

Convolutional neural network layers

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

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

```
In [1]: 1  ## Import and setups
2
3  import time
4  import numpy as np
5  import matplotlib.pyplot as plt
6  from nndl.conv_layers import *
7  from cs231n.data_utils import get_CIFAR10_data
8  from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
9  from cs231n.solver import Solver
10
11 %matplotlib inline
12 plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
13 plt.rcParams['image.interpolation'] = 'nearest'
14 plt.rcParams['image.cmap'] = 'gray'
15
16 # for auto-reloading external modules
17 # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
18 %load_ext autoreload
19 %autoreload 2
20
21 def rel_error(x, y):
22     """ returns relative error """
23     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]: 1  x_shape = (2, 3, 4, 4)
2  w_shape = (3, 3, 4, 4)
3  x = np.linspace(-0.1, 0.5, num=np.prod(x_shape)).reshape(x_shape)
4  w = np.linspace(-0.2, 0.3, num=np.prod(w_shape)).reshape(w_shape)
5  b = np.linspace(-0.1, 0.2, num=3)
6
7  conv_param = {'stride': 2, 'pad': 1}
8  out, _ = conv_forward_naive(x, w, b, conv_param)
9  correct_out = np.array([[[[-0.08759809, -0.10987781],
10                             [-0.18387192, -0.2109216 ]],
11                             [ 0.21027089,  0.21661097],
12                             [ 0.22847626,  0.23004637]],
13                             [ 0.50813986,  0.54309974],
14                             [ 0.64082444,  0.67101435]]],
15                             [[[-0.98053589, -1.03143541],
16                             [-1.19128892, -1.24695841]],
17                             [ 0.69108355,  0.66880383],
18                             [ 0.59480972,  0.56776003]],
19                             [ 2.36270298,  2.36904306],
20                             [ 2.38090835,  2.38247847]]]])
21
22 # Compare your output to ours; difference should be around 1e-8
23 print('Testing conv_forward_naive')
24 print('difference: ', rel_error(out, correct_out))
```

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

Convolutional backward pass

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

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

```
In [3]: 1 x = np.random.randn(4, 3, 5, 5)
2 w = np.random.randn(2, 3, 3, 3)
3 b = np.random.randn(2,)
4 dout = np.random.randn(4, 2, 5, 5)
5 conv_param = {'stride': 1, 'pad': 1}
6
7 out, cache = conv_forward_naive(x,w,b,conv_param)
8
9 dx_num = eval_numerical_gradient_array(lambda x: conv_forward_naive(x, w, b, conv_param)[0], x, dout)
10 dw_num = eval_numerical_gradient_array(lambda w: conv_forward_naive(x, w, b, conv_param)[0], w, dout)
11 db_num = eval_numerical_gradient_array(lambda b: conv_forward_naive(x, w, b, conv_param)[0], b, dout)
12
13 out, cache = conv_forward_naive(x, w, b, conv_param)
14 dx, dw, db = conv_backward_naive(dout, cache)
15
16 # Your errors should be around 1e-9'
17 print('Testing conv_backward_naive function')
18 print('dx error: ', rel_error(dx, dx_num))
19 print('dw error: ', rel_error(dw, dw_num))
20 print('db error: ', rel_error(db, db_num))
```

```
Testing conv_backward_naive function
dx error: 3.5280686226730303e-09
dw error: 2.72961002762906e-09
db error: 1.8822378684961035e-11
```

Max pool forward pass

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

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

```
In [4]: 1 x_shape = (2, 3, 4, 4)
2 x = np.linspace(-0.3, 0.4, num=np.prod(x_shape)).reshape(x_shape)
3 pool_param = {'pool_width': 2, 'pool_height': 2, 'stride': 2}
4
5 out, _ = max_pool_forward_naive(x, pool_param)
6
7 correct_out = np.array([[[[-0.26315789, -0.24842105],
8                             [-0.20421053, -0.18947368]],
9                             [[-0.14526316, -0.13052632],
10                             [-0.08631579, -0.07157895]],
11                             [[-0.02736842, -0.01263158],
12                             [ 0.03157895,  0.04631579]]],
13                             [[[ 0.09052632,  0.10526316],
14                             [ 0.14947368,  0.16421053]],
15                             [[ 0.20842105,  0.22315789],
16                             [ 0.26736842,  0.28210526]],
17                             [[ 0.32631579,  0.34105263],
18                             [ 0.38526316,  0.4          ]]]]])
19
20 # Compare your output with ours. Difference should be around 1e-8.
21 print('Testing max_pool_forward_naive function:')
22 print('difference: ', rel_error(out, correct_out))
```

```
Testing max_pool_forward_naive function:
difference: 4.1666665157267834e-08
```

Max pool backward pass

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

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

```
In [5]: 1 x = np.random.randn(3, 2, 8, 8)
2 dout = np.random.randn(3, 2, 4, 4)
3 pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}
4
5 dx_num = eval_numerical_gradient_array(lambda x: max_pool_forward_naive(x, pool_param)[0], x, dout)
6
7 out, cache = max_pool_forward_naive(x, pool_param)
8 dx = max_pool_backward_naive(dout, cache)
9
10 # Your error should be around 1e-12
11 print('Testing max_pool_backward_naive function:')
12 print('dx error: ', rel_error(dx, dx_num))
```

Testing max_pool_backward_naive function:
dx error: 3.2756189518889396e-12

Fast implementation of the CNN layers

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

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

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

NOTE: The fast implementation for pooling will only perform optimally if the pooling regions are non-overlapping and tile the input. If these conditions are not met then the fast pooling implementation will not be much faster than the naive implementation.

You can compare the performance of the naive and fast versions of these layers by running the cell below.

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

```
In [8]: 1 from cs231n.fast_layers import conv_forward_fast, conv_backward_fast
2 from time import time
3
4 x = np.random.randn(100, 3, 31, 31)
5 w = np.random.randn(25, 3, 3, 3)
6 b = np.random.randn(25,)
7 dout = np.random.randn(100, 25, 16, 16)
8 conv_param = {'stride': 2, 'pad': 1}
9
10 t0 = time()
11 out_naive, cache_naive = conv_forward_naive(x, w, b, conv_param)
12 t1 = time()
13 out_fast, cache_fast = conv_forward_fast(x, w, b, conv_param)
14 t2 = time()
15
16 print('Testing conv_forward_fast:')
17 print('Naive: %fs' % (t1 - t0))
18 print('Fast: %fs' % (t2 - t1))
19 print('Speedup: %fx' % ((t1 - t0) / (t2 - t1)))
20 print('Difference: ', rel_error(out_naive, out_fast))
21
22 t0 = time()
23 dx_naive, dw_naive, db_naive = conv_backward_naive(dout, cache_naive)
24 t1 = time()
25 dx_fast, dw_fast, db_fast = conv_backward_fast(dout, cache_fast)
26 t2 = time()
27
28 print('\nTesting conv_backward_fast:')
29 print('Naive: %fs' % (t1 - t0))
30 print('Fast: %fs' % (t2 - t1))
31 print('Speedup: %fx' % ((t1 - t0) / (t2 - t1)))
32 print('dx difference: ', rel_error(dx_naive, dx_fast))
33 print('dw difference: ', rel_error(dw_naive, dw_fast))
34 print('db difference: ', rel_error(db_naive, db_fast))
```

Testing conv_forward_fast:
Naive: 3.939321s
Fast: 0.009231s
Speedup: 426.744925x
Difference: 3.0126382996499765e-10

Testing conv_backward_fast:
Naive: 8.115447s
Fast: 0.006112s
Speedup: 1327.767670x
dx difference: 6.301465715688952e-11
dw difference: 8.522502018796885e-13
db difference: 4.689242667466879e-15

```
In [16]: 1 from cs231n.fast_layers import max_pool_forward_fast, max_pool_backward_fast
2
3 x = np.random.randn(100, 3, 32, 32)
4 dout = np.random.randn(100, 3, 16, 16)
5 pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}
6
7 t0 = time()
8 out_naive, cache_naive = max_pool_forward_naive(x, pool_param)
9 t1 = time()
10 out_fast, cache_fast = max_pool_forward_fast(x, pool_param)
11 t2 = time()
12
13 print('Testing pool_forward_fast:')
14 print('Naive: %fs' % (t1 - t0))
15 print('fast: %fs' % (t2 - t1))
16 print('speedup: %fx' % ((t1 - t0) / (t2 - t1)))
17 print('difference: ', rel_error(out_naive, out_fast))
18
19 t0 = time()
20 dx_naive = max_pool_backward_naive(dout, cache_naive)
21 t1 = time()
22 dx_fast = max_pool_backward_fast(dout, cache_fast)
23 t2 = time()
24
25 print('\nTesting pool_backward_fast:')
26 print('Naive: %fs' % (t1 - t0))
27 print('speedup: %fx' % ((t1 - t0) / (t2 - t1)))
28 print('dx difference: ', rel_error(dx_naive, dx_fast))
```

