

Spatial batch normalization 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 $(NH \times W, C)$ 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 `nn1/` 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 nn1.conv_layers import *  
from utils.data_utils import get_CIFAR10_data  
from utils.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array  
from utils.solver import Solver  
  
%matplotlib inline  
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots  
plt.rcParams['image.interpolation'] = 'nearest'  
plt.rcParams['image.cmap'] = 'gray'  
  
# for auto-reloading external modules  
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython  
%load_ext autoreload  
%autoreload 2  
  
def rel_error(x, y):  
    """returns relative error"""  
    return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

Spatial batch normalization forward pass

Implement the forward pass, `spatial_batchnorm_forward` in `nn1/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: [10.31759087  9.2986021  9.87311898]  
Stds: [3.26801225  3.91439358  4.01802861]  
After spatial batch normalization:  
Shape: (2, 3, 4, 5)  
Means: [-6.10622664e-16  5.82867089e-17  1.44328993e-16]  
Stds: [0.99999553  0.99999867  0.99999869]  
After spatial batch normalization (nontrivial gamma, beta):  
Shape: (2, 3, 4, 5)  
Means: [6.  7.  8.]  
Stds: [2.9999986  3.99999869  4.99999845]
```

Spatial batch normalization backward pass

Implement the backward pass, `spatial_batchnorm_backward` in `nn1/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 x: spatial_batchnorm_forward(x, gamma, beta, bn_param)[0]  
fb = lambda b: spatial_batchnorm_backward(x, gamma, beta, bn_param)[0]  
  
dx_num = eval_numerical_gradient_array(fx, x, dout)  
dx_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(dgamma, dgamma))  
print('dbeta error: ', rel_error(db_num, dbeta))  
  
dx error: 2.519303783801587e-07  
dgamma error: 1.75594065756475e-11  
dbeta error: 3.279603776544007e-12
```

