


```
import numpy as np

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

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

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

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

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

    x_res = x.reshape((x.shape[0], w.shape[0])) # Shape of N * D
    out = x_res.dot(w) + b.reshape((1, b.shape[0])) # Shape of N * M

    # =====
    # END YOUR CODE HERE

    # 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 the same shape as x
    - cache: x

    """
    # =====
    # YOUR CODE HERE
    # Implement the ReLU forward pass.
    # =====

    out = np.maximum(0, x)

    # =====
    # END YOUR CODE HERE

    return out, cache

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

    Inputs:
    - dout: Upstream derivatives, of shape (N, M)
    - cache: Tuple of:
      - 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:
    - 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.
    # Notice:
    # - dx should be N x d_1 x ... x d_k; it relates to dout through multiplication with w, which is D x M
    # - dw should be D x M; it relates to dout through multiplication with x, which is N x D after reshaping
    # - db should be M; it is just the sum over dout examples
    # =====

    x_res = np.reshape(x, (x.shape[0], w.shape[0]))
    dx_res = dout.dot(w.T)
    dx = np.reshape(dx_res, x.shape) # Shape of N * D
    dw = x_res.T.dot(dout) # Shape of D * M
    db = dout.T.dot(np.ones(x.shape[0])) # Shape of 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).

    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

    return out, cache

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

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

    Returns a tuple of:
    - dx: Gradient with respect to x
    - cache: x

    """
    # =====
    # 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 mean and variance using
    an exponential decay based on the momentum parameter:

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

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

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

    Returns a tuple of:
    - out: of shape (N, D)
    - cache: A tuple of values needed in the backward pass
    """
    mode = bn_param['mode']
    eps = bn_param.get('eps', 1e-5)
    momentum = bn_param.get('momentum', 0.9)

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

    out, cache = None, None
    if mode == 'train':
        # =====
        # YOUR CODE HERE
        # A few steps here:
        # (1) Calculate the running mean and variance of the minibatch.
        # (2) Normalize the activations with the running mean and variance.
        # (3) Scale and shift the normalized activations. Store this
        #     as the variable 'out'
        # (4) Store any variables you may need for the backward pass in
        #     the 'cache' variable.
        # =====

        mean_minibatch = np.mean(x, axis=0)
        var_minibatch = np.var(x, axis=0)
        x_normalizd = (x - mean_minibatch) / np.sqrt(var_minibatch + eps)
        out = gamma * x_normalizd + beta

        running_mean = running_mean + momentum * x_normalizd
        running_var = running_var + momentum * var_minibatch
        bn_param['running_mean'] = running_mean
        bn_param['running_var'] = running_var

        cache = (
            'minibatch_mean': mean_minibatch,
            'x_normalizd': x_normalizd,
            '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

    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_normalizd = cache.get('x_normalizd')
    gamma = cache.get('gamma')
    eps = cache.get('eps')

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

    # calculate dgamma and dbeta
    dgamma = np.sum(dout, axis=0)
    dbeta = np.sum(dout * x_normalizd, 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):
    """
    Perform the backward pass for (inverted) dropout.

    Inputs:
    - dout: Upstream derivatives, of any shape
    - cache: (dropout_param, mask) from dropout_forward.

    """
    dropout_param, mask = cache
    mode = dropout_param['mode']

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

        dx = 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)
    k = np.zeros_like(y)
    dx[margins > 0, y] = -num_pos
    dx[np.arange(N), y] = num_pos
    return loss, dx

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

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

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

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

```
In [ ]:

from layers import *

def affine_relu_forward(x, w, b):
    """
    Convenience layer that performs an affine transform followed by a ReLU

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

    Returns a tuple of:
    - out: Output from the ReLU
    - cache: Object to give to the backward pass

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

def affine_relu_backward(dout, cache):
    """
    Backward pass for the affine-relu convenience layer

    Inputs:
    - fc_cache, relu_cache: cache
    - da = relu_backward(dout, relu_cache)
    - dx, dw, db = affine_backward(da, fc_cache)

    """
    In [ ]:

import numpy as np
from layers import *
from layer_utils import *

class TwoLayerNet(object):
    """
    A two-layer fully-connected neural network with ReLU nonlinearity and
    softmax loss that uses a modular layer design. We assume an input dimension
    of D, a hidden dimension of H, and perform classification over C classes.

    The architecture should be affine - relu - affine - softmax.

    Note that this class does not implement gradient descent; instead, it
    will interact with a separate Solver object that is responsible for running
    optimization.

    The learnable parameters of the model are stored in the dictionary
    self.params that maps parameter names to numpy arrays.
    """
    def __init__(self, input_dim=3*32*32, hidden_dim=100, num_classes=10,
                 dropout=0, weight_scale=1e-3, reg=0.0):
        """
        Initialize a new network.

        Inputs:
        - input_dim: An integer giving the size of the input
        - hidden_dim: An integer giving the size of the hidden layer
        - num_classes: An integer giving the number of classes to classify
        - dropout: Scalar between 0 and 1 giving dropout strength.
        - weight_scale: Scalar giving the standard deviation for random
          initialization of the weights.
        - reg: Scalar giving L2 regularization strength.

        self.params = {}
        self.reg = reg

        """
        # YOUR CODE HERE:
        # Initialize W1, W2, b1, and b2. Store these as self.params['W1'],
        # self.params['W2'], self.params['b1'], and self.params['b2']. The
        # biases are initialized to zero and the weights are initialized
        # so that each parameter has mean 0 and standard deviation weight_scale.
        # The dimensions of W1 should be (input_dim, hidden_dim) and the
        # dimensions of W2 should be (hidden_dim, num_classes)
        # =====

        W1_size = (input_dim, hidden_dim)
        W2_size = (hidden_dim, num_classes)

        self.params['W1'] = np.random.normal(1e-08, weight_scale, size=W1_size)
        self.params['b1'] = np.zeros(hidden_dim)
        self.params['W2'] = np.random.normal(1e-08, weight_scale, size=W2_size)
        self.params['b2'] = np.zeros(num_classes)

        # =====
        # END YOUR CODE HERE

    def loss(self, X, y=None):
        """
        Compute loss and gradient for a minibatch of data.

        Inputs:
        - X: Array of input data of shape (N, d_1, ..., d_k)
        - y: Array of labels, of shape (N,) y[i] gives the label for X[i].

        Returns:
        If y is None, then run a test-time forward pass of the model and return:
        - scores: Array of shape (N, C) giving classification scores; C is number
          of classes, where C is the classification score for X[i] and class c.

        If y is not None, then run a training-time forward and backward pass and
        return a tuple of:
        - loss: Scalar value giving the loss
        - grads: Dictionary with the same keys as self.params, mapping parameter
          names to gradients of the loss with respect to those parameters.

        """
        scores = None

        # =====
        # YOUR CODE HERE
        # Implement the forward pass of the two-layer neural network. Store
        # the class scores as the variable 'scores'. Be sure to use the layers
        # you prior implemented.
        # =====

        W1 = self.params['W1']
        b1 = self.params['b1']
        W2 = self.params['W2']
        b2 = self.params['b2']

        Z, cache_h = affine_relu_forward(X, W1, b1)
        Z, cache_z = affine_relu_forward(Z, W2, b2)

        scores = Z

        # =====
        # END YOUR CODE HERE

    # If y is None then we are in test mode so just return scores
    if y is None:
        return scores

    # =====
    # YOUR CODE HERE
    # Implement the backward pass of the two-layer neural net. Store
    # the loss as the variable 'loss' and store the gradients in the
    # 'scores' dictionary. For the grads dictionary, grads['W1'] holds
    # the gradient for W1, grads['b1'] holds the gradient for b1, etc.
    # i.e., grads[k] holds the gradient for self.params[k].
    # =====
    # Add L2 regularization, where there is an added cost 0.5*self.reg*W^2
    # for each W. Be sure to include the 0.5 multiplying factor to
    # match our implementation.
    # And be sure to use the layers you prior implemented.
    # =====

    loss, dz = softmax_loss(scores, y)
    loss += 0.5*self.reg*(np.sum(W1**2) + np.sum(W2**2))

    dh, dw2, db2 = affine_backward(dz, cache_z)
    dx, dw1, db1 = affine_relu_backward(dh, cache_h)

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

    # =====
    # END YOUR CODE HERE

    return loss, grads

class FullyConnectedNet(object):
    """
    A fully-connected neural network with an arbitrary number of hidden layers,
    ReLU nonlinearities, and a softmax loss function. This will also implement
    dropout and batch normalization as options. For a network with L layers,
    the architecture will be

    (affine - [batch norm] - relu - [dropout]) x (L - 1) - affine - softmax

    where batch normalization and dropout are optional, and the (...) block is
    repeated L - 1 times.

    Similar to the TwoLayerNet above, learnable parameters are stored in the
    self.params dictionary and will be learned using the Solver class.
    """
    def __init__(self, hidden_dims, input_dim=3*32*32, hidden_dim=100,
                 dropout=0, weight_scale=1e-2, dtype=np.float32, seed=None):
        """
        Initialize a new FullyConnectedNet.