Testing pool_forward_fast:

Naive: 0.264137s

fast: 0.001603s

speedup: 164.763682x

difference: 0.0

Testing pool_backward_fast:

Naive: 0.324832s

speedup: 39.004981x

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 [17]: 1 from nndl.conv_layer_utils import conv_relu_pool_forward, conv_relu_pool_backward
2
3 x = np.random.randn(2, 3, 16, 16)
4 w = np.random.randn(3, 3, 3, 3)
5 b = np.random.randn(3,)
6 dout = np.random.randn(2, 3, 8, 8)
7 conv_param = {'stride': 1, 'pad': 1}
8 pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}
9
10 out, cache = conv_relu_pool_forward(x, w, b, conv_param, pool_param)
11 dx, dw, db = conv_relu_pool_backward(dout, cache)
12
13 dx_num = eval_numerical_gradient_array(lambda x: conv_relu_pool_forward(x, w, b, conv_param, pool_param)[0], x, dout)
14 dw_num = eval_numerical_gradient_array(lambda w: conv_relu_pool_forward(x, w, b, conv_param, pool_param)[0], w, dout)
15 db_num = eval_numerical_gradient_array(lambda b: conv_relu_pool_forward(x, w, b, conv_param, pool_param)[0], b, dout)
16
17 print('Testing conv_relu_pool')
18 print('dx error: ', rel_error(dx_num, dx))
19 print('dw error: ', rel_error(dw_num, dw))
20 print('db error: ', rel_error(db_num, db))
```

Testing conv_relu_pool

dx error: 1.883715579397943e-08

dw error: 7.239767601200679e-10

db error: 2.0588404424174744e-11

```
In [19]: 1 from nndl.conv_layer_utils import conv_relu_forward, conv_relu_backward
2
3 x = np.random.randn(2, 3, 8, 8)
4 w = np.random.randn(3, 3, 3, 3)
5 b = np.random.randn(3,)
6 dout = np.random.randn(2, 3, 8, 8)
7 conv_param = {'stride': 1, 'pad': 1}
8
9 out, cache = conv_relu_forward(x, w, b, conv_param)
10 dx, dw, db = conv_relu_backward(dout, cache)
11
12 dx_num = eval_numerical_gradient_array(lambda x: conv_relu_forward(x, w, b, conv_param)[0], x, dout)
13 dw_num = eval_numerical_gradient_array(lambda w: conv_relu_forward(x, w, b, conv_param)[0], w, dout)
14 db_num = eval_numerical_gradient_array(lambda b: conv_relu_forward(x, w, b, conv_param)[0], b, dout)
15
16 print('Testing conv_relu:')
17 print('dx error: ', rel_error(dx_num, dx))
18 print('dw error: ', rel_error(dw_num, dw))
19 print('db error: ', rel_error(db_num, db))
```

Testing conv_relu:

dx error: 1.7866767774698863e-09

dw error: 2.2984323863082598e-10

db error: 2.8699806422645387e-11

What next?

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

Spatial batch normalization

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

Normally batch-normalization accepts inputs of shape (N, D) and produces outputs of shape (N, D) , where we normalize across the minibatch dimension N . For data coming from convolutional layers, batch normalization accepts inputs of shape (N, C, H, W) and produces outputs of shape (N, C, H, W) where the N dimension gives the minibatch size and the (H, W) dimensions give the spatial size of the feature map.

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

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

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

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

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

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

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

```
In [1]: 1  ## Import and setups
        2
        3  import time
        4  import numpy as np
        5  import matplotlib.pyplot as plt
        6  from nndl.conv_layers import *
        7  from cs231n.data_utils import get_CIFAR10_data
        8  from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
        9  from cs231n.solver import Solver
        10
        11  %matplotlib inline
        12  plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
        13  plt.rcParams['image.interpolation'] = 'nearest'
        14  plt.rcParams['image.cmap'] = 'gray'
        15
        16  # for auto-reloading external modules
        17  # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
        18  %load_ext autoreload
        19  %autoreload 2
        20
        21  def rel_error(x, y):
        22      """ returns relative error """
        23      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]: 1 # Check the training-time forward pass by checking means and variances
2 # of features both before and after spatial batch normalization
3
4 N, C, H, W = 2, 3, 4, 5
5 x = 4 * np.random.randn(N, C, H, W) + 10
6
7 print('Before spatial batch normalization:')
8 print('  Shape: ', x.shape)
9 print('  Means: ', x.mean(axis=(0, 2, 3)))
10 print('  Stds: ', x.std(axis=(0, 2, 3)))
11
12 # Means should be close to zero and stds close to one
13 gamma, beta = np.ones(C), np.zeros(C)
14 bn_param = {'mode': 'train'}
15 out, _ = spatial_batchnorm_forward(x, gamma, beta, bn_param)
16 print('After spatial batch normalization:')
17 print('  Shape: ', out.shape)
18 print('  Means: ', out.mean(axis=(0, 2, 3)))
19 print('  Stds: ', out.std(axis=(0, 2, 3)))
20
21 # Means should be close to beta and stds close to gamma
22 gamma, beta = np.asarray([3, 4, 5]), np.asarray([6, 7, 8])
23 out, _ = spatial_batchnorm_forward(x, gamma, beta, bn_param)
24 print('After spatial batch normalization (nontrivial gamma, beta):')
25 print('  Shape: ', out.shape)
26 print('  Means: ', out.mean(axis=(0, 2, 3)))
27 print('  Stds: ', out.std(axis=(0, 2, 3)))
```

```
Before spatial batch normalization:
  Shape: (2, 3, 4, 5)
  Means: [10.45102498 10.2588575 10.073784 ]
  Stds: [3.98071961 3.61635859 3.38114213]
After spatial batch normalization:
  Shape: (2, 3, 4, 5)
  Means: [-3.95516953e-16 1.16573418e-16 8.77076189e-16]
  Stds: [0.99999968 0.99999962 0.99999956]
After spatial batch normalization (nontrivial gamma, beta):
  Shape: (2, 3, 4, 5)
  Means: [6. 7. 8.]
  Stds: [2.99999905 3.99999847 4.99999781]
```

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 [4]: 1 N, C, H, W = 2, 3, 4, 5
2 x = 5 * np.random.randn(N, C, H, W) + 12
3 gamma = np.random.randn(C)
4 beta = np.random.randn(C)
5 dout = np.random.randn(N, C, H, W)
6
7 bn_param = {'mode': 'train'}
8 fx = lambda x: spatial_batchnorm_forward(x, gamma, beta, bn_param)[0]
9 fg = lambda a: spatial_batchnorm_forward(x, gamma, beta, bn_param)[0]
10 fb = lambda b: spatial_batchnorm_forward(x, gamma, beta, bn_param)[0]
11
12 dx_num = eval_numerical_gradient_array(fx, x, dout)
13 da_num = eval_numerical_gradient_array(fg, gamma, dout)
14 db_num = eval_numerical_gradient_array(fb, beta, dout)
15
16 _, cache = spatial_batchnorm_forward(x, gamma, beta, bn_param)
17 dx, dgamma, dbeta = spatial_batchnorm_backward(dout, cache)
18 print('dx error: ', rel_error(dx_num, dx))
19 print('dgamma error: ', rel_error(da_num, dgamma))
20 print('dbeta error: ', rel_error(db_num, dbeta))
```

```
dx error: 8.654635401986262e-09
dgamma error: 3.764431527971484e-12
dbeta error: 5.380953961242855e-12
```

```
In [ ]: 1
```

Convolutional neural networks

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

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

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

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

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

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

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

```
In [1]: 1 # As usual, a bit of setup
2
3 import numpy as np
4 import matplotlib.pyplot as plt
5 from nndl.cnn import *
6 from cs231n.data_utils import get_CIFAR10_data
7 from cs231n.gradient_check import eval_numerical_gradient_array, eval_numerical_gradient
8 from nndl.layers import *
9 from nndl.conv_layers import *
10 from cs231n.fast_layers import *
11 from cs231n.solver import Solver
12
13 %matplotlib inline
14 plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
15 plt.rcParams['image.interpolation'] = 'nearest'
16 plt.rcParams['image.cmap'] = 'gray'
17
18 # for auto-reloading external modules
19 # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
20 %load_ext autoreload
21 %autoreload 2
22
23 def rel_error(x, y):
24     """ returns relative error """
25     return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

```
In [2]: 1 # Load the (preprocessed) CIFAR10 data.
2
3 data = get_CIFAR10_data()
4 for k in data.keys():
5     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 [5]: 1 num_inputs = 2
2 input_dim = (3, 16, 16)
3 reg = 0.0
4 num_classes = 10
5 X = np.random.randn(num_inputs, *input_dim)
6 y = np.random.randint(num_classes, size=num_inputs)
7
8 model = ThreeLayerConvNet(num_filters=3, filter_size=3,
9                             input_dim=input_dim, hidden_dim=7,
10                             dtype=np.float64)
11 loss, grads = model.loss(X, y)
12 for param_name in sorted(grads):
13     f = lambda _: model.loss(X, y)[0]
14     param_grad_num = eval_numerical_gradient(f, model.params[param_name], verbose=False, h=1e-6)
15     e = rel_error(param_grad_num, grads[param_name])
16     print('{} max relative error: {}'.format(param_name, rel_error(param_grad_num, grads[param_name])))