NNDL py files


```
## CNN.py

import numpy as np
from math import exp

from ndml.layers import *
from ndml.conv_layers import *
from ndml.fast_layers import *
from ndml.layer_utils import *
from ndml.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_dim1=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_dim1: 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
        - reg: Regularizer strength
        - dtype: numpy datatype to use for computation.
        - use_batchnorm: If True, use batch normalization.
        """
        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 with zero with entries
        #     drawn from a Gaussian distribution with a mean and
        #     standard deviation given by weight scale.
        # =====
        # END YOUR CODE HERE

        C, H, W = input_dim

        size_W1 = num_filters, C, filter_size, filter_size
        size_b1 = num_filters

        Cnn_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".
        # =====
        # END YOUR CODE HERE

        layer1_out, combined_cache = conv_relu_pool_forward(X, W1, b1, conv_param, pool_param)
        fcl_out, fcl_cache = affine_relu_forward(layer1_out, W2, b2)
        scores, fc2_cache = affine_forward(fcl_out, 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 "loss", the
        #   self.params[k] will be grads[k]), store the loss as "loss", and
        #   don't forget to add regularization on ALL weight matrices.
        # =====
        # END YOUR CODE HERE
        # =====

        loss, dscores = softmax_loss(scores, y)
        loss += self.reg * 0.5 * (np.sum(np.square(W1)) + np.sum(np.square(W2)) + np.sum(np.square(W3)))

        dx2, dx3, db3 = affine_backward(dscores, fc2_cache)
        dx1, db2, db1 = affine_relu_backward(dx2, fcl_cache)
        dx1, db1, dx0 = conv_relu_pool_backward(dx1, combined_cache)

        grads['W2'] = grads['b2'] = dx2 + self.reg * W2, db2
        grads['W3'] = grads['b3'] = dx3 + self.reg * W3, db3
        grads['W1'] = grads['b1'] = dx1 + self.reg * W1, db1

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

        return loss, grads

# Conv_layers.py

import numpy as np
from ndml.layers import *
import pdb

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

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

    Inputs:
    - 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': This number of pixels that will be used to zero-pad the input.

    Returns a tuple of:
    - out: Output data, of shape (N, F, H', W') where H' and W' are given by
      H' = 1 + (H + 2 * pad - HH) / stride
      W' = 1 + (W + 2 * pad - WW) / stride
    - cache: (x, w, b, conv_param)

    out = None
    pad = conv_param['pad']
    stride = conv_param['stride']

    # =====
    # YOUR CODE HERE:
    # - Implement the forward pass of a convolutional neural network.
    # - Store the output as 'out'.
    # - Hint: to pad the array, you can use the function np.pad.
    # =====
    # END YOUR CODE HERE
    # =====

    N, C, H, W = x.shape
    F, C, HH, WW = w.shape

    padded_x = np.pad(x, [(0, 0), (0, 0), (pad, pad), (pad, pad)], 'constant')
    out_height = np.int((H + 2 * pad - HH) / stride + 1)
    out_width = np.int((W + 2 * pad - WW) / stride + 1)
    out = np.zeros((N, F, out_height, out_width))

    for img_idx in range(N):
        for row_idx in range(F):
            for col_idx in range(out_width):
                out[img_idx, row_idx, col_idx] = np.sum(w[kernal, ...] * \
                    padded_x[img_idx, row*stride:row*stride+HH, col*stride:col*stride+WW])

    # =====
    # 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']
    dxpad, dx = conv_backward_naive(dout, cache)
    num_filters, _, 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.
    # =====
    # END YOUR CODE HERE
    # =====

    N, C, H, W = x.shape
    F, C, HH, WW = w.shape

    dx_temp = np.zeros_like(xpad) # initial to all zeros
    dx2 = np.zeros_like(w)
    db = np.zeros_like(b)

    # Calculate db.
    for kernal_idx in range(F):
        dx[kernal_idx, :, :, :] = np.sum(dout[:, kernal_idx, :, :]) # sum all N img's kernal -> [32, 32], then sum all 32x32 elem

    # Calculate dw.
    for img_idx in range(N):
        for kernal_idx in range(F):
            for row_idx in range(out_height):
                for col_idx in range(out_width):
                    dx[kernal_idx, row_idx, col_idx] = np.sum(dout[img_idx, kernal_idx, row_idx, col_idx] * xpad[img_idx, row_idx, col_idx])

    # Calculate db.
    for img_idx in range(N):
        for kernal_idx in range(F):
            for row_idx in range(out_height):
                for col_idx in range(out_width):
                    dx[kernal_idx, row_idx, col_idx] = np.sum(dout[img_idx, kernal_idx, row_idx, col_idx] * xpad[img_idx, row_idx, col_idx])

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

    pool_height = pool_param.get('pool_height')
    pool_width = pool_param.get('pool_width')
    stride = pool_param.get('stride')

    N, C, H, W = x.shape

    out_height = np.int((H - pool_height) / stride + 1)
    out_width = np.int((W - pool_width) / stride + 1)
    out = np.zeros((N, C, out_height, out_width))

    for img_idx in range(N):
        for channel_idx in range(C):
            for row_idx in range(out_height):
                for col_idx in range(out_width):
                    out[img_idx, channel_idx, row_idx, col_idx] = np.max(x[img_idx, channel_idx, row*stride:row*stride+pool_height, col*stride:col*stride+pool_width])

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

    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

    N, C, H, W = x.shape
    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.
    # =====
    # END YOUR CODE HERE
    # =====

    N, C, H, W = x.shape
    N, C, H, W = dout.shape

    dx = np.zeros_like(x)
    for img_idx in range(N):
        for channel_idx in range(C):
            for row_idx in range(out_height):
                for col_idx in range(out_width):
                    dx[img_idx, channel_idx, row_idx, col_idx] = np.sum(dout[img_idx, channel_idx, row_idx, col_idx] * x[img_idx, channel_idx, row_idx, col_idx])

    # =====
    # 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: Spatial parameters, of shape (C,)
    - beta: Shift parameters, 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.
    # =====
    # END YOUR CODE HERE
    # =====

    N, C, H, W = x.shape
    N, C, H, W = x.transpose(0, 2, 3, 1)
    x_reshape = x.transpose(0, 2, 3, 1)
    dx2, dx = batchnorm_forward(x_reshape, gamma, beta, bn_param)
    out_2d, cache = batchnorm_backward(x_reshape, gamma, beta, bn_param)
    out = dx2_2d.reshape((N, H, W, C)).transpose(0, 3, 1, 2) # reshape back

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

    N, C, H, W = dout.shape