        Inputs:
        - hidden_dims: A list of integers giving the size of each hidden layer
        - input_dim: An integer giving the size of the input
        - num_classes: An integer giving the number of classes to classify
        - dropout: Scalar between 0 and 1 giving dropout strength.
        - weight_scale: Scalar giving the standard deviation for random
          initialization of the weights.
        - dtype: A numpy datatype object; all computations will be performed using
          this datatype. float32 is faster but less accurate, so you should use
          float64 for numeric gradient checking.
        - seed: If not None, then pass this random seed to the dropout layers. This
          will make the dropout behavior deterministic so we can gradient check the
          model.

        """
        self.use_batchnorm = use_batchnorm
        self.use_dropout = dropout > 0
        self.reg = reg
        self.num_layers = 1 + len(hidden_dims)
        self.dtype = dtype
        self.params = {}

        # =====
        # YOUR CODE HERE
        # Initialize all parameters of the network in the self.params dictionary.
        # The weights and biases of layer 1 are W1 and b1; and in general the
        # weights and biases of layer i are Wi and bi. The
        # biases are initialized to zero and the weights are initialized
        # so that each parameter has mean 0 and standard deviation weight_scale.
        # =====
        # BATCHNORM: Initialize the gammas of each layer to 1 and the beta
        # parameters to zero. The gamma and beta parameters for layer l should
        # be self.params['gamma'+str(l)] and self.params['beta'+str(l)]. For layer 1, they
        # should be gamma2 and beta2, etc. Only use batchnorm if self.use_batchnorm
        # is true and DO NOT do batch normalize the output scores.
        # =====

        cur_dim = input_dim
        for idx, hidden_dim in enumerate(hidden_dims):
            self.params['W'+str(idx+1)] = np.random.randn(cur_dim, hidden_dim) * weight_scale
            self.params['b'+str(idx+1)] = np.zeros(hidden_dim)

            if self.use_batchnorm:
                self.params['gamma'+str(idx+1)] = np.ones(hidden_dim)
                self.params['beta'+str(idx+1)] = np.zeros(hidden_dim)

            cur_dim = hidden_dim

        self.params['W'+str(self.num_layers)] = np.random.randn(cur_dim, num_classes) * weight_scale
        self.params['b'+str(self.num_layers)] = np.zeros(num_classes)

        # =====
        # END YOUR CODE HERE

    # When using dropout we need to pass a dropout_param dictionary to each
    # dropout layer so that the layer knows the dropout probability and the mode
    # (train / test). You can pass the same dropout_param to each dropout layer.
    self.dropout_param = {}
    if self.use_dropout:
        self.dropout_param = {'mode': 'train', 'p': dropout}
    if seed is not None:
        self.dropout_param['seed'] = seed

    # With batch normalization we need to keep track of running means and
    # variances, so we need to pass a special bn_param object to each batch
    # normalization layer. You should pass self.bn_params[0] to the forward pass
    # of the first batch normalization layer, self.bn_params[1] to the forward
    # pass of the second batch normalization layer, etc.
    self.bn_params = []
    if self.use_batchnorm:
        self.bn_params = [{'mode': 'train'} for i in np.arange(self.num_layers - 1)]

    # Cast all parameters to the correct datatype
    for k, v in self.params.items():
        self.params[k] = v.astype(dtype)

    def loss(self, X, y=None):
        """
        Compute loss and gradient for the fully-connected net.

        Input / output: Same as TwoLayerNet above.

        """
        X = X.astype(self.dtype)
        mode = 'test' if y is None else 'train'

        # Set train/test mode for batchnorm params and dropout param since they
        # behave differently during training and testing.
        if self.dropout_param is not None:
            self.dropout_param['mode'] = mode
        if self.use_batchnorm:
            for bn_param in self.bn_params:
                bn_param['mode'] = mode

        scores = None

        # =====
        # YOUR CODE HERE
        # Implement the forward pass of the FC net and store the output
        # scores as the variable "scores".
        # =====
        # BATCHNORM: If self.use_batchnorm is true, insert a batchnorm layer
        # between the affine forward and linear layers. You may
        # also write an affine_batchnorm_relu function in layer_utils.py.
        # DROPOUT: If dropout is non-zero, insert a dropout layer after
        # every ReLU layer.
        # =====

        fc_cache = {}
        relu_cache = {}
        batchnorm_cache = {}
        dropout_cache = {}

        X = np.reshape(X, (X.shape[0], -1))

        # go through all layers
        for i in range(self.num_layers - 1):
            # fc layer
            fc_out, fc_cache[str(i+1)] = affine_forward(X, self.params['W'+str(i+1)], self.params['b'+str(i+1)])

            # batchnorm layer
            relu_input = fc_out
            if self.use_batchnorm:
                batchnorm_out, batchnorm_cache[str(i+1)] = batchnorm_forward(fc_out, self.params['gamma'+str(i+1)], self.params['beta'+str(i+1)], self.bn_params[i])
            else:
                relu_input = batchnorm_out
            relu_out, relu_cache[str(i+1)] = relu_forward(relu_input)

            # dropout layer
            if self.use_dropout:
                relu_out, dropout_cache[str(i+1)] = dropout_forward(relu_out, self.dropout_param)

            # cache X
            cache_X = relu_out.copy()

            # output final affine with no relu
            scores, final_cache = affine_forward(X, self.params['W'+str(self.num_layers)], self.params['b'+str(self.num_layers)], self.params['gamma'+str(self.num_layers)], self.params['beta'+str(self.num_layers)], self.bn_params[-1])

            # =====
            # END YOUR CODE HERE

        # If test mode return early
        if mode == 'test':
            return scores

        # =====
        # YOUR CODE HERE
        # Implement the backward pass of the FC net and store the gradients
        # in the grads dict, so that grads[k] is the gradient of self.params[k]
        # Be sure your L2 regularization includes a 0.5 factor.
        # BATCHNORM: Incorporate the backward pass of the batchnorm.
        # DROPOUT: Incorporate the backward pass of dropout.
        # =====

        loss, dx = softmax_loss(scores, y)
        loss += 0.5*self.reg * (np.sum(np.square(self.params['W'+str(self.num_layers)])))
        dx_back, dw_back, db_back = affine_backward(dx, final_cache)
        grads['W'+str(self.num_layers)] = dw_back + self.reg * self.params['W'+str(self.num_layers)]
        grads['b'+str(self.num_layers)] = db_back

        # go backward all layers and update weights, bias, gammas and betas
        for i in range(self.num_layers - 1, 0, -1):
            # dropout layer
            if self.use_dropout:
                dx_back, dw_back, dropout_cache[str(i)] = dropout_backward(dx_back, dropout_cache[str(i)])
            else:
                dx_back, dw_back, db_back = affine_backward(dx_back, relu_cache[str(i)])

            # batchnorm layer
            affine_backward_input = dx_relu
            if self.use_batchnorm:
                dx_bn, dgamma, dbeta = batchnorm_backward(dx_relu, batchnorm_cache[str(i)])
                grads['gamma'+str(i)] = dgamma
                grads['beta'+str(i)] = dbeta
            affine_backward_input = dx_bn
            dx_back, dw_back, db_back = affine_backward(affine_backward_input, fc_cache[str(i)])

            grads['W'+str(i)] = dw_back + self.reg * self.params['W'+str(i)]
            grads['b'+str(i)] = db_back
            loss += 0.5 * self.reg * (np.sum(np.square(self.params['W'+str(i)])))

        # =====
        # END YOUR CODE HERE

        return loss, grads
```



```
import numpy as np
```

```
"""
```

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

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

```
Inputs:
```

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

```
Returns:
```

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

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

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

```
"""
```

```
def sgd(w, dw, config=None):
```

```
"""
```

Performs vanilla stochastic gradient descent.

```
config format:
```

- `learning_rate`: Scalar learning rate.

```
if config is None: config = {}
```

```
config.setdefault('learning_rate', 1e-2)
```

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

```
def sgd_momentum(w, dw, config=None):
```

Performs stochastic gradient descent with momentum.

```
config format:
```

- `learning_rate`: Scalar learning rate.
- `momentum`: Scalar between 0 and 1 giving the momentum value. Setting momentum = 0 reduces to sgd.
- `velocity`: A numpy array of the same shape as `w` and `dw` used to store a moving average of the gradients.

```
if config is None: config = {}
config.setdefault('learning_rate', 1e-2)
config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.

# =====
# YOUR CODE HERE:
# Implement the momentum update formula. Return the updated weights
# as next_w, and the updated velocity as v.
# =====
```

```
v = config['momentum'] * v + config['learning_rate'] * dw
next_w = v + w
```

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

```
config['velocity'] = v
```

```
return next_w, config
```

```
def sgd_nesterov_momentum(w, dw, config=None):
```

Performs stochastic gradient descent with Nesterov momentum.

```
config format:
```

- `learning_rate`: Scalar learning rate.
- `momentum`: Scalar between 0 and 1 giving the momentum value.
- `decay_rate`: Scalar between 0 and 1 reducing to sgd.
- `velocity`: A numpy array of the same shape as `w` and `dw` used to store a moving average of the gradients.

```
if config is None: config = {}
config.setdefault('learning_rate', 1e-2)
config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.