W1 max relative error: 0.0011335262589003132
W2 max relative error: 0.0036604111993140593
W3 max relative error: 3.2625481963891465e-05
b1 max relative error: 3.933808619994069e-05
b2 max relative error: 4.6257258586985625e-06
b3 max relative error: 1.2767301506972432e-09
```

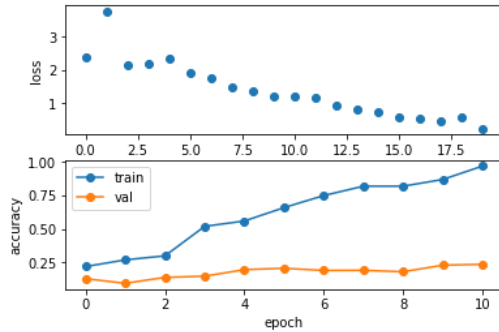
Overfit small dataset

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

```
In [8]: 1 num_train = 100
2 small_data = {
3     'X_train': data['X_train'][:num_train],
4     'y_train': data['y_train'][:num_train],
5     'X_val': data['X_val'],
6     'y_val': data['y_val'],
7 }
8
9 model = ThreeLayerConvNet(weight_scale=1e-2)
10
11 solver = Solver(model, small_data,
12                 num_epochs=10, batch_size=50,
13                 update_rule='adam',
14                 optim_config={
15                     'learning_rate': 1e-3,
16                 },
17                 verbose=True, print_every=1)
18 solver.train()

(Iteration 1 / 20) loss: 2.376132
(Epoch 0 / 10) train acc: 0.220000; val_acc: 0.129000
(Iteration 2 / 20) loss: 3.756471
(Epoch 1 / 10) train acc: 0.270000; val_acc: 0.095000
(Iteration 3 / 20) loss: 2.157400
(Iteration 4 / 20) loss: 2.196071
(Epoch 2 / 10) train acc: 0.300000; val_acc: 0.139000
(Iteration 5 / 20) loss: 2.340025
(Iteration 6 / 20) loss: 1.901690
(Epoch 3 / 10) train acc: 0.520000; val_acc: 0.148000
(Iteration 7 / 20) loss: 1.768312
(Iteration 8 / 20) loss: 1.470933
(Epoch 4 / 10) train acc: 0.560000; val_acc: 0.197000
(Iteration 9 / 20) loss: 1.377451
(Iteration 10 / 20) loss: 1.208964
(Epoch 5 / 10) train acc: 0.660000; val_acc: 0.207000
(Iteration 11 / 20) loss: 1.215419
(Iteration 12 / 20) loss: 1.187998
(Epoch 6 / 10) train acc: 0.750000; val_acc: 0.191000
(Iteration 13 / 20) loss: 0.956437
(Iteration 14 / 20) loss: 0.826714
(Epoch 7 / 10) train acc: 0.820000; val_acc: 0.192000
(Iteration 15 / 20) loss: 0.751160
(Iteration 16 / 20) loss: 0.580117
(Epoch 8 / 10) train acc: 0.820000; val_acc: 0.181000
(Iteration 17 / 20) loss: 0.555041
(Iteration 18 / 20) loss: 0.451633
(Epoch 9 / 10) train acc: 0.870000; val_acc: 0.230000
(Iteration 19 / 20) loss: 0.605399
(Iteration 20 / 20) loss: 0.245934
(Epoch 10 / 10) train acc: 0.970000; val_acc: 0.237000
```

```
In [9]: 1 plt.subplot(2, 1, 1)
2 plt.plot(solver.loss_history, 'o')
3 plt.xlabel('iteration')
4 plt.ylabel('loss')
5
6 plt.subplot(2, 1, 2)
7 plt.plot(solver.train_acc_history, '-o')
8 plt.plot(solver.val_acc_history, '-o')
9 plt.legend(['train', 'val'], loc='upper left')
10 plt.xlabel('epoch')
11 plt.ylabel('accuracy')
12 plt.show()
```



Train the network

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

```
In [10]: 1 model = ThreeLayerConvNet(weight_scale=0.001, hidden_dim=500, reg=0.001)
2
3 solver = Solver(model, data,
4                 num_epochs=1, batch_size=50,
5                 update_rule='adam',
6                 optim_config={
7                     'learning_rate': 1e-3,
8                 },
9                 verbose=True, print_every=20)
10 solver.train()
```

```
(Iteration 1 / 980) loss: 2.304721
(Epoch 0 / 1) train acc: 0.143000; val_acc: 0.146000
(Iteration 21 / 980) loss: 2.226032
(Iteration 41 / 980) loss: 1.903721
(Iteration 61 / 980) loss: 2.024781
(Iteration 81 / 980) loss: 1.947458
(Iteration 101 / 980) loss: 1.664826
(Iteration 121 / 980) loss: 1.865075
(Iteration 141 / 980) loss: 1.998707
(Iteration 161 / 980) loss: 1.880870
(Iteration 181 / 980) loss: 1.845643
(Iteration 201 / 980) loss: 1.738873
(Iteration 221 / 980) loss: 1.726869
(Iteration 241 / 980) loss: 1.855340
(Iteration 261 / 980) loss: 1.745509
(Iteration 281 / 980) loss: 1.721090
(Iteration 301 / 980) loss: 1.816408
(Iteration 321 / 980) loss: 1.650176
(Iteration 341 / 980) loss: 1.617935
(Iteration 361 / 980) loss: 1.904243
(Iteration 381 / 980) loss: 1.795476
(Iteration 401 / 980) loss: 1.459891
(Iteration 421 / 980) loss: 1.944634
(Iteration 441 / 980) loss: 1.988634
(Iteration 461 / 980) loss: 1.434619
(Iteration 481 / 980) loss: 1.801080
(Iteration 501 / 980) loss: 1.665090
(Iteration 521 / 980) loss: 1.725914
(Iteration 541 / 980) loss: 1.858884
(Iteration 561 / 980) loss: 1.624216
(Iteration 581 / 980) loss: 1.828506
(Iteration 601 / 980) loss: 1.385467
(Iteration 621 / 980) loss: 1.615376
(Iteration 641 / 980) loss: 1.677582
(Iteration 661 / 980) loss: 1.493979
(Iteration 681 / 980) loss: 1.557426
(Iteration 701 / 980) loss: 1.650439
(Iteration 721 / 980) loss: 1.488394
(Iteration 741 / 980) loss: 1.675061
(Iteration 761 / 980) loss: 1.746590
(Iteration 781 / 980) loss: 1.644568
(Iteration 801 / 980) loss: 1.731493
(Iteration 821 / 980) loss: 1.351853
(Iteration 841 / 980) loss: 1.614824
(Iteration 861 / 980) loss: 2.077891
(Iteration 881 / 980) loss: 2.098706
(Iteration 901 / 980) loss: 1.751987
(Iteration 921 / 980) loss: 1.686586
(Iteration 941 / 980) loss: 1.525075
(Iteration 961 / 980) loss: 1.623031
(Epoch 1 / 1) train acc: 0.484000; val_acc: 0.475000
```

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 [11]: 1 # ===== #
2 # YOUR CODE HERE:
3 # Implement a CNN to achieve greater than 65% validation accuracy
4 # on CIFAR-10.
5 # ===== #
6
7 model = ThreeLayerConvNet(weight_scale = 0.001,
8                             hidden_dim = 500,
9                             reg = 0.001,
10                             num_filters = 64,
11                             filter_size = 3)
12
13 solver = Solver(model, data,
14                 num_epochs=10, batch_size=500,
15                 update_rule='adam',
16                 optim_config={
17                     'learning_rate': 1e-3,
18                 },
19                 lr_decay=0.9,
20                 verbose=True, print_every=15)
21 solver.train()
22
23 # print out the validation and test accuracy
24 y_val_max = np.argmax(model.loss(data['X_val']), axis=1)
25 y_test_max = np.argmax(model.loss(data['X_test']), axis=1)
26 print('Validation set accuracy: {}'.format(np.mean(y_val_max == data['y_val'])))
27 print('Test set accuracy: {}'.format(np.mean(y_test_max == data['y_test'])))
28
29 # ===== #
30 # END YOUR CODE HERE
31 # ===== #
32