    N, C, H, W = dout.transpose(0, 2, 3, 1)
    dout_reshape = dout.transpose(0, 2, 3, 1)
    dx2, dx = batchnorm_backward(dout_reshape, gamma, beta, bn_param)
    dx = dx2_2d.reshape((N, H, W, C)).transpose(0, 3, 1, 2) # reshape back

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

    return dx, dgamma, dbeta

# Layers.py

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 dimension D = 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.
    # =====
    # END YOUR CODE HERE
    # =====

    x_reshape = x.reshape((x.shape[0], x.shape[0])) # N * D
    out = np.dot(x_reshape, w) + b
    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, d_1, ..., d_k)
    - dw: Gradient with respect to w, of shape (D, M)
    - db: Gradient with respect to b, of shape (M,)
    """
    dx, dw, db = None, None, None

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

    x_reshape = np.reshape(x, (x.shape[0], x.shape[0]))
    dx_reshape = np.dot(dout, x_reshape)
    dx = dx_reshape.reshape(x.shape) # N * D
    dw = x_reshape.T.dot(dout) # D * M
    db = dout.T.dot(np.ones(x.shape[0])) # M * 1

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

    return dx, dw, db

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

    Inputs:
    - 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.
    # =====
    # END YOUR CODE HERE
    # =====

    out = np.maximum(0, x)
    cache = x
    return out, cache

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

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

    Returns:
    - dx: Gradient with respect to x
    """
    # =====
    # YOUR CODE HERE:
    # - Implement the ReLU backward pass
    # =====
    # END YOUR CODE HERE
    # =====

    dx = (x > 0) * (dout)

    # =====
    # 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 time step 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 torch implementation
    of batch normalization also uses running averages.

    Inputs:
    - x: Input 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: Output 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:
        # - Calculate the sample mean and variance of the minibatch.
        # - (1) Calculate the running mean and variance of the minibatch.
        # - (2) Normalize the activations with the running mean and variance.
        # - (3) Store and normalize the normalized activations. Store this
        #   as the variable 'out'.
        # - (4) Store any variables you may need for the backward pass in
        #   the 'cache' variable.
        # =====
        # END YOUR CODE HERE
        # =====

        minibatch_mean = np.mean(x, axis=0)
        minibatch_var = np.var(x, axis=0)
        x_normalized = (x - minibatch_mean) / np.sqrt(minibatch_var + eps)
        running_mean = momentum * running_mean + (1 - momentum) * minibatch_mean
        running_var = momentum * running_var + (1 - momentum) * minibatch_var
        bn_param['running_mean'] = running_mean
        bn_param['running_var'] = running_var

        cache = [
            'minibatch_var', minibatch_var,
            'x_normalized', x_normalized,
            'gamma', gamma,
            'eps', eps
        ]

        # =====
        # 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'.
        # =====
        # END YOUR CODE HERE
        # =====

        out = gamma * (x - running_mean) / np.sqrt(running_var + eps) + beta

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

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

    N = dout.shape[0]
    minibatch_var = cache.get('minibatch_var')
    x_normalized = cache.get('x_normalized')
    gamma = cache.get('gamma')
    eps = cache.get('eps')

    # Calculate dx
    dxdt = dout * gamma
    dxdt_var = np.sqrt(minibatch_var + eps)
    dxdt_var = np.sqrt(minibatch_var + eps)
    dxdt_var = dxdt_var * x_normalized, axis=0 / (sqrt_var**2)
    dxdt = 2 * x_normalized * dxdt_var * np.ones_like(dout) / N
    dx1 = dxdt + dxdt2
    dx2 = np.sum(dx1, axis=0) * np.ones_like(dout) / N
    dx = dx1 + dx2

    # Calculate dbeta and dgamma
    dbeta = np.sum(dout, axis=0)
    dgamma = np.sum(dout * x_normalized, axis=0)

    # =====
    # YOUR CODE HERE:
    # - Implement the batchnorm backward pass, calculating dx, dgamma, and dbeta.
    # =====
    # 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 probability. We should keep neuron output with probability p.
      - mode: 'train' or 'test'; if the mode is train, then perform dropout;
      - seed: seed for the random number generator. Passing seed makes this
        function deterministic, which is needed for gradient checking but not in
        real networks.

    Returns:
    - out: Array of the same shape as x
    - mask: A tuple of (mask, np.random.seed())
    """
    p = dropout_param['p']
    mode = dropout_param['mode']
    seed = dropout_param['seed']
    np.random.seed(dropout_param['seed'])

    mask = 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'.
        # =====
        # END YOUR CODE HERE
        # =====

        mask = np.random.rand(x.shape[0] * p) / (1 - p)
        out = x * mask

        # =====
        # YOUR CODE HERE:
        # - Implement the inverted dropout forward pass during test time.
        # =====
        # END YOUR CODE HERE
        # =====

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

    return out, cache

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

    Inputs:
    - dout: Upstream derivatives, of any shape
    - cache: (dropout_param, mask) from dropout_forward.
    - dropout_param: A dictionary with the following keys:
      - p: Dropout probability. We should keep neuron output with probability p.
      - mode: 'train' or 'test'; if the mode is train, then perform dropout;
      - seed: seed for the random number generator. Passing seed makes this
        function deterministic, which is needed for gradient checking but not in
        real networks.

    Returns:
    - dx: Array of the same shape as x
    """
    # =====
    # YOUR CODE HERE:
    # - Implement the inverted dropout backward pass during training time.
    # =====
    # END YOUR CODE HERE
    # =====

    dx = dout * mask

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

    return dx

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

    Inputs:
    - x: Input data, of shape (N, C) where x[i, :] 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.sum(margins, axis=1)
    loss = np.sum(margins) / N
    dx = np.zeros_like(x)
    dx[margins > 0] = 1
    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, :] 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 = probs / np.sum(probs, axis=1, keepdims=True)
    loss = -np.sum(np.log(probs[np.arange(N), y])) / N
    dx = probs.copy()
    dx[np.arange(N), y] -= 1
    return loss, dx
```