# =====
# YOUR CODE HERE:
# Implement the momentum update formula. Return the updated weights
# as next_w, and the updated velocity as v.
# =====
```

```
v_0 = v
v = config['momentum'] * v + config['learning_rate'] * dw
w += v + config['momentum'] * (v - v_0)
next_w = w
```

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

```
config['velocity'] = v
```

```
return next_w, config
```

```
def rmsprop(w, dw, config=None):
```

Uses the RMSProp update rule, which uses a moving average of squared gradient values to set adaptive per-parameter learning rates.

```
config format:
```

- `learning_rate`: Scalar learning rate.
- `decay_rate`: Scalar between 0 and 1 giving the decay rate for the squared gradient cache.
- `epsilon`: Small scalar used for smoothing to avoid dividing by zero.
- `beta`: Moving average of second moments of gradients.

```
if config is None: config = {}
config.setdefault('learning_rate', 1e-2)
config.setdefault('decay_rate', 0.99)
config.setdefault('epsilon', 1e-8)
config.setdefault('beta', np.zeros_like(w))
```

```
next_w = None
```

```
# =====
# YOUR CODE HERE:
# Implement RMSProp. Store the next value of w as next_w. You need
# to also store in config['s'] the moving average of the second
# moment gradients, so they can be used for future gradients. Concretely,
# config['s'] corresponds to "s" in the lecture notes.
# =====
```

```
config['s'] = config['decay_rate'] * config['s'] + (1 - config['decay_rate']) * (dw ** 2)
next_w = w - config['learning_rate'] * dw / (np.sqrt(config['s']) + config['epsilon'])
```

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

```
return next_w, config
```

```
def adam(w, dw, config=None):
```

```
"""
```

Uses the Adam update rule, which incorporates moving averages of both the gradient and its square and a bias correction term.

```
config format:
```

- `learning_rate`: Scalar learning rate.
- `beta1`: Decay rate for moving average of first moment of gradient.
- `beta2`: Decay rate for moving average of second moment of gradient.
- `epsilon`: Small scalar used for smoothing to avoid dividing by zero.
- `m`: Moving average of gradient.
- `v`: Moving average of squared gradient.
- `t`: Iteration number.

```
"""
```

```
if config is None: config = {}
config.setdefault('learning_rate', 1e-3)
config.setdefault('beta1', 0.9)
config.setdefault('beta2', 0.999)
config.setdefault('epsilon', 1e-8)
config.setdefault('v', np.zeros_like(w))
config.setdefault('m', np.zeros_like(w))
config.setdefault('t', 0)
```

```
next_w = None
```

```
# =====
# YOUR CODE HERE:
# Implement Adam. Store the next value of w as next_w. You need
# to also store in config['s'] the moving average of the second
# moment gradients, and in config['v'] the moving average of the
# first moments. Finally, store in config['t'] the increasing time.
# =====
```

```
beta1 = config['beta1']
beta2 = config['beta2']
t = config['t'] + 1
```

```
v = beta1 * config['v'] + (1 - beta1) * dw
a = beta2 * config['s'] + (1 - beta2) * (dw ** 2)
v_corrected = v / (1 - beta1 ** t)
a_corrected = a / (1 - beta2 ** t)
next_w = w - config['learning_rate'] * v_corrected / (np.sqrt(a_corrected) + config['epsilon'])
```

```
config['v'] = v
config['s'] = a
config['t'] = t
```

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

```
return next_w, config
```


ECE C147/247 HW4 Q2: Batch Normalization

In this notebook, you will implement the batch normalization layers of a neural network to increase its performance. Please review the details of batch normalization from the lecture notes.

`utils` has built a solid API for building these modular frameworks and training them, and we will use their very well implemented framework as opposed to "reinventing the wheel." This includes using the Solver, various utility functions, and their layer structure. This also includes `nnrl.fc_net`, `nnrl.layers`, and `nnrl.layer_utils`.

```
In [1]: ## Import and Setup

import time
import numpy as np
import matplotlib.pyplot as plt
from nnrl.fc_net import *
from nnrl.layers import *
from nnrl.data_utils import get_CIFAR10_data
from utils.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
from nnrl.solver import Solver

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

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
%load_ext autoreload
%autoreload 2

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

In [2]: # Load the (preprocessed) CIFAR10 data.

data = get_CIFAR10_data()
for k in data.keys():
    print(' %s: %s' % (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,)
```

Batchnorm forward pass

Implement the training time batchnorm forward pass, `batchnorm_forward`, in `nnrl/layers.py`. After that, test your implementation by running the following cell.

```
In [3]: # Check the training-time forward pass by checking means and variances
# of features both before and after batch normalization

# Simulate the forward pass for a two-layer network
N, D1, D2, D3 = 200, 50, 60, 3
X = np.random.randn(N, D1)
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)
a = np.maximum(0, X.dot(W1)).dot(W2)

print('Before batch normalization:')
print(' means: ', a.mean(axis=0))
print(' stds: ', a.std(axis=0))

# Means should be close to zero and stds close to one
print('After batch normalization (gamma=1, beta=0)')
a_norm, _ = batchnorm_forward(a, np.ones(D3), np.zeros(D3), {'mode': 'train'})
print(' mean: ', a_norm.mean(axis=0))
print(' std: ', a_norm.std(axis=0))

# Now means should be close to beta and stds close to gamma
gamma = np.asarray([1.0, 2.0, 3.0])
beta = np.asarray([11.0, 12.0, 13.0])
a_norm = batchnorm_forward(a, gamma, beta, {'mode': 'train'})
print('After batch normalization (nontrivial gamma, beta)')
print(' means: ', a_norm.mean(axis=0))
print(' stds: ', a_norm.std(axis=0))

Before batch normalization:
means: [-28.84907154 -25.98395893 -2.11862989]
stds: [31.89265604 38.41952822 28.18587946]
After batch normalization (gamma=1, beta=0)
mean: [-2.22044605e-17 2.22044605e-17 -4.21884749e-17]
std: [1. 0.99999999]
After batch normalization (nontrivial gamma, beta)
means: [11. 12. 13.]
stds: [1. 1.99999999 2.99999998]
```

Implement the testing time batchnorm forward pass, `batchnorm_forward`, in `nnrl/layers.py`. After that, test your implementation by running the following cell.

```
In [4]: # Check the test-time forward pass by running the training-time
# forward pass many times to warm up the running averages, and then
# checking the means and variances of activations after a test-time
# forward pass.

N, D1, D2, D3 = 200, 50, 60, 3
X = np.random.randn(N, D1)
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)

bn_param = {'mode': 'train'}
gamma = np.ones(D3)
beta = np.zeros(D3)
for i in np.arange(50):
    X = np.random.randn(N, D1)
    a = np.maximum(0, X.dot(W1)).dot(W2)
    batchnorm_forward(a, gamma, beta, bn_param)
bn_param['mode'] = 'test'
X = np.random.randn(N, D1)
a = np.maximum(0, X.dot(W1)).dot(W2)
a_norm, _ = batchnorm_forward(a, gamma, beta, bn_param)

# Means should be close to zero and stds close to one, but will be
print('After batch normalization (test-time):')
print(' means: ', a_norm.mean(axis=0))
print(' stds: ', a_norm.std(axis=0))

After batch normalization (test-time):
means: [1.569444183941 -11.55932501338 4320.14438343]
stds: [9330.88699977 9424.10712598 9575.23599007]
```

Batchnorm backward pass

Implement the backward pass for the batchnorm layer, `batchnorm_backward` in `nnrl/layers.py`. Check your implementation by running the following cell.

```
In [5]: # Gradient check batchnorm backward pass

N, D = 4, 5
x = 5 * np.random.randn(N, D) + 12
gamma = np.random.randn(D)
beta = np.random.randn(D)
dout = np.random.randn(N, D)

bn_param = {'mode': 'train'}
fx = lambda x: batchnorm_forward(x, gamma, beta, bn_param)[0]
fy = lambda a: batchnorm_backward(x, gamma, beta, bn_param)[0]
fb = lambda b: batchnorm_backward(x, gamma, beta, bn_param)[0]

dx_num = eval_numerical_gradient_array(fx, x, dout)
da_num = eval_numerical_gradient_array(fy, gamma, dout)
db_num = eval_numerical_gradient_array(fb, beta, dout)

_, cache = batchnorm_forward(x, gamma, beta, bn_param)
dx, dgamma, dbeta = batchnorm_backward(dout, cache)
print('dx error: ', rel_error(dx_num, dx))
print('dgamma error: ', rel_error(da_num, dgamma))
print('dbeta error: ', rel_error(db_num, dbeta))

dx error: 1.7703768271105096e-09
dgamma error: 3.707945705184271e-11
dbeta error: 3.275692151309245e-12
```

Implement a fully connected neural network with batchnorm layers

Modify the `FullyConnectedNet()` class in `nnrl/fc_net.py` to incorporate batchnorm layers. You will need to modify the class in the following areas:

- (1) The gammas and betas need to be initialized to 1's and 0's respectively in `__init__`.
- (2) The `batchnorm_forward` layer needs to be inserted between each affine and relu layer (except in the output layer) in a forward pass computation in `loss`. You may find it helpful to write an `affine_batchnorm_relu()` layer in `nnrl/layer_utils.py` although this is not necessary.
- (3) The `batchnorm_backward` layer has to be appropriately inserted when calculating gradients.

After you have done the appropriate modifications, check your implementation by running the following cell.

Note, while the relative error for W3 should be small, as we backprop gradients more, you may find the relative error increases. Our relative error for W1 is on the order of $1e-4$.