```

```

(Iteration 1 / 980) loss: 2.306659
(Epoch 0 / 10) train acc: 0.090000; val_acc: 0.107000
(Iteration 16 / 980) loss: 1.941937
(Iteration 31 / 980) loss: 1.706596
(Iteration 46 / 980) loss: 1.509531
(Iteration 61 / 980) loss: 1.505432
(Iteration 76 / 980) loss: 1.430747
(Iteration 91 / 980) loss: 1.291258
(Epoch 1 / 10) train acc: 0.516000; val_acc: 0.531000
(Iteration 106 / 980) loss: 1.237608
(Iteration 121 / 980) loss: 1.339353
(Iteration 136 / 980) loss: 1.231740
(Iteration 151 / 980) loss: 1.318899
(Iteration 166 / 980) loss: 1.200950
(Iteration 181 / 980) loss: 1.192369
(Iteration 196 / 980) loss: 1.263459
(Epoch 2 / 10) train acc: 0.641000; val_acc: 0.579000
(Iteration 211 / 980) loss: 1.152634
(Iteration 226 / 980) loss: 1.080518
(Iteration 241 / 980) loss: 1.033657
(Iteration 256 / 980) loss: 1.094278
(Iteration 271 / 980) loss: 1.033186
(Iteration 286 / 980) loss: 1.144037
(Epoch 3 / 10) train acc: 0.692000; val_acc: 0.599000
(Iteration 301 / 980) loss: 0.925115
(Iteration 316 / 980) loss: 1.143499
(Iteration 331 / 980) loss: 1.054726
(Iteration 346 / 980) loss: 0.934795
(Iteration 361 / 980) loss: 0.987350
(Iteration 376 / 980) loss: 0.998432
(Iteration 391 / 980) loss: 0.855865
(Epoch 4 / 10) train acc: 0.729000; val_acc: 0.619000
(Iteration 406 / 980) loss: 0.869595
(Iteration 421 / 980) loss: 0.951113
(Iteration 436 / 980) loss: 0.829878
(Iteration 451 / 980) loss: 0.912258
(Iteration 466 / 980) loss: 0.867191
(Iteration 481 / 980) loss: 0.897211
(Epoch 5 / 10) train acc: 0.769000; val_acc: 0.668000
(Iteration 496 / 980) loss: 0.802053
(Iteration 511 / 980) loss: 0.829430
(Iteration 526 / 980) loss: 0.874852
(Iteration 541 / 980) loss: 0.781208
(Iteration 556 / 980) loss: 0.779746
(Iteration 571 / 980) loss: 0.749799
(Iteration 586 / 980) loss: 0.784848
(Epoch 6 / 10) train acc: 0.761000; val_acc: 0.653000
(Iteration 601 / 980) loss: 0.802043
(Iteration 616 / 980) loss: 0.730665
(Iteration 631 / 980) loss: 0.736980
(Iteration 646 / 980) loss: 0.687644
(Iteration 661 / 980) loss: 0.686610
(Iteration 676 / 980) loss: 0.639770
(Epoch 7 / 10) train acc: 0.794000; val_acc: 0.662000
(Iteration 691 / 980) loss: 0.639489

```

```
(Iteration 706 / 980) loss: 0.656530
(Iteration 721 / 980) loss: 0.660048
(Iteration 736 / 980) loss: 0.605011
(Iteration 751 / 980) loss: 0.683397
(Iteration 766 / 980) loss: 0.614072
(Iteration 781 / 980) loss: 0.591042
(Epoch 8 / 10) train acc: 0.857000; val_acc: 0.678000
(Iteration 796 / 980) loss: 0.615200
(Iteration 811 / 980) loss: 0.579629
(Iteration 826 / 980) loss: 0.574568
(Iteration 841 / 980) loss: 0.531007
(Iteration 856 / 980) loss: 0.521273
(Iteration 871 / 980) loss: 0.518840
(Epoch 9 / 10) train acc: 0.865000; val_acc: 0.657000
(Iteration 886 / 980) loss: 0.571132
(Iteration 901 / 980) loss: 0.597874
(Iteration 916 / 980) loss: 0.521601
(Iteration 931 / 980) loss: 0.442987
(Iteration 946 / 980) loss: 0.475231
(Iteration 961 / 980) loss: 0.453263
(Iteration 976 / 980) loss: 0.501769
(Epoch 10 / 10) train acc: 0.884000; val_acc: 0.674000
Validation set accuracy: 0.678
Test set accuracy: 0.666
```

In []:

1

```

1 import numpy as np
2 import pdb
3
4 """
5 This code was originally written for CS 231n at Stanford University
6 (cs231n.stanford.edu). It has been modified in various areas for use in the
7 ECE 239AS class at UCLA. This includes the descriptions of what code to
8 implement as well as some slight potential changes in variable names to be
9 consistent with class nomenclature. We thank Justin Johnson & Serena Yeung
10 for
11 permission to use this code. To see the original version, please visit
12 cs231n.stanford.edu.
13 """
14 def affine_forward(x, w, b):
15     """
16     Computes the forward pass for an affine (fully-connected) layer.
17
18     The input x has shape (N, d_1, ..., d_k) and contains a minibatch of N
19     examples, where each example x[i] has shape (d_1, ..., d_k). We will
20     reshape each input into a vector of dimension D = d_1 * ... * d_k, and
21     then transform it to an output vector of dimension M.
22
23     Inputs:
24     - x: A numpy array containing input data, of shape (N, d_1, ..., d_k)
25     - w: A numpy array of weights, of shape (D, M)
26     - b: A numpy array of biases, of shape (M,)
27
28     Returns a tuple of:
29     - out: output, of shape (N, M)
30     - cache: (x, w, b)
31     """
32
33     # ===== #
34     # YOUR CODE HERE:
35     # Calculate the output of the forward pass. Notice the dimensions
36     # of w are D x M, which is the transpose of what we did in earlier
37     # assignments.
38     # ===== #
39
40     x_reshape = x.reshape((x.shape[0], w.shape[0])) # N * D
41     out = x_reshape.dot(w) + b.reshape((1, b.shape[0])) # N * M
42
43     # ===== #
44     # END YOUR CODE HERE
45     # ===== #
46
47     cache = (x, w, b)
48     return out, cache
49
50 def affine_backward(dout, cache):
51     """
52     Computes the backward pass for an affine layer.
53
54     Inputs:
55     - dout: Upstream derivative, of shape (N, M)
56     - cache: Tuple of:
57       - x: Input data, of shape (N, d_1, ..., d_k)
58       - w: Weights, of shape (D, M)
59
60     Returns a tuple of:
61     - dx: Gradient with respect to x, of shape (N, d_1, ..., d_k)
62     - dw: Gradient with respect to w, of shape (D, M)
63     - db: Gradient with respect to b, of shape (M,)
64     """
65
66     x, w, b = cache
67     dx, dw, db = None, None, None
68
69     # ===== #
70     # YOUR CODE HERE:
71     # Calculate the gradients for the backward pass.
72     # ===== #
73
74     x_reshape = np.reshape(x, (x.shape[0], w.shape[0]))
75     dx_reshape = dout.dot(w.T)
76     dx = np.reshape(dx_reshape, x.shape) # N * D
77     dw = x_reshape.T.dot(dout) # D * M
78     db = dout.T.dot(np.ones(x.shape[0])) # M * 1
79
80     # ===== #
81     # END YOUR CODE HERE
82     # ===== #
83

```

```

84     return dx, dw, db
85
86 def relu_forward(x):
87     """
88     Computes the forward pass for a layer of rectified linear units (ReLUs).
89
90     Input:
91     - x: Inputs, of any shape
92
93     Returns a tuple of:
94     - out: Output, of the same shape as x
95     - cache: x
96     """
97     # ===== #
98     # YOUR CODE HERE:
99     #   Implement the ReLU forward pass.
100    # ===== #
101
102    out = np.maximum(0, x)
103
104    # ===== #
105    # END YOUR CODE HERE
106    # ===== #
107
108    cache = x
109    return out, cache
110
111 def relu_backward(dout, cache):
112     """
113     Computes the backward pass for a layer of rectified linear units (ReLUs).
114
115     Input:
116     - dout: Upstream derivatives, of any shape
117     - cache: Input x, of same shape as dout
118
119     Returns:
120     - dx: Gradient with respect to x
121     """
122
123     x = cache
124
125     # ===== #
126     # YOUR CODE HERE:
127     #   Implement the ReLU backward pass
128     # ===== #
129
130     dx = (x > 0) * (dout)
131
132     # ===== #
133     # END YOUR CODE HERE
134     # ===== #
135
136     return dx
137
138 def batchnorm_forward(x, gamma, beta, bn_param):
139     """
140     Forward pass for batch normalization.
141
142     During training the sample mean and (uncorrected) sample variance are
143     computed from minibatch statistics and used to normalize the incoming data.
144     During training we also keep an exponentially decaying running mean of the
145     mean and variance of each feature, and these averages are used to normalize data
146     at test-time.
147
148     At each timestep we update the running averages for mean and variance using
149     an exponential decay based on the momentum parameter:
150
151     running_mean = momentum * running_mean + (1 - momentum) * sample_mean
152     running_var = momentum * running_var + (1 - momentum) * sample_var
153
154     Note that the batch normalization paper suggests a different test-time
155     behavior: they compute sample mean and variance for each feature using a
156     large number of training images rather than using a running average. For
157     this implementation we have chosen to use running averages instead since
158     they do not require an additional estimation step; the torch7
159     implementation
160     of batch normalization also uses running averages.
161
162     Input:
163     - x: Data of shape (N, D)
164     - gamma: Scale parameter of shape (D,)
165     - beta: Shift parameter of shape (D,)
166     - bn_param: Dictionary with the following keys:

```



```

166 - mode: 'train' or 'test'; required
167 - eps: Constant for numeric stability
168 - momentum: Constant for running mean / variance.
169 - running_mean: Array of shape (D,) giving running mean of features
170 - running_var Array of shape (D,) giving running variance of features
171
172 Returns a tuple of:
173 - out: of shape (N, D)
174 - cache: A tuple of values needed in the backward pass
175 """
176 mode = bn_param['mode']
177 eps = bn_param.get('eps', 1e-5)
178 momentum = bn_param.get('momentum', 0.9)
179
180 N, D = x.shape
181 running_mean = bn_param.get('running_mean', np.zeros(D, dtype=x.dtype))
182 running_var = bn_param.get('running_var', np.zeros(D, dtype=x.dtype))
183
184 out, cache = None, None
185 if mode == 'train':
186
187     # ===== #
188     # YOUR CODE HERE:
189     #   A few steps here:
190     #   (1) Calculate the running mean and variance of the minibatch.
191     #   (2) Normalize the activations with the batch mean and variance.
192     #   (3) Scale and shift the normalized activations. Store this
193     #       as the variable 'out'
194     #   (4) Store any variables you may need for the backward pass in
195     #       the 'cache' variable.
196     # ===== #
197
198     minibatch_mean = np.mean(x, axis=0)
199     minibatch_var = np.var(x, axis=0)
200     x_normalize = (x - minibatch_mean) / np.sqrt(minibatch_var + eps)
201     out = gamma * x_normalize + beta
202
203     running_mean = momentum * running_mean + (1 - momentum) * minibatch_mean
204     running_var = momentum * running_var + (1 - momentum) * minibatch_var
205     bn_param['running_mean'] = running_mean
206     bn_param['running_var'] = running_var
207
208     cache = {
209         'minibatch_var': minibatch_var,
210         'x_centralize': (x - minibatch_mean),
211         'x_normalize': x_normalize,
212         'gamma': gamma,
213         'eps': eps
214     }
215
216     # ===== #
217     # END YOUR CODE HERE
218     # ===== #
219
220 elif mode == 'test':
221
222     # ===== #
223     # YOUR CODE HERE:
224     #   Calculate the testing time normalized activations. Normalize using
225     #   the running mean and variance, and then scale and shift
226     #   appropriately.
227     #   Store the output as 'out'.
228     # ===== #
229
230     out = gamma * (x - running_mean) / np.sqrt(running_var + eps) + beta
231
232     # ===== #
233     # END YOUR CODE HERE
234     # ===== #
235
236 else:
237     raise ValueError('Invalid forward batchnorm mode "%s"' % mode)
238
239 # Store the updated running means back into bn_param
240 bn_param['running_mean'] = running_mean
241 bn_param['running_var'] = running_var
242
243 return out, cache
244
245 def batchnorm_backward(dout, cache):
246     """
247     Backward pass for batch normalization.
248
249     For this implementation, you should write out a computation graph for

```

```

249 batch normalization on paper and propagate gradients backward through
250 intermediate nodes.
251
252 Inputs:
253 - dout: Upstream derivatives, of shape (N, D)
254 - cache: Variable of intermediates from batchnorm_forward.
255
256 Returns a tuple of:
257 - dx: Gradient with respect to inputs x, of shape (N, D)
258 - dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
259 - dbeta: Gradient with respect to shift parameter beta, of shape (D,)
260 """
261 dx, dgamma, dbeta = None, None, None
262
263 # ===== #
264 # YOUR CODE HERE:
265 # Implement the batchnorm backward pass, calculating dx, dgamma, and
266 dbeta.
267 # ===== #
268
269 # get parameters from cache
270 N = dout.shape[0]
271 minibatch_var = cache.get('minibatch_var')
272 x_centralize = cache.get('x_centralize')
273 x_normalize = cache.get('x_normalize')
274 gamma = cache.get('gamma')
275 eps = cache.get('eps')
276
277 # calculate dx
278 dxhat = dout * gamma
279 dxmu1 = dxhat / np.sqrt(minibatch_var + eps)
280 sqrt_var = np.sqrt(minibatch_var + eps)
281 dsqrt_var = -np.sum(dxhat * x_centralize, axis=0) / (sqrt_var**2)
282 dvar = dsqrt_var * 0.5 / sqrt_var
283 dxmu2 = 2 * x_centralize * dvar * np.ones_like(dout) / N
284 dx1 = dxmu1 + dxmu2
285 dx2 = -np.sum(dx1, axis=0) * np.ones_like(dout) / N
286 dx = dx1 + dx2
287
288 # calculate dbeta and dgamma
289 dbeta = np.sum(dout, axis=0)
290 dgamma = np.sum(dout * x_normalize, axis=0)
291
292 # ===== #
293 # END YOUR CODE HERE
294 # ===== #
295
296 return dx, dgamma, dbeta
297
298 def dropout_forward(x, dropout_param):
299     """
300     Performs the forward pass for (inverted) dropout.
301
302     Inputs:
303     - x: Input data, of any shape
304     - dropout_param: A dictionary with the following keys:
305       - p: Dropout parameter. We drop each neuron output with probability p.
306       - mode: 'test' or 'train'. If the mode is train, then perform dropout;
307       if the mode is test, then just return the input.
308       - seed: Seed for the random number generator. Passing seed makes this
309       function deterministic, which is needed for gradient checking but not
310       in
311       real networks.
312
313     Outputs:
314     - out: Array of the same shape as x.
315     - cache: A tuple (dropout_param, mask). In training mode, mask is the
316     dropout
317     mask that was used to multiply the input; in test mode, mask is None.
318     """
319     p, mode = dropout_param['p'], dropout_param['mode']
320     if 'seed' in dropout_param:
321         np.random.seed(dropout_param['seed'])
322
323     mask = None
324     out = None
325
326     if mode == 'train':
327         # ===== #
328         # YOUR CODE HERE:
329         # Implement the inverted dropout forward pass during training time.
330         # Store the masked and scaled activations in out, and store the
331         # dropout mask as the variable mask.
332         # ===== #

```

```

330
331     mask = (np.random.random_sample(x.shape) >= p) / (1 - p)
332     out = x * mask
333
334     # ===== #
335     # END YOUR CODE HERE
336     # ===== #
337
338 elif mode == 'test':
339
340     # ===== #
341     # YOUR CODE HERE:
342     # Implement the inverted dropout forward pass during test time.
343     # ===== #
344
345     out = x
346
347     # ===== #
348     # END YOUR CODE HERE
349     # ===== #
350
351     cache = (dropout_param, mask)
352     out = out.astype(x.dtype, copy=False)
353
354     return out, cache
355
356 def dropout_backward(dout, cache):
357     """
358     Perform the backward pass for (inverted) dropout.
359
360     Inputs:
361     - dout: Upstream derivatives, of any shape
362     - cache: (dropout_param, mask) from dropout_forward.
363     """
364     dropout_param, mask = cache
365     mode = dropout_param['mode']
366
367     dx = None
368     if mode == 'train':
369         # ===== #
370         # YOUR CODE HERE:
371         # Implement the inverted dropout backward pass during training time.
372         # ===== #
373
374         dx = dout * mask
375
376         # ===== #
377         # END YOUR CODE HERE
378         # ===== #
379     elif mode == 'test':
380         # ===== #
381         # YOUR CODE HERE:
382         # Implement the inverted dropout backward pass during test time.
383         # ===== #
384
385         dx = dout
386
387         # ===== #
388         # END YOUR CODE HERE
389         # ===== #
390     return dx
391
392 def svm_loss(x, y):
393     """
394     Computes the loss and gradient using for multiclass SVM classification.
395
396     Inputs:
397     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth
398     class for the ith input.
399     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
400     0 <= y[i] < C
401
402     Returns a tuple of:
403     - loss: Scalar giving the loss
404     - dx: Gradient of the loss with respect to x
405     """
406     N = x.shape[0]
407     correct_class_scores = x[np.arange(N), y]
408     margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
409     margins[np.arange(N), y] = 0
410     loss = np.sum(margins) / N
411     num_pos = np.sum(margins > 0, axis=1)
412     dx = np.zeros_like(x)

