

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 `nnl/` 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 nnl.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 `nnl/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 `nnl/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.layers import *
from ndml.layers import *
from ndml.layers 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 L weight
    channels.
    """
    def __init__(self, input_dim=(3, 32, 32), num_filters=32, filter_size=7,
                 hidden_dim=(10, 10), 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
        - 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 with zero with entries
        #   drawn from a Gaussian distribution with a 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

        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.

        Inputs:
        - X: Input data of shape (N, C, H, W)
        - 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
        - grads: Dictionary of gradients for each parameter.
        """
        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".
        # =====
        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

        # =====
        # 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 "grad" of
        # self.params[k] will be grads[k]). Store the loss as "loss", and
        # don't forget to add regularization on ALL weight matrices.
        # =====
        loss, 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)))

        dx1, dx2, dx3 = affine_backward(dscores, fc2_cache)
        dx1, dx2, dx3 = affine_relu_backward(dx1, fcl_cache)
        dx1, dx2, dx3 = conv_relu_pool_backward(dx1, combined_cache)

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

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

    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']
    dx, dw, db = np.zeros((N, C, H, W)), np.zeros((F, C, HH, WW)), np.zeros((F,))
    num_filters = F
    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.
    # =====

    N, C, H, W = x.shape
    F, C, HH, WW = w.shape
    dx, dw, db = np.zeros_like(xpad), np.zeros_like(dw), np.zeros_like(db)

    # 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*stride:row*stride+HH, col*stride:col*stride+WW])

    # 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*stride:row*stride+HH, col*stride:col*stride+WW])

    dx = dx.reshape((N, C, H, W))
    dw = dw.reshape((F, C, HH, WW))
    db = db.reshape((F,))

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

    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, 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
    # =====

    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
    N, C, H, W = dout.shape

    dx, dw, db = np.zeros_like(x), np.zeros_like(dw), np.zeros_like(db)

    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[kernal_idx, row_idx, col_idx] = np.sum(dout[img_idx, kernal_idx, row*stride:row*stride+pool_height, col*stride:col*stride+pool_width])

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

    N, C, H, W = x.shape
    x_transpose = x.transpose(0, 2, 3, 1)
    dx, dw, db = np.zeros_like(x), np.zeros_like(dw), np.zeros_like(db)
    out_2d, cache = batchnorm_forward(x_transpose, 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.
    # =====

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

    dout_transpose = dout.transpose(0, 2, 3, 1)
    dx_2d, dgamma, dbeta = batchnorm_backward(dout_transpose, cache)
    dx = dx_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.
    # =====

    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,)

    """
    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], x.shape[0]))
    dx_reshape = np.dot(dout, w.T)
    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

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

    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 (ReLUs).

    Inputs:
    - 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 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:
        # (1) Calculate the running mean and variance of the minibatch.
        # (2) Normalize the activations with the running mean and variance.
        # (3) Store the normalized and scaled 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)
        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_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_normalize = cache.get('x_normalize')
    gamma = cache.get('gamma')
    eps = cache.get('eps')

    # Calculate dx
    dx_hat = dout * gamma
    sqrt_var = np.sqrt(minibatch_var + eps)
    dx_normalize = dx_hat / sqrt_var
    dx = dx_normalize * sqrt_var
    dx_hat = dx_normalize * sqrt_var
    dx2 = np.sum(dx_normalize * dx_normalize, axis=0) / N
    dx = dx1 * dx2

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

    # =====
    # 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 drop neuron output with probability p.
        - 'mode': 'train' or 'test'; 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.

    Returns:
    - out: Array of the same shape as x.
    - mask: A tuple of (mask, cache). 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 is not None:
        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
        # 'p' dropout mask as the variable mask.
        # =====

        mask = np.random.rand(x.shape[0]) >= p / (1 - p)
        out = x * 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
    """
    probs = np.exp(x - np.max(x, axis=1), keepdims=True)
    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

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

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

    Returns a tuple of:
    - loss: Scalar giving the loss
    - dx: Gradient of the loss with respect to x
    """
    probs = np.exp(x - np.max(x, axis=1), keepdims=True)
    probs /= np.sum(probs, axis=1, keepdims=True)
    loss = -np.sum(np.log(probs[np.arange(N), y])) / N
    dx = probs.copy()
    dx[np.arange(N), y] -= 1
    return loss, dx
```