Convolutional neural network layers

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

```
In [2]: ## Import and setup

import time

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

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

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-python
%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 `nnutils/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 `nnutils/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 [3]: x_shape = (2, 3, 4, 4)
w_shape = (3, 3, 4, 4)
w = np.linspace(-0.1, 0.5, numnp.prod(x_shape)).reshape(x_shape)
w = np.linspace(-0.2, 0.3, numnp.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([[[[-0.0875909, -0.0987781],
                        [-0.18387192, -0.2109216 ]],
                        [[ 0.21027089,  0.21610197],
                        [ 0.22847626,  0.23004637]],
                        [[ 0.50813986,  0.54309974],
                        [ 0.64082444,  0.67101435]]],
                        [[[-0.98053599, -1.0314344],
                        [-1.19128892, -1.24695841]],
                        [[ 0.69108355,  0.66880383],
                        [ 0.59480972,  0.56776003]],
                        [[ 2.36272098,  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 `nnutils/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, 2)
dout = np.random.randn(4, 2, 3, 5)
pool_param = {'stride': 1, 'pad': 1}
out, cache = conv_backward_naive(x, w, b, conv_param)

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

out, cache = conv_backward_naive(dout, cache)
dx, dw, db = conv_backward_naive(dout, cache)

# Your errors should be around 1e-7
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: 3.767230441055007e-09
dw error: 4.433983702039593e-10
db error: 1.2560968119626294e-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 `nnutils/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, numnp.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]],
                        [[-0.14926316, -0.13952662],
                        [-0.08631579, -0.07157895]],
                        [[-0.02736842, -0.01263158],
                        [ 0.03157895,  0.04631579]],
                        [[ 0.09052632,  0.10526316],
                        [ 0.14947369,  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: 4.166665157267834e-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 `nnutils/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_param)[0], x, dout)

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

# Your error should be around 1e-12
print('Testing max_pool_backward_naive function:')
print('dx error: ', rel_error(dx, dx_num))

Testing max_pool_backward_naive function:
dx error: 3.2756311759753567e-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()
dx_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('\nTesting 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: 3.308228s
Fast: 0.00954s
Speedup: 346.520341x
Difference: 3.630933663400736e-11

Testing conv_backward_fast:
Naive: 4.533300s
Fast: 0.00824s
Speedup: 598.633300x
dx difference: 8.00848880409912e-12
dw difference: 1.4601170835547198e-12
db difference: 1.775874036379349e-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: 0.27421s
Fast: 0.003771s
Speedup: 72.697901x
Difference: 0.0

Testing pool_backward_fast:
Naive: 0.338737s
Speedup: 40.866565x
dx difference: 0.0
```