```

```

413 dx[margins > 0] = 1
414 dx[np.arange(N), y] -= num_pos
415 dx /= N
416 return loss, dx
417
418
419 def softmax_loss(x, y):
420     """
421     Computes the loss and gradient for softmax classification.
422
423     Inputs:
424     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth
425       class for the ith input.
426     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
427       0 <= y[i] < C
428
429     Returns a tuple of:
430     - loss: Scalar giving the loss
431     - dx: Gradient of the loss with respect to x
432     """
433
434     probs = np.exp(x - np.max(x, axis=1, keepdims=True))
435     probs /= np.sum(probs, axis=1, keepdims=True)
436     N = x.shape[0]
437     loss = -np.sum(np.log(probs[np.arange(N), y])) / N
438     dx = probs.copy()
439     dx[np.arange(N), y] -= 1
440     dx /= N
441     return loss, dx
442

```

```

1 import numpy as np
2
3 """
4 This code was originally written for CS 231n at Stanford University
5 (cs231n.stanford.edu). It has been modified in various areas for use in the
6 ECE 239AS class at UCLA. This includes the descriptions of what code to
7 implement as well as some slight potential changes in variable names to be
8 consistent with class nomenclature. We thank Justin Johnson & Serena Yeung
9 for
10 permission to use this code. To see the original version, please visit
11 cs231n.stanford.edu.
12 """
13
14 """
15 This file implements various first-order update rules that are commonly used
16 for
17 training neural networks. Each update rule accepts current weights and the
18 gradient of the loss with respect to those weights and produces the next set
19 of
20 weights. Each update rule has the same interface:
21
22 def update(w, dw, config=None):
23
24     Inputs:
25     - w: A numpy array giving the current weights.
26     - dw: A numpy array of the same shape as w giving the gradient of the
27           loss with respect to w.
28     - config: A dictionary containing hyperparameter values such as learning
29               rate,
30               momentum, etc. If the update rule requires caching values over many
31               iterations, then config will also hold these cached values.
32
33     Returns:
34     - next_w: The next point after the update.
35     - config: The config dictionary to be passed to the next iteration of the
36               update rule.
37
38     NOTE: For most update rules, the default learning rate will probably not
39     perform
40     well; however the default values of the other hyperparameters should work
41     well
42     for a variety of different problems.
43
44     For efficiency, update rules may perform in-place updates, mutating w and
45     setting next_w equal to w.
46 """
47
48 def sgd(w, dw, config=None):
49     """
50     Performs vanilla stochastic gradient descent.
51
52     config format:
53     - learning_rate: Scalar learning rate.
54     """
55     if config is None: config = {}
56     config.setdefault('learning_rate', 1e-2)
57
58     w -= config['learning_rate'] * dw
59     return w, config
60
61 def sgd_momentum(w, dw, config=None):
62     """
63     Performs stochastic gradient descent with momentum.
64
65     config format:
66     - learning_rate: Scalar learning rate.
67     - momentum: Scalar between 0 and 1 giving the momentum value.
68                 Setting momentum = 0 reduces to sgd.
69     - velocity: A numpy array of the same shape as w and dw used to store a
70                 moving
71                 average of the gradients.
72     """
73     if config is None: config = {}
74     config.setdefault('learning_rate', 1e-2)
75     config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
76     v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets
77     it to zero.
78
79     # ===== #
80     # YOUR CODE HERE:
81     # Implement the momentum update formula. Return the updated weights
82     # as next_w, and store the updated velocity as v.

```

```

77 # ===== #
78
79 v = config['momentum'] * v - config['learning_rate'] * dw
80 next_w = w + v
81
82 # ===== #
83 # END YOUR CODE HERE
84 # ===== #
85
86 config['velocity'] = v
87
88 return next_w, config
89
90 def sgd_nesterov_momentum(w, dw, config=None):
91     """
92     Performs stochastic gradient descent with Nesterov momentum.
93
94     config format:
95     - learning_rate: Scalar learning rate.
96     - momentum: Scalar between 0 and 1 giving the momentum value.
97       Setting momentum = 0 reduces to sgd.
98     - velocity: A numpy array of the same shape as w and dw used to store a
99       moving
100       average of the gradients.
101     """
102     if config is None: config = {}
103     config.setdefault('learning_rate', 1e-2)
104     config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
105     v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets
106     it to zero.
107
108     # ===== #
109     # YOUR CODE HERE:
110     # Implement the momentum update formula. Return the updated weights
111     # as next_w, and store the updated velocity as v.
112     # ===== #
113
114     v_old = v
115     v = config['momentum'] * v - config['learning_rate'] * dw
116     w += v + config['momentum'] * (v - v_old)
117     next_w = w
118
119     # ===== #
120     # END YOUR CODE HERE
121     # ===== #
122
123     config['velocity'] = v
124
125     return next_w, config
126
127 def rmsprop(w, dw, config=None):
128     """
129     Uses the RMSProp update rule, which uses a moving average of squared
130     gradient
131     values to set adaptive per-parameter learning rates.
132
133     config format:
134     - learning_rate: Scalar learning rate.
135     - decay_rate: Scalar between 0 and 1 giving the decay rate for the squared
136       gradient cache.
137     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
138     - beta: Moving average of second moments of gradients.
139     """
140     if config is None: config = {}
141     config.setdefault('learning_rate', 1e-2)
142     config.setdefault('decay_rate', 0.99)
143     config.setdefault('epsilon', 1e-8)
144     config.setdefault('a', np.zeros_like(w))
145
146     next_w = None
147
148     # ===== #
149     # YOUR CODE HERE:
150     # Implement RMSProp. Store the next value of w as next_w. You need
151     # to also store in config['a'] the moving average of the second
152     # moment gradients, so they can be used for future gradients. Concretely,
153     # config['a'] corresponds to "a" in the lecture notes.
154     # ===== #
155
156     config['a'] = config['decay_rate'] * config['a'] + (1 -
157 config['decay_rate']) * (dw**2)
158     next_w = w - config['learning_rate'] * dw / (np.sqrt(config['a']) +
159 config['epsilon'])

```

```

156 # ===== #
157 # END YOUR CODE HERE
158 # ===== #
159
160 return next_w, config
161
162
163 def adam(w, dw, config=None):
164     """
165     Uses the Adam update rule, which incorporates moving averages of both the
166     gradient and its square and a bias correction term.
167
168     config format:
169     - learning_rate: Scalar learning rate.
170     - beta1: Decay rate for moving average of first moment of gradient.
171     - beta2: Decay rate for moving average of second moment of gradient.
172     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
173     - m: Moving average of gradient.
174     - v: Moving average of squared gradient.
175     - t: Iteration number.
176     """
177     if config is None: config = {}
178     config.setdefault('learning_rate', 1e-3)
179     config.setdefault('beta1', 0.9)
180     config.setdefault('beta2', 0.999)
181     config.setdefault('epsilon', 1e-8)
182     config.setdefault('v', np.zeros_like(w))
183     config.setdefault('a', np.zeros_like(w))
184     config.setdefault('t', 0)
185
186     next_w = None
187
188     # ===== #
189     # YOUR CODE HERE:
190     # Implement Adam. Store the next value of w as next_w. You need
191     # to also store in config['a'] the moving average of the second
192     # moment gradients, and in config['v'] the moving average of the
193     # first moments. Finally, store in config['t'] the increasing time.
194     # ===== #
195
196     beta1 = config['beta1']
197     beta2 = config['beta2']
198     t = config['t'] + 1
199
200     v = beta1 * config['v'] + (1 - beta1) * dw
201     a = beta2 * config['a'] + (1 - beta2) * (dw**2)
202     v_corrected = v / (1 - beta1**t)
203     a_corrected = a / (1 - beta2**t)
204     next_w = w - config['learning_rate'] * v_corrected / (np.sqrt(a_corrected)
+ config['epsilon'])
205
206     config['v'] = v
207     config['a'] = a
208     config['t'] = t
209
210     # ===== #
211     # END YOUR CODE HERE
212     # ===== #
213
214     return next_w, config
215
216
217
218
219
220