```
In [6]: N, D, H1, H2, C = 2, 15, 20, 30, 10
X = np.random.randn(N, D)
y = np.random.randint(C, size=(N,))

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

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

    for name in sorted(grads):
        f = lambda x: model.loss(X, y)[0]
        grad_num = eval_numerical_gradient(f, model.params[name], verbose=False, h=1e-5)
        print(' %s relative error: %s' % (name, rel_error(grad_num, grads[name])))
    if reg == 0: print('\n')

Running check with reg = 0
Initial loss: 2.3428296959216444
W1 relative error: 0.0018034962629854376
W2 relative error: 3.96234040407421e-05
W3 relative error: 3.51907905294395e-10
b1 relative error: 0.002220443273692751
b2 relative error: 1.1102230246251565e-08
b3 relative error: 1.7625948184968003e-10
beta1 relative error: 3.2409224864496935e-08
beta2 relative error: 1.5593687895310636e-08
gamma1 relative error: 3.1089775687156094e-08
gamma2 relative error: 1.200086744520607e-08

Running check with reg = 3.14
Initial loss: 6.95488371352358
W1 relative error: 0.0001695121156961158
W2 relative error: 4.82658633053028e-06
W3 relative error: 7.172263548226862e-09
b1 relative error: 5.551115123125783e-09
b2 relative error: 4.440892098500626e-08
b3 relative error: 2.706977503819273e-10
beta1 relative error: 1.9583724009991693e-08
beta2 relative error: 1.6570164039795128e-08
gamma1 relative error: 4.618920964413653e-08
gamma2 relative error: 1.728582850964246e-08
```

Training a deep fully connected network with batch normalization.

To see if batchnorm helps, let's train a deep neural network with and without batch normalization.

```
In [7]: # Try training a very deep net with batchnorm
hidden_dims = [100, 100, 100, 100]

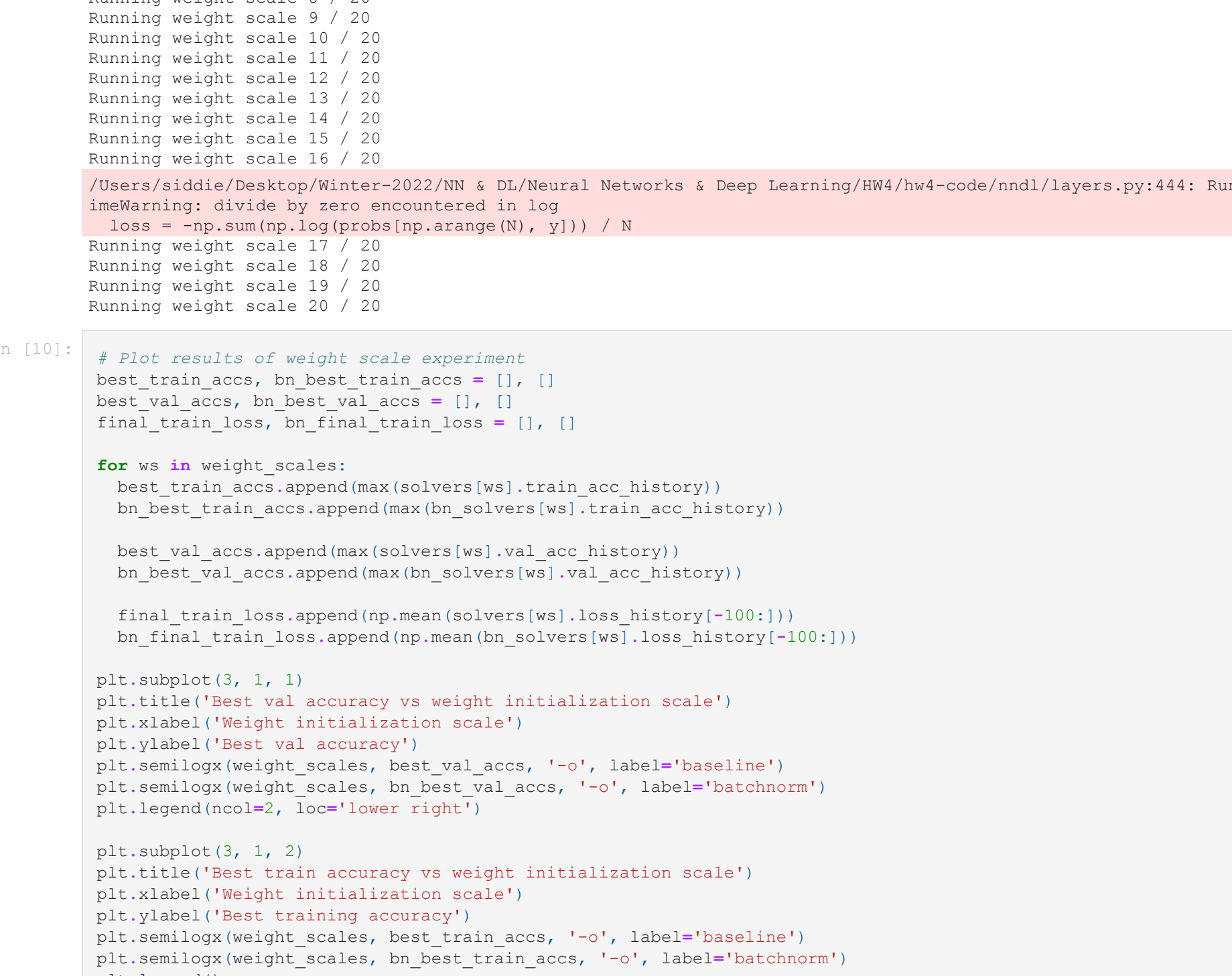
num_train = 1000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

weight_scale = 2e-2
bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=True)
model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=False)

bn_solver = Solver(bn_model, small_data,
                  num_epochs=10, batch_size=50,
                  update_rule='adam',
                  optim_config={
                      'learning_rate': 1e-3,
                  },
                  verbose=True, print_every=200)
bn_solver.train()

solver = Solver(model, small_data,
                num_epochs=10, batch_size=50,
                update_rule='adam',
                optim_config={
                    'learning_rate': 1e-3,
                },
                verbose=True, print_every=200)
solver.train()

(iteration 1 / 200) loss: 2.300179
(Epoch 0 / 10) train acc: 0.119000; val_acc: 0.104000
(Epoch 1 / 10) train acc: 0.322000; val_acc: 0.268000
(Epoch 2 / 10) train acc: 0.436000; val_acc: 0.311000
(Epoch 3 / 10) train acc: 0.535000; val_acc: 0.302000
(Epoch 4 / 10) train acc: 0.569000; val_acc: 0.337000
(Epoch 5 / 10) train acc: 0.607000; val_acc: 0.318000
(Epoch 6 / 10) train acc: 0.653000; val_acc: 0.306000
(Epoch 7 / 10) train acc: 0.739000; val_acc: 0.345000
(Epoch 8 / 10) train acc: 0.725000; val_acc: 0.299000
(Epoch 9 / 10) train acc: 0.769000; val_acc: 0.314000
(Epoch 10 / 10) train acc: 0.808000; val_acc: 0.338000
(iteration 1 / 200) loss: 2.303840
(Epoch 0 / 10) train acc: 0.146000; val_acc: 0.143000
(Epoch 1 / 10) train acc: 0.243000; val_acc: 0.210000
(Epoch 2 / 10) train acc: 0.284000; val_acc: 0.254000
(Epoch 3 / 10) train acc: 0.316000; val_acc: 0.246000
(Epoch 4 / 10) train acc: 0.374000; val_acc: 0.274000
(Epoch 5 / 10) train acc: 0.407000; val_acc: 0.293000
(Epoch 6 / 10) train acc: 0.436000; val_acc: 0.283000
(Epoch 7 / 10) train acc: 0.502000; val_acc: 0.306000
(Epoch 8 / 10) train acc: 0.528000; val_acc: 0.308000
(Epoch 9 / 10) train acc: 0.575000; val_acc: 0.304000
(Epoch 10 / 10) train acc: 0.598000; val_acc: 0.326000
```



Batchnorm and initialization

The following cells run an experiment where for a deep network, the initialization is varied. We do training for when batchnorm layers are and are not included.