Implementation of cascaded layers

We've provided the following functions in `nnutils/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 [9]: from nnutils.conv_layer_utils import conv_relu_pool_forward, conv_relu_pool_backward

x = np.random.randn(2, 3, 16, 16)
w = np.random.randn(3, 3, 3, 3)
b = np.random.randn(3,)
dout = np.random.randn(2, 3, 8, 8)
conv_param = {'stride': 1, 'pad': 1}
pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}

out, cache = conv_relu_pool_forward(x, w, b, conv_param, pool_param)
dx, dw, db = conv_relu_pool_backward(dout, cache)

dx_num = eval_numerical_gradient_array(lambda x: conv_relu_pool_forward(x, w, b, conv_param, pool_param)[0], x,
dw_num = eval_numerical_gradient_array(lambda w: conv_relu_pool_forward(x, w, b, conv_param)[0], w,
db_num = eval_numerical_gradient_array(lambda b: conv_relu_pool_forward(x, w, b, conv_param, pool_param)[0], 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: 2.648973494450611e-08
dw error: 3.6802080481655105e-10
db error: 2.3107608188380604e-11
```

```
In [10]: from nnutils.conv_layer_utils import conv_relu_forward, conv_relu_backward

x = np.random.randn(2, 3, 8, 8)
w = np.random.randn(3, 3, 3, 3)
b = np.random.randn(3,)
dout = np.random.randn(2, 3, 8, 8)
conv_param = {'stride': 1, 'pad': 1}

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

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

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

Testing conv_relu:
dx error: 3.0716994466208765e-09
dw error: 4.736204398019599e-10
db error: 2.982569184594821e-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.

NNDL py files

Convolutional neural networks

In this notebook, we'll put together our convolutional layers and pooling layers. Our model 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

layers in HW #4, you may collaborate with a classmate to use their layer implementations from HW #4. You may also visit the TA or 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.

our solutions.

- ```
import numpy as np
import matplotlib.pyplot as plt
from nnrl.cnn import *
from utils.data_utils import get_CIFAR10_data
from utils.gradient_check import eval_numerical_gradient_array, eval_numerical_gradients
from nnrl.layers import *
from nnrl.conv_layers import *
```

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

```
for auto-reloading external modules
see: https://stackoverflow.com/questions/1907993/autoreload-of-modules-in-python
%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 nn1/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 different 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 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,
 dropout_rate=0.4)
```

```

for param_name in sorted(grades):
 f = lambda _: model.loss(X, y[0])
 param_grad_num = eval('numerical_gradient(f, model.params[param_name], verbose=False, h=1e-6)')
 e = rel_error(param_grad_num, grads[param_name])
 print('({}) max relative error: {}'.format(param_name, rel_error(param_grad_num, grads[param_name])))

W1 max relative error: 0.003922954350389366
W2 max relative error: 0.00278989632466564695
b1 max relative error: 2.9185898998921736e=05
b2 max relative error: 1.076309395490432e=05
b3 max relative error: 2.229697351371742e=07
b4 max relative error: 6.23692469489459e=10

```

## 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_rules='adam',
 optim_config={
 'learning_rate': 1e-3,
 },
 show_images=True, print_every=100)

```

```
Solver.train()