```

```

1 import numpy as np
2 from nndl.layers import *
3 import pdb
4
5 """
6 This code was originally written for CS 231n at Stanford University
7 (cs231n.stanford.edu). It has been modified in various areas for use in the
8 ECE 239AS class at UCLA. This includes the descriptions of what code to
9 implement as well as some slight potential changes in variable names to be
10 consistent with class nomenclature. We thank Justin Johnson & Serena Yeung
11 for
12 permission to use this code. To see the original version, please visit
13 cs231n.stanford.edu.
14 """
15 def conv_forward_naive(x, w, b, conv_param):
16     """
17     A naive implementation of the forward pass for a convolutional layer.
18
19     The input consists of N data points, each with C channels, height H and
20     width
21     W. We convolve each input with F different filters, where each filter spans
22     all C channels and has height HH and width WW.
23
24     Input:
25     - x: Input data of shape (N, C, H, W)
26     - w: Filter weights of shape (F, C, HH, WW)
27     - b: Biases, of shape (F,)
28     - conv_param: A dictionary with the following keys:
29         - 'stride': The number of pixels between adjacent receptive fields in the
30           horizontal and vertical directions.
31         - 'pad': The number of pixels that will be used to zero-pad the input.
32
33     Returns a tuple of:
34     - out: Output data, of shape (N, F, H', W') where H' and W' are given by
35            $H' = 1 + (H + 2 * \text{pad} - \text{HH}) / \text{stride}$ 
36            $W' = 1 + (W + 2 * \text{pad} - \text{WW}) / \text{stride}$ 
37     - cache: (x, w, b, conv_param)
38     """
39     out = None
40     pad = conv_param['pad']
41     stride = conv_param['stride']
42
43     # ===== #
44     # YOUR CODE HERE:
45     # Implement the forward pass of a convolutional neural network.
46     # Store the output as 'out'.
47     # Hint: to pad the array, you can use the function np.pad.
48     # ===== #
49
50     N, C, H, W = x.shape # [N, 3, 32, 32]
51     F, C, HH, WW = w.shape # [32, 3, 7, 7]
52
53     padded_x = (np.pad(x, ((0, 0), (0, 0), (pad, pad), (pad, pad)),
54 'constant')) # [N, 3, 38, 38]
55     out_height = np.int(((H + 2 * pad - HH) / stride) + 1) # 32
56     out_width = np.int(((W + 2 * pad - WW) / stride) + 1) # 32
57     out = np.zeros((N, F, out_height, out_width)) # [N, 32, 32, 32]
58
59     for img in range(N): # for each image, do convolutional process
60         for kernal in range(F): # for each channel, there are 3 W (7x7) and 1 b
61             (scalar), linear sum together
62             for row in range(out_height): # from top to bottom
63                 for col in range(out_width): # from left to right
64                     # each kernel has 3 W (7x7), for each elements in W, multiply it
65                     with corresponding elements in original graph
66                     # then add up these 49 numbers together -> 1 scalar
67                     # then add up three scalar and 1 bias, as the current position's
68                     convolutional result
69                     out[img, kernal, row, col] = np.sum(w[kernal, :, :, :] * \
70 padded_x[img, :, :, :],
71 row*stride:row*stride+HH, col*stride:col*stride+WW) + b[kernal]
72
73     # ===== #
74     # END YOUR CODE HERE
75     # ===== #
76
77     cache = (x, w, b, conv_param)
78     return out, cache
79
80 def conv_backward_naive(dout, cache):
81     """
82     A naive implementation of the backward pass for a convolutional layer.

```



```

78
79 Inputs:
80 - dout: Upstream derivatives.
81 - cache: A tuple of (x, w, b, conv_param) as in conv_forward_naive
82
83 Returns a tuple of:
84 - dx: Gradient with respect to x
85 - dw: Gradient with respect to w
86 - db: Gradient with respect to b
87 """
88 dx, dw, db = None, None, None
89
90 N, F, out_height, out_width = dout.shape
91 x, w, b, conv_param = cache
92
93 stride, pad = [conv_param['stride'], conv_param['pad']]
94 xpad = np.pad(x, ((0,0), (0,0), (pad,pad), (pad,pad)), mode='constant')
95 num_filts, _, f_height, f_width = w.shape
96
97 # ===== #
98 # YOUR CODE HERE:
99 # Implement the backward pass of a convolutional neural network.
100 # Calculate the gradients: dx, dw, and db.
101 # ===== #
102
103 _, _, H, W = x.shape # [N, 3, 32, 32]
104 dx_temp = np.zeros_like(xpad) # initial to all zeros
105 dw = np.zeros_like(w)
106 db = np.zeros_like(b)
107
108 # Calculate dB.
109 for kernal in range(F):
110     db[kernal] += np.sum(dout[:, kernal, :, :]) # sum all N img's kernal ->
[32, 32], then sum all 32x32 elements -> 1 scalar
111
112 # Calculate dw.
113 for img in range(N): # for each image
114     for kernal in range(F): # for each kernal
115         for row in range(out_height): # from top to bottom
116             for col in range(out_width): # from left to right
117                 dw[kernal, ...] += dout[img, kernal, row, col] * xpad[img, :,
row*stride:row*stride+f_height, col*stride:col*stride+f_width]
118
119 # Calculate dx.
120 for img in range(N): # for each image
121     for kernal in range(F): # for each kernal
122         for row in range(out_height): # from top to bottom
123             for col in range(out_width): # from left to right
124                 dx_temp[img, :, row*stride:row*stride+f_height,
col*stride:col*stride+f_width] += dout[img, kernal, row, col] * w[kernal, ...]
125
126 dx = dx_temp[:, :, pad:H+pad, pad:W+pad]
127
128 # ===== #
129 # END YOUR CODE HERE
130 # ===== #
131
132 return dx, dw, db
133
134
135 def max_pool_forward_naive(x, pool_param):
136     """
137     A naive implementation of the forward pass for a max pooling layer.
138
139     Inputs:
140     - x: Input data, of shape (N, C, H, W)
141     - pool_param: dictionary with the following keys:
142         - 'pool_height': The height of each pooling region
143         - 'pool_width': The width of each pooling region
144         - 'stride': The distance between adjacent pooling regions
145
146     Returns a tuple of:
147     - out: Output data
148     - cache: (x, pool_param)
149     """
150     out = None
151
152     # ===== #
153     # YOUR CODE HERE:
154     # Implement the max pooling forward pass.
155     # ===== #
156
157     pool_height = pool_param.get('pool_height')
158     pool_width = pool_param.get('pool_width')

```

```

159 stride = pool_param.get('stride')
160 N, C, H, W = x.shape # [N, 3, 32, 32]
161
162 out_height = np.int(((H - pool_height) / stride) + 1) # calculate output
height
163 out_width = np.int(((W - pool_width) / stride) + 1) # calculate output
width
164 out = np.zeros([N, C, out_height, out_width])
165
166 for img in range(N): # for each image
167     for channel in range(C): # for each channel
168         for row in range(out_height): # from top to bottom
169             for col in range(out_width): # from left to right
170                 out[img, channel, row, col] = np.max(x[img, channel,
row*stride:row*stride+pool_height, col*stride:col*stride+pool_width])
171
172 # ===== #
173 # END YOUR CODE HERE
174 # ===== #
175 cache = (x, pool_param)
176 return out, cache
177
178 def max_pool_backward_naive(dout, cache):
179     """
180     A naive implementation of the backward pass for a max pooling layer.
181
182     Inputs:
183     - dout: Upstream derivatives
184     - cache: A tuple of (x, pool_param) as in the forward pass.
185
186     Returns:
187     - dx: Gradient with respect to x
188     """
189     dx = None
190     x, pool_param = cache
191     pool_height, pool_width, stride = pool_param['pool_height'],
pool_param['pool_width'], pool_param['stride']
192
193 # ===== #
194 # YOUR CODE HERE:
195 # Implement the max pooling backward pass.
196 # ===== #
197
198 N, C, H, W = x.shape # [N, 3, 32, 32]
199 _, _, dout_height, dout_width = dout.shape
200 dx = np.zeros_like(x)
201
202 for img in range(N): # for each image
203     for channel in range(C): # for each channel
204         for row in range(dout_height): # from top to bottom
205             for col in range(dout_width): # from left to right
206                 max_idx = np.argmax(x[img, channel,
row*stride:row*stride+pool_height, col*stride:col*stride+pool_width])
207                 max_position = np.unravel_index(max_idx, [pool_height, pool_width])
208                 dx[img, channel, row*stride:row*stride+pool_height,
col*stride:col*stride+pool_width][max_position] = dout[img, channel, row,
col]
209
210 # ===== #
211 # END YOUR CODE HERE
212 # ===== #
213
214 return dx
215
216 def spatial_batchnorm_forward(x, gamma, beta, bn_param):
217     """
218     Computes the forward pass for spatial batch normalization.
219
220     Inputs:
221     - x: Input data of shape (N, C, H, W)
222     - gamma: Scale parameter, of shape (C,)
223     - beta: Shift parameter, of shape (C,)
224     - bn_param: Dictionary with the following keys:
225         - mode: 'train' or 'test'; required
226         - eps: Constant for numeric stability
227         - momentum: Constant for running mean / variance. momentum=0 means that
228           old information is discarded completely at every time step, while
229           momentum=1 means that new information is never incorporated. The
230           default of momentum=0.9 should work well in most situations.
231         - running_mean: Array of shape (D,) giving running mean of features
232         - running_var: Array of shape (D,) giving running variance of features
233
234     Returns a tuple of:
235     - out: Output data, of shape (N, C, H, W)