```
In [9]: # Try training a very deep net with batchnorm
hidden_dims = [50, 50, 50, 50, 50, 50]

num_train = 1000
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'],
}

solvers = {}
weight_scales = np.logspace(-4, 0, num=20)
for i, weight_scale in enumerate(weight_scales):
    print('Running weight scale %f' % (weight_scale))
    bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=True)
    model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=False)

    bn_solver = Solver(bn_model, small_data,
                      num_epochs=10, batch_size=50,
                      update_rule='adam',
                      optim_config={
                          'learning_rate': 1e-3,
                      },
                      verbose=False, print_every=200)
    bn_solver.train()
    solvers[weight_scale] = bn_solver

    solver = Solver(model, small_data,
                    num_epochs=10, batch_size=50,
                    update_rule='adam',
                    optim_config={
                        'learning_rate': 1e-3,
                    },
                    verbose=False, print_every=200)
    solver.train()
    solvers[weight_scale] = solver

Running weight scale 1 / 20
Running weight scale 2 / 20
Running weight scale 3 / 20
Running weight scale 4 / 20
Running weight scale 5 / 20
Running weight scale 6 / 20
Running weight scale 7 / 20
Running weight scale 8 / 20
Running weight scale 9 / 20
Running weight scale 10 / 20
Running weight scale 11 / 20
Running weight scale 12 / 20
Running weight scale 13 / 20
Running weight scale 14 / 20
Running weight scale 15 / 20
Running weight scale 16 / 20

/Users/siddie/Desktop/Winter2022/NN & DL/Neural Networks & Deep Learning/HW4/hw4-code/nnrl/layers.py:444: RuntimeWarning: divide by zero encountered in log
    loss = -np.sum(np.log(probs[np.arange(N), y])) / N
Running weight scale 17 / 20
Running weight scale 18 / 20
Running weight scale 19 / 20
Running weight scale 20 / 20

In [10]: # Plot results of weight scale experiment
best_train_accs, bn_best_train_accs = [], []
best_val_accs, bn_best_val_accs = [], []
final_train_loss, bn_final_train_loss = [], []

for ws in weight_scales:
    best_train_accs.append(max(solvers[ws].train_acc_history))
    bn_best_train_accs.append(max(bn_solvers[ws].train_acc_history))
    best_val_accs.append(max(solvers[ws].val_acc_history))
    bn_best_val_accs.append(max(bn_solvers[ws].val_acc_history))

    final_train_loss.append(np.mean(solvers[ws].loss_history[-100:]))
    bn_final_train_loss.append(np.mean(bn_solvers[ws].loss_history[-100:]))

plt.subplot(3, 1, 1)
plt.title('Best val accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.semilogx(weight_scales, best_val_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_val_accs, '-o', label='batchnorm')
plt.legend(ncol=2, loc='lower right')

plt.subplot(3, 1, 2)
plt.title('Best train accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.semilogx(weight_scales, best_train_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_train_accs, '-o', label='batchnorm')
plt.legend()

plt.subplot(3, 1, 3)
plt.title('Final training loss vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.semilogx(weight_scales, bn_final_train_loss, '-o', label='batchnorm')
plt.show()

Best val accuracy vs weight initialization scale
Best train accuracy vs weight initialization scale
Final training loss vs weight initialization scale
```

Question:

In the cell below, summarize the findings of this experiment, and WHY these results make sense.

Answer:

You fill this in.

Batchnorm makes the model less sensitive to the initialization of the weights and bias parameters. As seen in the loss figure, the loss function of the model with batchnorm highlighted in orange is consistent with deviation no more than 0.3 units. But, in the case of baseline model, we can see the effect of random weight initialization on the loss (Chaning from 2.3 to 1). With batchnorm, the training and validation accuracies are less sensitive to weight initializations in comparison with baseline model.

ECE C147/247 HW4 Q3: Dropout

In this notebook, you will implement dropout. Then we will ask you to train a network with batchnorm and dropout, and achieve over 55% accuracy on CIFAR-10.

`utils` has built a solid API for building these modular frameworks and training them, and we will use their very well implemented framework as opposed to "reinventing the wheel". This includes using their Solver, various utility functions, and their layer structure. This also includes `nnml.fc_net`, `nnml.layers`, and `nnml.layer_utils`.

```
In [3]: ## Import and setup

import time
import numpy as np
import matplotlib.pyplot as plt
from nnml.fc_net import *
from nnml.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))))
```

```
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,)
```

Dropout forward pass

Implement the training and test time dropout forward pass, `dropout_forward`, in `nnml/layers.py`. After that, test your implementation by running the following cell.

```
In [3]: x = np.random.randn(500, 500) + 10

for p in [0.3, 0.6, 0.75]:
    out, _ = dropout_forward(x, {'mode': 'train', 'p': p})
    out_test, _ = dropout_forward(x, {'mode': 'test', 'p': p})

    print('Running tests with p = {}'.format(p))
    print('Mean of input: {}'.format(x.mean()))
    print('Mean of train-time output: {}'.format(out.mean()))
    print('Mean of test-time output: {}'.format(out_test.mean()))
    print('Fraction of train-time output set to zero: {}'.format((out == 0).mean()))
    print('Fraction of test-time output set to zero: {}'.format((out_test == 0).mean()))

Running tests with p = 0.3
Mean of input: 9.9985010214567
Mean of train-time output: 9.97354144793436
Mean of test-time output: 9.9985010214567
Fraction of train-time output set to zero: 0.301636
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.6
Mean of input: 9.9985010214567
Mean of train-time output: 9.995998684721517
Mean of test-time output: 9.9985010214567
Fraction of train-time output set to zero: 0.59936
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.75
Mean of input: 9.9985010214567
Mean of train-time output: 10.00209523824298
Mean of test-time output: 9.99450214567
Fraction of train-time output set to zero: 0.749776
Fraction of test-time output set to zero: 0.0
```

Dropout backward pass

Implement the backward pass, `dropout_backward`, in `nnml/layers.py`. After that, test your gradients by running the following cell:

```
In [4]: x = np.random.randn(10, 10) + 10
dout = np.random.randn(*x.shape)

dropout_param = {'mode': 'train', 'p': 0.8, 'seed': 123}
out, cache = dropout_forward(x, dropout_param)
dx = dropout_backward(dout, cache)
dx_num = eval_numerical_gradient_array(lambda x: dropout_forward(x, dropout_param)[0], x, dout)

print('dx relative error: {}'.format(dx_num))

dx relative error: 1.892906988930252e-11
```

Implement a fully connected neural network with dropout layers

Modify the `FullyConnectedNet` class in `nnml/fc_net.py` to incorporate dropout. A dropout layer should be incorporated after every ReLU layer. Concretely, there should be a dropout at the output layer since there is no ReLU at the output layer. You will need to modify the class in the following areas:

- (1) In the forward pass, you will need to incorporate a dropout layer after every ReLU layer.
- (2) In the backward pass, you will need to incorporate a dropout backward pass layer.

Check your implementation by running the following code. Our W1 gradient relative error is on the order of 1e-6 (the largest of all the relative errors).

```
In [5]: N, D, H1, H2, C = 2, 15, 20, 30, 10
X = np.random.randn(N, D)
y = np.random.randn(C, size=(N,))

for dropout in [0, 0.25, 0.5]:
    print('Running check with dropout = {}'.format(dropout))
    model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
                              weight_scale=5e-2, dtype=np.float64,
                              dropout=dropout, seed=123)

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

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

Running check with dropout = 0
Initial loss: 2.3051948273987857
W1 relative error: 2.3272570382129597e-07
W2 relative error: 1.50348484922161239e-05
W3 relative error: 1.847669394519802e-07
b1 relative error: 2.336957506391204e-06
b2 relative error: 5.0013939785120424e-08
b3 relative error: 1.1749467839205477e-08

Running check with dropout = 0.25
Initial loss: 2.3052077546540826
W1 relative error: 2.6138469533429284e-07
W2 relative error: 5.02205655497157e-07
W3 relative error: 4.4563160427353564e-08
b1 relative error: 1.30711745000869e-08
b2 relative error: 7.151679404650099e-10
b3 relative error: 1.003974732116764e-10

Running check with dropout = 0.5
Initial loss: 2.303567586595423
W1 relative error: 1.1401257489261862e-06
W2 relative error: 1.847669394519802e-07
W3 relative error: 6.5966195253431734e-09
b1 relative error: 7.16359514021063e-08
b2 relative error: 1.175510493920166e-09
b3 relative error: 1.4558471033827801e-10
```

Dropout as a regularizer

In class, we claimed that dropout acts as a regularizer by effectively bagging. To check this, we will train two small networks, one with dropout and one without dropout.