(iteration 1 / 20) loss: 2.325183
(Epoch 0 / 10) train acc: 0.2100000/ val_acc: 0.1280000
(iteration 2 / 20) loss: 2.895399
(Epoch 1 / 10) train acc: 0.2100000/ val_acc: 0.1060000
(iteration 3 / 20) loss: 2.679631
(iteration 4 / 20) loss: 2.547447
(Epoch 2 / 10) train acc: 0.2800000/ val_acc: 0.1340000
(iteration 5 / 20) loss: 2.490574
(iteration 6 / 20) loss: 1.855074
(Epoch 3 / 10) train acc: 0.3300000/ val_acc: 0.2070000
(iteration 7 / 20) loss: 1.787314
(iteration 8 / 20) loss: 1.830555
(Epoch 4 / 10) train acc: 0.5800000/ val_acc: 0.1880000
(iteration 9 / 20) loss: 1.225561
(iteration 10 / 20) loss: 1.211132
(Epoch 5 / 10) train acc: 0.5900000/ val_acc: 0.1600000
(iteration 11 / 20) loss: 1.245911
(iteration 12 / 20) loss: 1.075091
(Epoch 6 / 10) train acc: 0.6200000/ val_acc: 0.2160000
(iteration 13 / 20) loss: 1.038730
(iteration 14 / 20) loss: 1.117598
(Epoch 7 / 10) train acc: 0.6200000/ val_acc: 0.2330000
(iteration 15 / 20) loss: 0.750677
(iteration 16 / 20) loss: 0.685618
(Epoch 8 / 10) train acc: 0.6600000/ val_acc: 0.2280000
(iteration 17 / 20) loss: 0.566709
(iteration 18 / 20) loss: 0.590791
(Epoch 9 / 10) train acc: 0.6900000/ val_acc: 0.2120000
(iteration 19 / 20) loss: 0.537815
(iteration 20 / 20) loss: 0.680158
(Epoch 10 / 10) train acc: 0.9000000/ val_acc: 0.2000000
```

```
plt.subplot(2, 1, 1)
```

```
plt.plot(solver.train_acc_history, '-o')
plt.plot(solver.val_acc_history, '-o')
```