```

```

236 - cache: Values needed for the backward pass
237 """
238 out, cache = None, None
239
240 # ===== #
241 # YOUR CODE HERE:
242 # Implement the spatial batchnorm forward pass.
243 #
244 # You may find it useful to use the batchnorm forward pass you
245 # implemented in HW #4.
246 # ===== #
247
248 N, C, H, W = x.shape # [N, 3, 32, 32]
249 x_transpose = x.transpose(0, 2, 3, 1)
250 x_reshape = np.reshape(x_transpose, (N*H*W, C)) # reshape to 2D to do
batchnorm
251 out_2d, cache = batchnorm_forward(x_reshape, gamma, beta, bn_param)
252 out = out_2d.reshape((N, H, W, C)).transpose(0, 3, 1, 2) # reshape back
253
254 # ===== #
255 # END YOUR CODE HERE
256 # ===== #
257
258 return out, cache
259
260
261 def spatial_batchnorm_backward(dout, cache):
262     """
263     Computes the backward pass for spatial batch normalization.
264
265     Inputs:
266     - dout: Upstream derivatives, of shape (N, C, H, W)
267     - cache: Values from the forward pass
268
269     Returns a tuple of:
270     - dx: Gradient with respect to inputs, of shape (N, C, H, W)
271     - dgamma: Gradient with respect to scale parameter, of shape (C,)
272     - dbeta: Gradient with respect to shift parameter, of shape (C,)
273     """
274     dx, dgamma, dbeta = None, None, None
275
276     # ===== #
277     # YOUR CODE HERE:
278     # Implement the spatial batchnorm backward pass.
279     #
280     # You may find it useful to use the batchnorm forward pass you
281     # implemented in HW #4.
282     # ===== #
283
284     dx = np.zeros_like(dout)
285     N, C, H, W = dout.shape
286     dout_transpose = dout.transpose((0, 2, 3, 1))
287     dout_reshape = np.reshape(dout_transpose, (N*H*W, C)) # reshape to 2D to do
batchnorm
288     dx_2d, dgamma, dbeta = batchnorm_backward(dout_reshape, cache)
289     dx = dx_2d.reshape((N, H, W, C)).transpose(0, 3, 1, 2) # reshape back
290
291     # ===== #
292     # END YOUR CODE HERE
293     # ===== #
294
295     return dx, dgamma, dbeta

```

```

1 import numpy as np
2
3 from nndl.layers import *
4 from nndl.conv_layers import *
5 from cs231n.fast_layers import *
6 from nndl.layer_utils import *
7 from nndl.conv_layer_utils import *
8
9 import pdb
10
11 """
12 This code was originally written for CS 231n at Stanford University
13 (cs231n.stanford.edu). It has been modified in various areas for use in the
14 ECE 239AS class at UCLA. This includes the descriptions of what code to
15 implement as well as some slight potential changes in variable names to be
16 consistent with class nomenclature. We thank Justin Johnson & Serena Yeung
17 for
18 permission to use this code. To see the original version, please visit
19 cs231n.stanford.edu.
20 """
21
22 class ThreeLayerConvNet(object):
23     """
24     A three-layer convolutional network with the following architecture:
25
26     conv - relu - 2x2 max pool - affine - relu - affine - softmax
27
28     The network operates on minibatches of data that have shape (N, C, H, W)
29     consisting of N images, each with height H and width W and with C input
30     channels.
31     """
32     def __init__(self, input_dim=(3, 32, 32), num_filters=32, filter_size=7,
33                 hidden_dim=100, num_classes=10, weight_scale=1e-3, reg=0.0,
34                 dtype=np.float32, use_batchnorm=False):
35         """
36         Initialize a new network.
37
38         Inputs:
39         - input_dim: Tuple (C, H, W) giving size of input data
40         - num_filters: Number of filters to use in the convolutional layer
41         - filter_size: Size of filters to use in the convolutional layer
42         - hidden_dim: Number of units to use in the fully-connected hidden layer
43         - num_classes: Number of scores to produce from the final affine layer.
44         - weight_scale: Scalar giving standard deviation for random
45         initialization
46         of weights.
47         - reg: Scalar giving L2 regularization strength
48         - dtype: numpy datatype to use for computation.
49         """
50         self.use_batchnorm = use_batchnorm
51         self.params = {}
52         self.reg = reg
53         self.dtype = dtype
54
55         # ===== #
56         # YOUR CODE HERE:
57         # Initialize the weights and biases of a three layer CNN. To
58         initialize:
59         # - the biases should be initialized to zeros.
60         # - the weights should be initialized to a matrix with entries
61         #   drawn from a Gaussian distribution with zero mean and
62         #   standard deviation given by weight_scale.
63         # ===== #
64
65         C, H, W = input_dim
66         ## Initial 1st layer (conv): [32, 3, 7, 7]
67         ## there are 32 kernels in this layer, and each kernel has 3 dimensions
68         ## and each kernel is 7x7
69         stride = 1
70         pad = (filter_size - 1) / 2
71         out_conv_height = (H + 2 * pad - filter_size) / stride + 1
72         out_conv_width = (W + 2 * pad - filter_size) / stride + 1
73         self.params['W1'] = np.random.normal(0, weight_scale, [num_filters, C,
74 filter_size, filter_size])
75         self.params['b1'] = np.zeros([num_filters]) # each filter (kernel) has a
76         bias
77
78         ## Initial 2nd layer (fc): after conv and max pooling, the weight and
79         height are half, and channel changes from 3 to 32
80         ## so for fully connected layer, the shape would be [N, 16x16x32].

```

```

77     ## First, we will flatten the graph after conv from [32, 16, 16] to
[32x16x16],
78     ## then, do fc and reduce the dimension from 32x16x16 to 100
79     out_pool_height = int((out_conv_height - 2) / 2 + 1)
80     out_pool_width = int((out_conv_width - 2) / 2 + 1)
81     self.params['W2'] = np.random.normal(0, weight_scale,
[out_pool_height*out_pool_width*num_filters, hidden_dim])
82     self.params['b2'] = np.zeros([hidden_dim])
83
84     ## Initial 3rd layer (fc): keep reducing the dimension from 100 to 10
(cifar dataset has 10 classes)
85     self.params['W3'] = np.random.normal(0, weight_scale, [hidden_dim,
num_classes])
86     self.params['b3'] = np.zeros([num_classes])
87
88     # ===== #
89     # END YOUR CODE HERE
90     # ===== #
91
92     for k, v in self.params.items():
93         self.params[k] = v.astype(dtype)
94
95
96 def loss(self, X, y=None):
97     """
98     Evaluate loss and gradient for the three-layer convolutional network.
99
100    Input / output: Same API as TwoLayerNet in fc_net.py.
101    """
102    W1, b1 = self.params['W1'], self.params['b1']
103    W2, b2 = self.params['W2'], self.params['b2']
104    W3, b3 = self.params['W3'], self.params['b3']
105
106    # pass conv_param to the forward pass for the convolutional layer
107    filter_size = W1.shape[2]
108    conv_param = {'stride': 1, 'pad': (filter_size - 1) / 2}
109
110    # pass pool_param to the forward pass for the max-pooling layer
111    pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}
112
113    scores = None
114
115    # ===== #
116    # YOUR CODE HERE:
117    # Implement the forward pass of the three layer CNN. Store the output
118    # scores as the variable "scores".
119    # ===== #
120
121    ## forward: [N, 3, 32, 32] -> [N, 10]
122    layer1_out, combined_cache = conv_relu_pool_forward(X, W1, b1,
conv_param, pool_param) # 1st layer: conv + relu + 2x2 max pool
123    fc1_out, fc1_cache = affine_relu_forward(layer1_out, W2, b2) # 2nd layer:
fc and relu
124    scores, fc2_cache = affine_forward(fc1_out, W3, b3) # 3rd layer: fc
125
126    # ===== #
127    # END YOUR CODE HERE
128    # ===== #
129
130    if y is None:
131        return scores
132
133    loss, grads = 0, {}
134    # ===== #
135    # YOUR CODE HERE:
136    # Implement the backward pass of the three layer CNN. Store the grads
137    # in the grads dictionary, exactly as before (i.e., the gradient of
138    # self.params[k] will be grads[k]). Store the loss as "loss", and
139    # don't forget to add regularization on ALL weight matrices.
140    # ===== #
141
142    loss, dscores = softmax_loss(scores, y) # 3rd layer: softmax
143    loss += self.reg * 0.5 * (np.sum(np.square(W1)) + np.sum(np.square(W2)) +
np.sum(np.square(W3)))
144
145    ## backward: [N, 10] -> [N, 3, 32, 32]
146    dx3, dw3, db3 = affine_backward(dscores, fc2_cache) # oppo 3rd layer: fc
147    dx2, dw2, db2 = affine_relu_backward(dx3, fc1_cache) # oppo 2nd layer: fc
148    dx1, dw1, db1 = conv_relu_pool_backward(dx2, combined_cache) # oppo 1st
layer: conv + relu + pool
149
150    grads['W3'], grads['b3'] = dw3 + self.reg * W3, db3
151    grads['W2'], grads['b2'] = dw2 + self.reg * W2, db2
152    grads['W1'], grads['b1'] = dw1 + self.reg * W1, db1

```

```
153
154 # ===== #
155 # END YOUR CODE HERE
156 # ===== #
157
158     return loss, grads
159
160
161 pass
162
```