```
In [6]: # Train two identical nets, one with dropout and one without

num_train = 500
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'],
}

solvers = {}
dropout_choices = [0, 0.6]
for dropout in dropout_choices:
    model = FullyConnectedNet([100, 100, 100], dropout=dropout)

    solver = Solver(model, small_data,
                    num_epochs=25, batch_size=100,
                    update_rule='adam',
                    optim_config={
                        'learning_rate': 5e-4,
                    },
                    verbose=True, print_every=100)
    solver.train()
    solvers[dropout] = solver

(iteration 1 / 125) loss: 2.300804
Epoch 0 / 25) train acc: 0.220000; val_acc: 0.168000
Epoch 1 / 25) train acc: 0.188000; val_acc: 0.170000
Epoch 2 / 25) train acc: 0.266000; val_acc: 0.200000
Epoch 3 / 25) train acc: 0.338000; val_acc: 0.260000
Epoch 4 / 25) train acc: 0.278000; val_acc: 0.278000
Epoch 5 / 25) train acc: 0.428000; val_acc: 0.297000
Epoch 6 / 25) train acc: 0.468000; val_acc: 0.323000
Epoch 7 / 25) train acc: 0.494000; val_acc: 0.287000
Epoch 8 / 25) train acc: 0.566000; val_acc: 0.328000
Epoch 9 / 25) train acc: 0.570000; val_acc: 0.320000
Epoch 10 / 25) train acc: 0.622000; val_acc: 0.324000
Epoch 11 / 25) train acc: 0.670000; val_acc: 0.278000
Epoch 12 / 25) train acc: 0.718000; val_acc: 0.338000
Epoch 13 / 25) train acc: 0.746000; val_acc: 0.318000
Epoch 14 / 25) train acc: 0.792000; val_acc: 0.307000
Epoch 15 / 25) train acc: 0.834000; val_acc: 0.297000
Epoch 16 / 25) train acc: 0.876000; val_acc: 0.327000
Epoch 17 / 25) train acc: 0.886000; val_acc: 0.320000
Epoch 18 / 25) train acc: 0.918000; val_acc: 0.314000
Epoch 19 / 25) train acc: 0.922000; val_acc: 0.296000
Epoch 20 / 25) train acc: 0.848000; val_acc: 0.306000
Epoch 21 / 25) train acc: 0.948000; val_acc: 0.302000
Epoch 22 / 25) train acc: 0.978000; val_acc: 0.302000
Epoch 23 / 25) train acc: 0.976000; val_acc: 0.289000
Epoch 24 / 25) train acc: 0.986000; val_acc: 0.285000
Epoch 25 / 25) train acc: 0.978000; val_acc: 0.311000
(iteration 1 / 125) loss: 2.29816
Epoch 0 / 25) train acc: 0.320000; val_acc: 0.146000
Epoch 1 / 25) train acc: 0.118000; val_acc: 0.131000
Epoch 2 / 25) train acc: 0.220000; val_acc: 0.214000
Epoch 3 / 25) train acc: 0.206000; val_acc: 0.180000
Epoch 4 / 25) train acc: 0.220000; val_acc: 0.193000
Epoch 5 / 25) train acc: 0.264000; val_acc: 0.239000
Epoch 6 / 25) train acc: 0.268000; val_acc: 0.203000
Epoch 7 / 25) train acc: 0.266000; val_acc: 0.212000
Epoch 8 / 25) train acc: 0.282000; val_acc: 0.256000
Epoch 9 / 25) train acc: 0.310000; val_acc: 0.250000
Epoch 10 / 25) train acc: 0.320000; val_acc: 0.267000
Epoch 11 / 25) train acc: 0.338000; val_acc: 0.273000
Epoch 12 / 25) train acc: 0.346000; val_acc: 0.278000
Epoch 13 / 25) train acc: 0.332000; val_acc: 0.273000
Epoch 14 / 25) train acc: 0.328000; val_acc: 0.284000
Epoch 15 / 25) train acc: 0.354000; val_acc: 0.271000
Epoch 16 / 25) train acc: 0.386000; val_acc: 0.277000
Epoch 17 / 25) train acc: 0.388000; val_acc: 0.297000
Epoch 18 / 25) train acc: 0.402000; val_acc: 0.280000
Epoch 19 / 25) train acc: 0.388000; val_acc: 0.274000
Epoch 20 / 25) train acc: 0.386000; val_acc: 0.274000
(iteration 101 / 125) loss: 1.313649
Epoch 21 / 25) train acc: 0.402000; val_acc: 0.272000
Epoch 22 / 25) train acc: 0.440000; val_acc: 0.286000
Epoch 23 / 25) train acc: 0.458000; val_acc: 0.295000
Epoch 24 / 25) train acc: 0.462000; val_acc: 0.310000
Epoch 25 / 25) train acc: 0.446000; val_acc: 0.297000
```

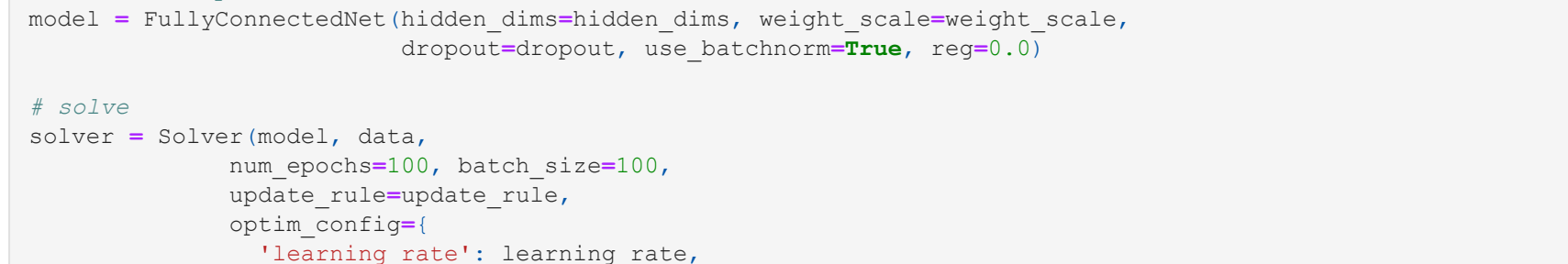
```
In [7]: # Plot train and validation accuracies of the two models

train_accs = []
val_accs = []
for dropout in dropout_choices:
    solver = solvers[dropout]
    train_accs.append(solver.train_acc_history[-1])
    val_accs.append(solver.val_acc_history[-1])

plt.subplot(3, 1, 1)
for dropout in dropout_choices:
    plt.plot(solvers[dropout].train_acc_history, 'o', label='%2f dropout' % dropout)
plt.title('Train accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend(ncol=2, loc='lower right')

plt.subplot(3, 1, 2)
for dropout in dropout_choices:
    plt.plot(solvers[dropout].val_acc_history, 'o', label='%2f dropout' % dropout)
plt.title('Val accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend(ncol=2, loc='lower right')

plt.gcf().set_size_inches(15, 15)
plt.show()
```



Question

Based off the results of this experiment, is dropout performing regularization? Explain your answer.

Answer:

Yes, the dropout layer does perform regularization. Although the model with and without dropout have similar validation accuracies, the model without dropout (blue) has significantly higher training accuracy compared to the model with dropout (orange). This means that the additional training accuracy the model without dropout has is overfitting and the model with dropout, in fact, regularize it. Also, dropout can be thought of regularizing each hidden unit to work well in many different contexts.

Final part of the assignment

Get over 55% validation accuracy on CIFAR-10 by using the layers you have implemented. You will be graded according to the following equation:

$$\min(\text{floor}((X - 32\%) / 28\%, 1))$$

where if you get 60% or higher validation accuracy, you get full points.

```
In [9]: # ===== #
# YOUR CODE HERE:
# Implement a FC-net that achieves at least 55% validation accuracy
# on CIFAR-10.
# ===== #

hidden_dims = [600, 600, 600, 600]
learning_rate = 2e-3
weight_scale = 0.01
lr_decay = 0.95
dropout = 0.55
update_rule = 'adam'

# create FullyConnectedNet
model = FullyConnectedNet(hidden_dims=hidden_dims, weight_scale=weight_scale,
                          dropout=dropout, use_batchnorm=True, reg=0.0)

# solve
solver = Solver(model, data,
                num_epochs=100, batch_size=100,
                update_rule=update_rule,
                optim_config={
                    'learning_rate': learning_rate,
                    'lr_decay': lr_decay,
                },
                verbose=True, print_every=100)
solver.train()

# print out the validation accuracy
y_test_pred = np.argmax(model.loss(data['X_test']), axis=1)
y_val_pred = np.argmax(model.loss(data['X_val']), axis=1)
print('Validation set accuracy: {}'.format(np.mean(y_val_pred == data['y_val'])))
print('Test set accuracy: {}'.format(np.mean(y_test_pred == data['y_test'])))

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


[illegible]


```

import numpy as np

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

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

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

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

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

    x_res = x.reshape((x.shape[0], w.shape[0])) # Shape of N * D
    out = x_res.dot(w) + b.reshape((1, b.shape[0])) # Shape of N * M

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

    cache = (x, w, b, out)
    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: 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:
    - 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.
    # Notice:
    #   - dx should be N x d_1 x ... x d_k; it relates to dout through multiplication with w, which is D x M
    #   - dw should be D x M; it is just the sum over dout examples
    #   - db should be M; it is just the sum over dout examples
    # ===== #

    x_res = np.reshape(x, (x.shape[0], w.shape[0]))
    dx = np.reshape(dx_res, (w.T))
    dw = np.reshape(dw_res, (w.shape)) # Shape of N * D
    db = x_res.T.dot(dout) # Shape of D * M
    db = dout.T.dot(np.ones(x.shape[0])) # Shape of 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, out)
    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 the same shape as dout

    Returns a tuple of:
    - 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 mean and variance using
    an exponential decay based on the momentum parameter:

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

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

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

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

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

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

    out, cache = None, None
    if mode == 'train':
        # ===== #
        # YOUR CODE HERE:
        # A few steps here:
        #   (1) Calculate the running mean and variance of the minibatch.
        #   (2) Normalize the activations with the running mean and variance.
        #   (3) Scale and shift the normalized activations. Store this
        #       as the variable 'out'
        #   (4) Store any variables you may need for the backward pass in
        #       the 'cache' variable.
        # ===== #

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

        mean_running = running_mean
        var_running = running_var

        mean_running = momentum * mean_running + (1 - momentum) * mean_minibatch
        var_running = momentum * var_running + (1 - momentum) * var_minibatch
        bn_param['running_mean'] = mean_running
        bn_param['running_var'] = var_running

        cache = (
            'minibatch_var': var_minibatch,
            'x_centralize': (x - mean_minibatch),
            '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
    # ===== #

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

# Store the updated running mean 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
    ddxat = dout * np.sqrt(minibatch_var + eps)
    sqrt_var = np.sqrt(minibatch_var + eps)
    dsqrt_var = np.sum(ddxat * x_centralize, axis=0) / (sqrt_var**2)
    dvar = dsqrt_var * 0.5 / sqrt_var
    dx1 = ddxat + dvar * x_centralize * dx * np.ones_like(dout) / N
    dx2 = np.sum(dx1, axis=0) * np.ones_like(dout) / N
    dx = dx1 + dx2

    # calculate dgamma 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):
    """
    Perform the backward pass for (inverted) dropout.

    Inputs:
    - dout: Upstream derivatives, of any shape
    - cache: (dropout_param, mask) from dropout_forward.

    """
    dropout_param, mask = cache
    mode = dropout_param['mode']

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

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

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

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

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

    return loss, dx

