)

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

```

```

    return dx, dw, db

def relu_forward(x):
    """
    Computes the forward pass for a layer of rectified linear units (ReLU).

    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.
    # ===== #
    relu = lambda x: x * (x > 0)
    out = relu(x)
    # ===== #
    # 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 (ReLU).

    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 * (x >= 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:

- x: Data of shape (N, D)
- gamma: Scale parameter of shape (D,)
- beta: Shift parameter of shape (D,)
- bn\_param: Dictionary with the following keys:
- mode: 'train' or 'test'; required
- eps: Constant for numeric stability
- momentum: Constant for running mean / variance.
- running\_mean: Array of shape (D,) giving running mean of features
- running\_var: Array of shape (D,) giving running variance of features

Returns a tuple of:

- out: of shape (N, D)
- cache: A tuple of values needed in the backward pass

```

mode = bn_param['mode']
eps = bn_param.get('eps', 1e-5)
momentum = bn_param.get('momentum', 0.9)

N, D = x.shape
running_mean = bn_param.get('running_mean', np.zeros(D, dtype=x.dtype))
running_var = bn_param.get('running_var', np.zeros(D, dtype=x.dtype))

out, cache = None, None
if mode == 'train':

    # ===== #
    # YOUR CODE HERE:
    #   A few steps here:
    #   (1) Calculate the running mean and variance of the minibatch.
    #   (2) Normalize the activations with the running mean and variance.
    #   (3) Scale and shift the normalized activations. Store this
    #       as the variable 'out'
    #   (4) Store any variables you may need for the backward pass in
    #       the 'cache' variable.
    # ===== #

    Mean_minibatch = x.mean(axis=0)
    Variance_minibatch = np.var(x, axis=0)

    running_mean = momentum * running_mean + (1.0 - momentum) * Mean_minibatch
    running_var = momentum * running_var + (1.0 - momentum) * Variance_minibatch

    Mean_minibatch = np.expand_dims(Mean_minibatch, axis=0)
    Variance_minibatch = np.expand_dims(Variance_minibatch, axis=0)

    x_normal = (x - Mean_minibatch) / (np.sqrt(Variance_minibatch + eps))
    out = np.expand_dims(gamma, axis=0) * x_normal + np.expand_dims(beta, axis=0)
    cache = (Mean_minibatch, Variance_minibatch, x_normal, gamma, beta, x, eps)

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

elif mode == 'test':

    # ===== #
    # YOUR CODE HERE:
    #   Calculate the testing time normalized activation. Normalize using
    #   the running mean and variance, and then scale and shift appropriately.
    #   Store the output as 'out'.
    # ===== #

    x_norm = (x - running_mean) / (np.sqrt(running_var + eps))
    out = np.expand_dims(gamma, axis=0) * x_norm + np.expand_dims(beta, axis=0)

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

```

```

else:
    raise ValueError('Invalid forward batchnorm mode "%s"' % mode)

# Store the updated running means back into bn_param
bn_param['running_mean'] = running_mean
bn_param['running_var'] = running_var

return out, cache

def batchnorm_backward(dout, cache):
    """
    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.
    # ===== #

    Mean_minibatch = cache[0]
    Variance_minibatch = cache[1]
    x_normal = cache[2]
    gamma = cache[3]
    beta = cache[4]
    x = cache[5]
    eps = cache[6]

    M = x_normal.shape[0]
    std = np.sqrt(Variance_minibatch + eps)

    dbeta = dout.sum(axis=0)
    dgamma = np.sum(x_normal * dout, axis = 0)
    dx_hat = gamma * dout
    da = (1.0 / std) * dx_hat
    dMu = np.sum(-da, axis = 0)
    dVar = np.sum((-1.0 / (2*np.power(std,3))) * dx_hat * (x - Mean_minibatch), axis = 0)
    dx = (1.0 / std) * dx_hat + (1.0 / M) * dMu + (2.0 / M) * dVar * (x - Mean_minibatch)

    # ===== #
    # 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)
    out = mask * x

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

elif mode == 'test':
    # ===== #
    # YOUR CODE HERE:
    # Implement the inverted dropout forward pass during test time.
    # ===== #

    out = x

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

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

return out, cache

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

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

    dx = None
    if mode == 'train':
        # ===== #
        # YOUR CODE HERE:
        # Implement the inverted dropout backward pass during training time.
        # ===== #

        dx = mask * dout

        # ===== #
        # END YOUR CODE HERE
        # ===== #
    elif mode == 'test':
        # ===== #
        # YOUR CODE HERE:
        # Implement the inverted dropout backward pass during test time.
        # ===== #

```

```

    dx = dout

    # ===== #
    # END YOUR CODE HERE
    # ===== #
    return dx

def svm_loss(x, y):
    """
    Computes the loss and gradient using for multiclass SVM 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
    """
    N = x.shape[0]
    correct_class_scores = x[np.arange(N), y]
    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

```

```
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
# ===== #

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_old = v
    v = config['momentum'] * v_old - config['learning_rate'] * dw
    next_w = w + v + config['momentum'] * (v - v_old)

    # ===== #
    # 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
coeff = np.ones_like(w) / ( np.sqrt(config['a']) + config['epsilon'] * np.ones_like(w) )
next_w = w - config['learning_rate'] * (dw * coeff)
```

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

```
# END YOUR CODE HERE
```

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

```
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['v'] = config['beta1'] * config['v'] + (1 - config['beta1']) * dw
config['a'] = config['beta2'] * config['a'] + (1 - config['beta2']) * dw * dw
config['t'] += 1

a_corrected = np.sqrt( config['a'] / (1 - np.power(config['beta2'],config['t'])) )
v_corrected = config['v'] / (1 - np.power(config['beta1'],config['t']))

coeff = np.ones_like(w) / ( a_corrected + (config['epsilon'] * np.ones_like(w)) )
next_w = w - ( config['learning_rate'] * v_corrected * coeff )

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

return next_w, config
```