```
plt.xlabel('epoch')
plt.ylabel('accuracy')
plt.show()
```



```
(Iteration 741 / 980)
(Iteration 761 / 980)
(Iteration 781 / 980)
```

```
(Iteration 801 / 980) loss: 1.917589
(iteration 821 / 980) loss: 1.501216
```

## 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 learning 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 [10]:
=====
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,
 num_filters = 64,
 filter_size = 3)

solver = Solver(model, data,
 num_epochs=10, batch_size=500,
 update_rule='adam',
 optim_config=[
 'learning_rate': 1e-3,
 'lr_decay': 0.9,
 'verbose': True, print_every=15])

solver.train()

print out the validation and test accuracy
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 / 980) loss: 2.36750
(Epoch 0 / 10) train acc: 0.194000/ val_acc: 0.112000
(iteration 16 / 980) loss: 1.875961
(iteration 31 / 980) loss: 1.727335
(iteration 46 / 980) loss: 1.612002
(iteration 61 / 980) loss: 1.385337
(Epoch 1 / 1) train acc: 0.459000/ val_acc: 0.472000
```

```
(Iteration 76 / 980) loss: 1.309349
(Iteration 91 / 980) loss: 1.281215
(Epoch 1 / 10) train acc: 0.581000; val acc: 0.560000
(Iteration 106 / 980) loss: 1.282992
```

```
(Iteration 121 / 980) loss: 1.291973
```

```
(Iteration 138 / 980) 1c
(Iteration 151 / 980) 1c
(Iteration 166 / 980) 1c
(Iteration 181 / 980) 1c
```

- ```
(Iteration 184 / 980) train acc: 1.17236e
(Iteration 194 / 980) loss: 1.11847e
(Epoch 2 / 10) train acc: 0.620200; val_acc: 0.630300
(Iteration 214 / 980) loss: 1.06670e
(Iteration 226 / 980) loss: 1.08402e
(Iteration 244 / 980) loss: 1.08324e
(Iteration 256 / 980) loss: 1.11462e
(Iteration 274 / 980) loss: 1.01887e
(Iteration 286 / 980) loss: 1.01412e
(Epoch 3 / 10) train acc: 0.661000; val_acc: 0.631000
(Iteration 304 / 980) loss: 1.03939e
(Iteration 316 / 980) loss: 0.965458
(Iteration 334 / 980) loss: 0.99249e
(Iteration 346 / 980) loss: 0.98237e
(Iteration 364 / 980) loss: 0.904714
(Iteration 376 / 980) loss: 0.926157
(Iteration 394 / 980) loss: 0.904662
(Epoch 4 / 10) train acc: 0.726000; val_acc: 0.634000
(Iteration 406 / 980) loss: 0.950047
(Iteration 424 / 980) loss: 0.885359
(Iteration 436 / 980) loss: 0.863559
(Iteration 454 / 980) loss: 0.946079
(Iteration 466 / 980) loss: 0.789447
(Iteration 484 / 980) loss: 0.820704
(Epoch 5 / 10) train acc: 0.769000; val_acc: 0.640000
(Iteration 496 / 980) loss: 0.825014
```

```
{Iteration 52
{Iteration 54
```

```
(Iteration 556 / 980) loss: 0.704718
(Iteration 571 / 980) loss: 0.757800
(Iteration 586 / 980) loss: 0.783915
(Epoch 6 / 10) train_acc: 0.761000; val_acc: 0.653000
(Iteration 601 / 980) loss: 0.718778
(Iteration 616 / 980) loss: 0.746557
(Iteration 631 / 980) loss: 0.739426
(Iteration 646 / 980) loss: 0.729680
(Iteration 661 / 980) loss: 0.708834
(Iteration 676 / 980) loss: 0.851334
```

```

(iteration 691 / 980) loss: 0.657606
(iteration 706 / 980) loss: 0.658193
(iteration 721 / 980) loss: 0.660211
(iteration 736 / 980) loss: 0.663333
(iteration 751 / 980) loss: 0.630992
(iteration 766 / 980) loss: 0.613800
(iteration 781 / 980) loss: 0.610026
Epoch 8 / 10: train acc: 0.850000; val_acc: 0.672000
(iteration 796 / 980) loss: 0.684875
(iteration 811 / 980) loss: 0.645196
(iteration 826 / 980) loss: 0.555597
(iteration 842 / 980) loss: 0.593875
(iteration 856 / 980) loss: 0.594800
(iteration 871 / 980) loss: 0.525813
Epoch 9 / 10: train acc: 0.831000; val_acc: 0.666000
(iteration 886 / 980) loss: 0.531366
(iteration 901 / 980) loss: 0.499404
(iteration 916 / 980) loss: 0.545275
(iteration 931 / 980) loss: 0.495110
(iteration 946 / 980) loss: 0.449716
(iteration 961 / 980) loss: 0.481002
(iteration 976 / 980) loss: 0.520237
Epoch 10 / 10: train acc: 0.883000; val_acc: 0.676000
Validation set accuracy: 0.678
Test set accuracy: 0.661

```

CNN.py

```

In [ ]:
import numpy as np

from nnlib.layers import *
from nnlib.conv_layers import *
from utils.fast_layers import *
from nnlib.layer_utils import *
from nnlib.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

        conv_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: Output data

```

out = None

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

pool_height = pool_param.get('pool_height')
pool_width = pool_param.get('pool_width')
stride = pool_param.get('stride')
N, C, H, W = x.shape

out_height = np.int(((H - pool_height) / stride) + 1)
out_width = np.int(((W - pool_width) / stride) + 1)
out = np.zeros((N, C, out_height, out_width))

for img in range(N):
    for channel in range(C):
        for row in range(out_height):
            for col in range(out_width):
                out[img, channel, row, col] = np.max(x[img, channel, row*stride:row*stride+pool_height, col*stride:col*stride+pool_width])

# ===== #
# 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
    _, dout_height, dout_width = dout.shape
    dx = np.zeros_like(x)

    for img in range(N):
        for channel in range(C):
            for row in range(dout_height):
                for col in range(dout_width):
                    max_idx = np.argmax(x[img, channel, row*stride:row*stride+pool_height, col*stride:col*stride+pool_width])
                    max_position = np.unravel_index(max_idx, (pool_height, pool_width))
                    dx[img, channel, row*stride:row*stride+pool_height, col*stride:col*stride+pool_width][max_position] += dout[img, channel, row, col]