```

```

In [ ]:

import numpy as np
from layer_utils import *

class TwoLayerNet(object):
    """
    A two-layer fully-connected neural network with ReLU nonlinearity and
    softmax loss that uses a modular layer design. We assume an input dimension
    of D, a hidden dimension of H, and perform classification over C classes.

    The architecture should be affine - relu - affine - softmax.

    Note that this class does not implement data augmentation. In training, it
    will interact with a separate Solver object that is responsible for running
    optimization.

    The learnable parameters of the model are stored in the dictionary
    self.params that maps parameter names to numpy arrays.

    """
    def __init__(self, input_dim=3*32*32, hidden_dim=100, num_classes=10,
                 dropout=0, weight_scale=1e-3, reg=0):
        """
        Initialize a new network.

        Inputs:
        - input_dim: An integer giving the size of the input
        - hidden_dim: An integer giving the size of the hidden layer
        - num_classes: An integer giving the number of classes to classify
        - dropout: Scalar between 0 and 1 giving dropout strength.
        - weight_scale: Scalar giving the standard deviation for random
          initialization of the weights.
        - reg: Scalar giving L2 regularization strength.

        """
        self.params = {}
        self.reg = reg

        # YOUR CODE HERE:
        # Initialize W1, W2, b1, and b2. Store these as self.params['W1'],
        # self.params['W2'], self.params['b1'], and self.params['b2']. The
        # biases are initialized to zero and the weights are initialized
        # so that each parameter has mean 0 and standard deviation weight_scale.
        # The dimensions of W1 should be (input_dim, hidden_dim) and the
        # dimensions of W2 should be (hidden_dim, num_classes)
        # ===== #

        W1_size = (input_dim, hidden_dim)
        W2_size = (hidden_dim, num_classes)

        self.params['W1'] = np.random.normal(loc=0.0, scale=weight_scale, size=W1_size)
        self.params['W2'] = np.random.normal(loc=0.0, scale=weight_scale, size=W2_size)
        self.params['b1'] = np.zeros(hidden_dim)
        self.params['b2'] = np.zeros(num_classes)

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

    def loss(self, X, y=None):
        """
        Compute loss and gradient for a minibatch of data.

        Inputs:
        - X: Array of input data of shape (N, d_1, ..., d_k)
        - y: Array of labels, of shape (N,) y[i] gives the label for X[i].

        Returns:
        If y is None, then run a test-time forward pass of the model and return:
        - scores: Array of shape (N, C) giving classification scores, where
          scores[i, c] is the classification score for X[i] and class c.

        If y is not None, then run a training-time forward and backward pass and
        return a tuple of:
        - loss: Scalar value giving the loss
        - grads: Dictionary with the same keys as self.params, mapping parameter
          names to gradients of the loss with respect to those parameters.

        """
        scores = None

        # ===== #
        # YOUR CODE HERE:
        # Implement the forward pass of the two-layer neural network. Store
        # the class scores as the variable 'scores'. Be sure to use the layers
        # you prior implemented.
        # ===== #

        W1 = self.params['W1']
        b1 = self.params['b1']
        W2 = self.params['W2']
        b2 = self.params['b2']

        H, cache_h = affine_relu_forward(X, W1, b1)
        D, cache_d = affine_relu_forward(H, W2, b2)

        scores = Z

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

        # If y is None then we are in test mode so just return scores
        if y is None:
            return scores

        loss, grads = 0, {}

        # YOUR CODE HERE:
        # Implement the backward pass of the two-layer neural net. Store
        # the loss as the variable 'loss' and store the gradients in the
        # 'grads' dictionary. For the grads dictionary, grads['W1'] holds
        # the gradient for W1, grads['b1'] holds the gradient for b1, etc.
        # I.e., grads[k] holds the gradient for self.params[k].
        # Add L2 regularization, where there is an added cost 0.5*self.reg*W^2
        # for each W. Be sure to include the 0.5 multiplying factor.
        # And be sure to use the layers you prior implemented.
        # ===== #

        loss, dz = softmax_loss(scores, y)
        loss += 0.5*self.reg*(np.sum(W1**2) + np.sum(W2**2))

        dx, dw2, db2 = affine_backward(dz, cache_h)
        dx, dw1, db1 = affine_backward(dx, cache_d)

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

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

        return loss, grads

class FullyConnectedNet(object):
    """
    A fully-connected neural network with an arbitrary number of hidden layers,
    ReLU nonlinearities, and a softmax loss function. This will also implement
    dropout and batch normalization as options. For a network with L layers,
    the architecture will be

    (affine - [batch norm] - relu - [dropout]) x (L - 1) - affine - softmax

    where batch normalization and dropout are optional, and the (...) block is
    repeated L - 1 times.

    Similar to the TwoLayerNet above, learnable parameters are stored in the
    self.params dictionary and will be learned using the Solver class.

    """
    def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,
                 dropout=0, use_batchnorm=False, reg=0.0,
                 weight_scale=1e-2, dtype=np.float32, seed=None):
        """
        Initialize a new FullyConnectedNet.

        Inputs:
        - hidden_dims: A list of integers giving the size of each hidden layer.
        - num_classes: An integer giving the number of classes to classify.
        - dropout: Scalar between 0 and 1 giving dropout strength. If dropout=0 then
          the network should not use dropout at all.
        - use_batchnorm: Whether or not the network should use batch normalization.
        - reg: Scalar giving L2 regularization strength.
        - weight_scale: Scalar giving the standard deviation for random
          initialization of the weights.
        - dtype: A numpy datatype object, all computations will be performed using
          float64 for numeric gradient checking.
        - seed: If not None, this specifies a random seed to use to initialize the
          weights. For reproducibility, this is used to initialize the weights. This
          will make the dropout layers deterministic so we can gradient check the
          model.

        """
        self.use_batchnorm = use_batchnorm
        self.use_dropout = dropout > 0
        self.reg = reg
        self.num_layers = 1 + len(hidden_dims)
        self.dtype = dtype
        self.params = {}

        # ===== #
        # YOUR CODE HERE:
        # Initialize all parameters of the network in the self.params dictionary.
        # The weights and biases of layer l are Wl and bl, and in general the
        # weights should be gamma_l and beta_l. For layer 1, the weights are W1 and b1,
        # and the biases are gamma1 and beta1, etc. Only use batchnorm if self.use_batchnorm
        # is true and DO NOT do batch normalization the output scores.
        # ===== #

        cur_dim = input_dim
        for idx, hidden_dim in enumerate(hidden_dims):
            self.params['W'+str(idx+1)] = np.random.randn(cur_dim, hidden_dim) * weight_scale
            self.params['b'+str(idx+1)] = np.zeros(hidden_dim)
            self.params['gamma'+str(idx+1)] = np.ones(hidden_dim)
            self.params['beta'+str(idx+1)] = np.zeros(hidden_dim)
            cur_dim = hidden_dim

        self.params['W'+str(self.num_layers)] = np.random.randn(cur_dim, num_classes) * weight_scale
        self.params['b'+str(self.num_layers)] = np.zeros(num_classes)

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

        # When using dropout we need to pass a dropout_param dictionary to each
        # layer (train / test). You can pass the same dropout_param to each dropout layer.
        self.dropout_param = {}
        if self.use_dropout:
            self.dropout_param = {'mode': 'train', 'p': dropout}
        if seed is not None:
            self.dropout_param['seed'] = seed

        # With batch normalization we need to keep track of running means and
        # variances, so we need to pass a special bn_param object to each batch
        # normalization layer. You should pass self.bn_params[0] to the forward pass
        # of the first batch normalization layer, self.bn_params[1] to the forward
        # pass of the second batch normalization layer, etc.
        self.bn_params = []
        if self.use_batchnorm:
            self.dropout_param['bn_param'] = self.bn_params

        # Set all parameters to 'train' for the current layer
        for k in self.params:
            self.params[k] = v.astype(dtype)

    def loss(self, X, y=None):
        """
        Compute loss and gradient for the fully-connected net.

        Input / output: Same as TwoLayerNet above.

        """
        X = X.astype(self.dtype)
        mode = 'test' if y is None else 'train'

        # Set train/test mode for batchnorm params and dropout param since they
        # behave differently during training and testing.
        if self.dropout_param:
            self.dropout_param['mode'] = mode

        if self.use_batchnorm:
            for k in self.params:
                bn_param = self.bn_params

        scores = None

        # ===== #
        # YOUR CODE HERE:
        # Implement the forward pass of the FC net and store the output
        # scores as the variable "scores".
        # ===== #

        # BATCHNORM: If self.use_batchnorm is true, insert a batchnorm layer
        # between the affine forward and relu forward layers. You may
        # also write an affine_batchnorm_relu function in layer_utils.py.
        # ===== #

        # DROPOUT: If dropout is non-zero, insert a dropout layer after
        # every ReLU layer.
        # ===== #

        fc_cache = {}
        relu_cache = {}
        dropout_cache = {}

        # flatten image
        X = np.reshape(X, (X.shape[0], -1))

        # go through all layers
        for i in range(self.num_layers - 1):
            # fc layer
            fc_output, fc_cache[str(i+1)] = affine_forward(X, self.params['W'+str(i+1)], self.params['b'+str(i+1)])

            # batchnorm layer
            relu_input = fc_output
            if self.use_batchnorm:
                relu_output, batchnorm_cache[str(i+1)] = batchnorm_forward(fc_output, self.params['gamma'+str(i+1)], self.params['beta'+str(i+1)], self.bn_params[i])
            else:
                relu_output, relu_cache[str(i+1)] = relu_forward(relu_input)

            # dropout layer
            if self.use_dropout:
                relu_out, dropout_cache[str(i+1)] = dropout_forward(relu_out, self.dropout_param)
            else:
                relu_out = relu_output.copy()

            # update X
            X = relu_out.copy()

        scores, final_cache = affine_forward(X, self.params['W'+str(self.num_layers)], self.params['b'+str(self.num_layers)])

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

        # If test mode return early
        if y is None:
            return scores

        loss, grads = 0, {}

        # YOUR CODE HERE:
        # Implement the backwards pass of the FC net and store the gradients
        # in the grads dict, so that grads[k] is the gradient of self.params[k]
        # Be sure your L2 regularization includes a 0.5 factor.
        # ===== #

        # BATCHNORM: Incorporate the backward pass of the batchnorm.
        # ===== #

        # DROPOUT: Incorporate the backward pass of dropout.
        # ===== #

        loss, dx = softmax_loss(scores, y)
        loss += 0.5 * self.reg * (np.sum(np.square(self.params['W'+str(self.num_layers)])))
        dx_back, dx_fc, db_back, db_fc = affine_backward(dx, final_cache)
        grads['W'+str(self.num_layers)] = dx_back + self.reg * self.params['W'+str(self.num_layers)]
        grads['b'+str(self.num_layers)] = db_back

        # go backward all layers and update weights, bias, gammas and betas
        for i in range(self.num_layers - 1, 0, -1):
            # dropout layer
            dx_reli = relu_backward(dx_back, dropout_cache[str(i)])
            # batchnorm layer
            affine_backward_input = dx_reli
            if self.use_batchnorm:
                dx_bn, dgamma, dbeta = batchnorm_backward(dx_reli, batchnorm_cache[str(i)])
                grads['gamma'+str(i)] = dgamma
                grads['beta'+str(i)] = dbeta
            else:
                affine_backward_input = dx_bn
            dx_back, dx_fc, db_back, db_fc = affine_backward(affine_backward_input, fc_cache[str(i)])
            grads['W'+str(i)] = dx_fc + self.reg * self.params['W'+str(i)]
            grads['b'+str(i)] = db_fc
            loss += 0.5 * self.reg * (np.sum(np.square(self.params['W'+str(i)])))

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

        return loss, grads

```



```
import numpy as np
```

```
"""
```

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

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

```
Inputs:
```

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

```
Returns:
```

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

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

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

```
"""
```

```
def sgd(w, dw, config=None):  
    """
```

Performs vanilla stochastic gradient descent.

```
    config format:
```

- learning_rate: Scalar learning rate.

```
    if config is None: config = {}  
    config.setdefault('learning_rate', 1e-2)
```

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

```
def sgd_momentum(w, dw, config=None):  
    """
```

Performs stochastic gradient descent with momentum.

```
    config format:
```

- learning_rate: Scalar learning rate.
- momentum: Scalar between 0 and 1 giving the momentum value. Setting momentum = 0 reduces to sgd.
- velocity: A numpy array of the same shape as w and dw used to store a moving average of the gradients.

```
    if config is None: config = {}  
    config.setdefault('learning_rate', 1e-2)  
    config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there  
    v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.
```

```
    # =====  
    # YOUR CODE HERE:  
    # Implement the momentum update formula. Return the updated weights  
    # as next_w, and the updated velocity as v.  
    # =====
```

```
    v = config['momentum'] * v + config['learning_rate'] * dw  
    next_w = v + w
```

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

```
    config['velocity'] = v
```

```
    return next_w, config
```

```
def sgd_nesterov_momentum(w, dw, config=None):  
    """
```

Performs stochastic gradient descent with Nesterov momentum.

```
    config format:
```

- learning_rate: Scalar learning rate.
- momentum: Scalar between 0 and 1 giving the momentum value.
- decay_rate: Scalar between 0 and 1 reducing to sgd.
- velocity: A numpy array of the same shape as w and dw used to store a moving average of the gradients.

```
    if config is None: config = {}  
    config.setdefault('learning_rate', 1e-2)  
    config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there  
    v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.
```

```
    # =====  
    # YOUR CODE HERE:  
    # Implement the momentum update formula. Return the updated weights  
    # as next_w, and the updated velocity as v.  
    # =====
```

```
    v_0 = v  
    v = config['momentum'] * v + config['learning_rate'] * dw  
    w += v + config['momentum'] * (v - v_0)  
    next_w = w
```

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

```
    config['velocity'] = v
```

```
    return next_w, config
```

```
def rmsprop(w, dw, config=None):  
    """
```

Uses the RMSProp update rule, which uses a moving average of squared gradient values to set adaptive per-parameter learning rates.

```
    config format:
```

- learning_rate: Scalar learning rate.
- decay_rate: Scalar between 0 and 1 giving the decay rate for the squared gradient cache.
- epsilon: Small scalar used for smoothing to avoid dividing by zero.
- beta: Moving average of second moments of gradients.

```
    if config is None: config = {}  
    config.setdefault('learning_rate', 1e-2)  
    config.setdefault('decay_rate', 0.99)  
    config.setdefault('epsilon', 1e-8)  
    config.setdefault('beta', np.zeros_like(w))
```

```
    next_w = None
```

```
    # =====  
    # YOUR CODE HERE:  
    # Implement RMSProp. Store the next value of w as next_w. You need  
    # to also store in config['s'] the moving average of the second  
    # moment gradients, so they can be used for future gradients. Concretely,  
    # config['s'] corresponds to "s" in the lecture notes.  
    # =====
```

```
    config['s'] = config['decay_rate'] * config['s'] + (1 - config['decay_rate']) * (dw**2)  
    next_w = w - config['learning_rate'] * dw / (np.sqrt(config['s']) + config['epsilon'])
```

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

```
    return next_w, config
```

```
def adam(w, dw, config=None):  
    """
```

Uses the Adam update rule, which incorporates moving averages of both the gradient and its square and a bias correction term.

```
    config format:
```

- learning_rate: Scalar learning rate.
- beta1: Decay rate for moving average of first moment of gradient.
- beta2: Decay rate for moving average of second moment of gradient.
- epsilon: Small scalar used for smoothing to avoid dividing by zero.
- m: Moving average of gradient.
- v: Moving average of squared gradient.
- t: Iteration number.

```
    """
```

```
    if config is None: config = {}  
    config.setdefault('learning_rate', 1e-3)  
    config.setdefault('beta1', 0.9)  
    config.setdefault('beta2', 0.999)  
    config.setdefault('epsilon', 1e-8)  
    config.setdefault('m', np.zeros_like(w))  
    config.setdefault('v', 0)
```

```
    next_w = None
```

```
    # =====  
    # YOUR CODE HERE:  
    # Implement Adam. Store the next value of w as next_w. You need  
    # to also store in config['s'] the moving average of the second  
    # moment gradients, and in config['v'] the moving average of the  
    # first moments. Finally, store in config['t'] the increasing time.  
    # =====
```

```
    beta1 = config['beta1']  
    beta2 = config['beta2']  
    t = config['t'] + 1  
  
    m = beta1 * config['m'] + (1 - beta1) * dw  
    s = beta2 * config['s'] + (1 - beta2) * (dw**2)  
    v_corrected = m / (1 - beta1**t)  
    a_corrected = s / (1 - beta2**t)  
    next_w = w - config['learning_rate'] * v_corrected / (np.sqrt(a_corrected) + config['epsilon'])
```

```
    config['v'] = v  
    config['s'] = s  
    config['t'] = t
```

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

```
    return next_w, config
```