    # ===== #
    # 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 # (N, 3, 32, 32)
    x_transpose = x.transpose(0, 2, 3, 1)
    x_reshape = np.reshape(x_transpose, (NH*HW, C)) # reshape to 2D to do batchnorm
    out_2d, cache = batchnorm_forward(x_reshape, gamma, beta, bn_param)
    out = out_2d.reshape((N, H, W, C)).transpose(0, 3, 1, 2) # reshape back

    # ===== #
    # 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.
    # ===== #

    dx = np.zeros_like(dout)
    N, C, H, W = dout.shape
    dout_transpose = dout.transpose((0, 2, 3, 1))
    dout_reshape = np.reshape(dout_transpose, (NH*HW, C)) # reshape to 2D to do batchnorm
    dx_2d, dgamma, dbeta = batchnorm_backward(dout_reshape, cache)
    dx = dx_2d.reshape((H, H, W, C)).transpose(0, 3, 1, 2) # reshape back

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

```

```
return dx
```



```
import numpy as np
import pdb

def affine_forward(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_reshape = x.reshape((x.shape[0], w.shape[0])) # N * D
    out = x_reshape.dot(w) + b.reshape((1, b.shape[0])) # N * M

    # =====
    # 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, d_1, ..., 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_reshape = np.reshape(x, (x.shape[0], w.shape[0]))
    dx_reshape = dout.dot(w.T)
    dx = np.reshape(dx_reshape, x.shape) # N * D
    dw = x_reshape.T.dot(dout) # D * M
    db = dout.T.dot(np.ones(x.shape[0])) # M * 1

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

    return dx, dw, db

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

    Input:
    - x: Inputs, of any shape

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

    out = np.maximum(0, 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's).

    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 = (x > 0) * (dout)

    # =====
    # 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 each feature 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
      - 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=1e-5
    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.
        # =====

        minibatch_mean = np.mean(x, axis=0)
        minibatch_var = np.var(x, axis=0)
        x_normalize = (x - minibatch_mean) / np.sqrt(minibatch_var + eps)
        out = gamma * x_normalize + beta

        running_mean = momentum * running_mean + (1 - momentum) * minibatch_mean
        running_var = momentum * running_var + (1 - momentum) * minibatch_var
        bn_param['running_mean'] = running_mean
        bn_param['running_var'] = running_var

        cache = {
            'minibatch_var': minibatch_var,
            'x_normalize': (x - minibatch_mean),
            'x_normalize': x_normalize,
            'gamma': gamma,
            'eps': 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'.
        # =====

        out = gamma * (x - running_mean) / np.sqrt(running_var + eps) + beta

        # =====
        # 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.
    # =====

    N = dout.shape[0]
    minibatch_var = cache.get('minibatch_var')
    x_centralize = cache.get('x_centralize')
    x_normalize = cache.get('x_normalize')
    gamma = cache.get('gamma')
    eps = cache.get('eps')

    # calculate dx
    dxhat = dout * gamma
    dxhat = dxhat / np.sqrt(minibatch_var + eps)
    sqrt_var = np.sqrt(minibatch_var + eps)
    dvar = dvar = -np.sum(dxhat * x_centralize, axis=0) / (sqrt_var**2)
    dvar = dvar * 0.5 / sqrt_var
    dxm2 = 2 * x_centralize * dvar * np.ones_like(dout) / N
    dx1 = dxm2 + dxm2
    dx2 = np.sum(dx1, axis=0) * np.ones_like(dout) / N
    dx = dx1 * dx2

    # calculate dbeta and dgamma
    dbeta = np.sum(dout, axis=0)
    dgamma = np.sum(dout * x_normalize, axis=0)

    # =====
    # 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.random_sample(x.shape) >= p) / (1 - p)
        out = x * mask

        # =====
        # 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):
    """
    Performs 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 = dout * mask

        # =====
        # 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 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, 0] = 1
    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))
    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
```