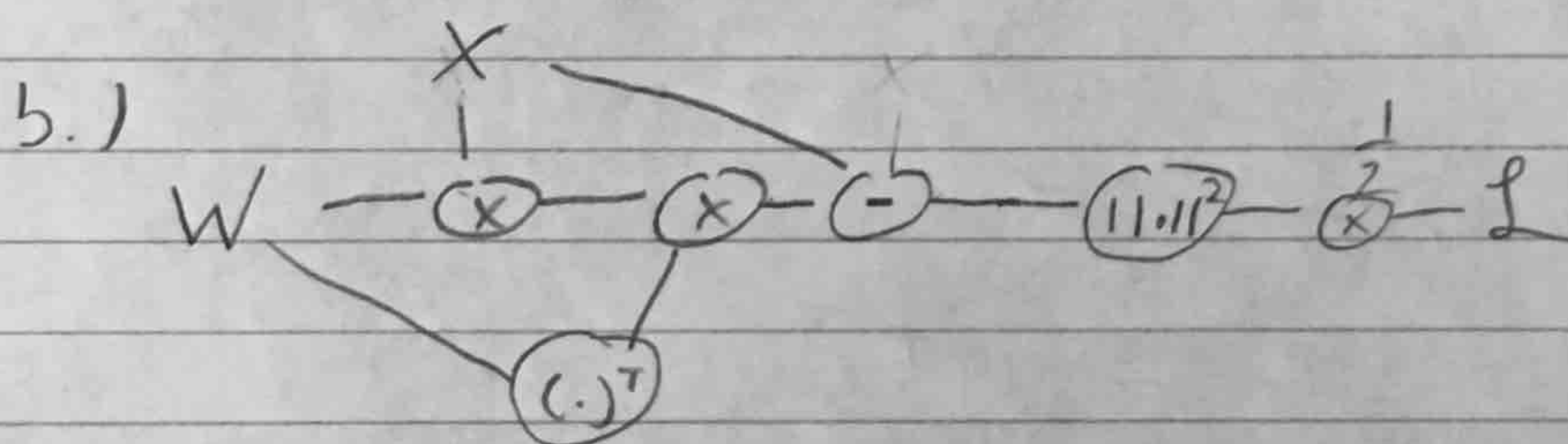


# ECE 239AS HW # 3

$$1. \mathcal{L} = \frac{1}{2} \|W^T W x - x\|^2$$

- a.) If  $W \in \mathbb{R}^{m \times n}$ , then  $W^T W \in \mathbb{R}^{n \times n}$ .  $Wx$  represents the dimensionally reduced input data.  $W^T W x$  then represents the reconstructed data. Minimizing the difference between this term and the original vector will result in a matrix  $W$  that retains information within  $W$ .



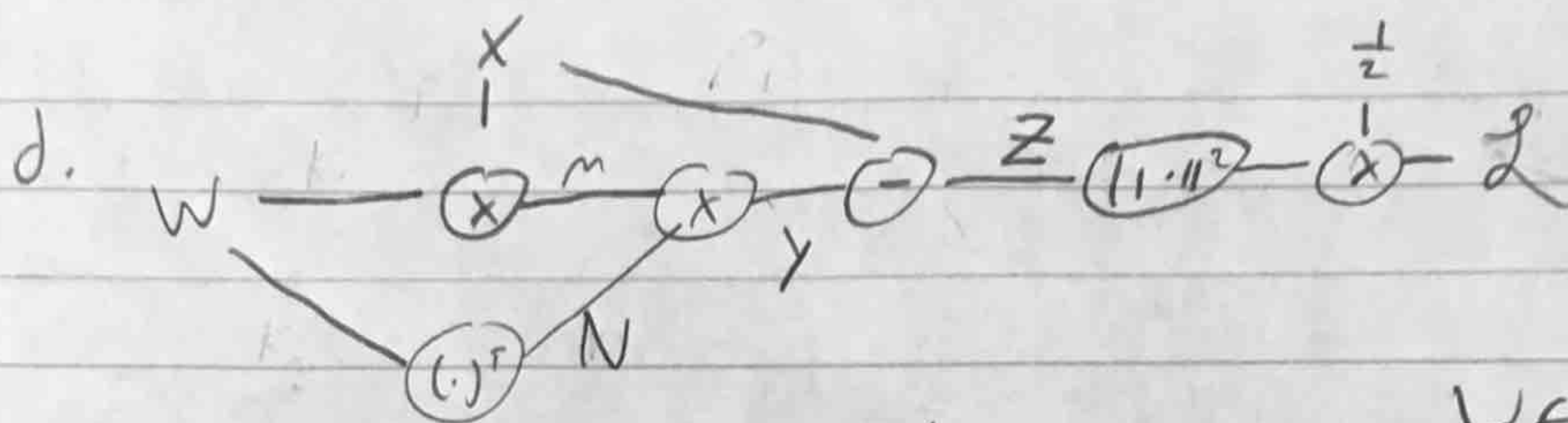
- c.) The gradient of  $\mathcal{L}$  w.r.t.  $W$  must be a summation of all paths leading from  $W$ . This is shown by the definition of a "total derivative" which is defined as:

$$\frac{df}{dt} = \frac{\partial f}{\partial t} + \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

where  $x$  and  $y$  are variables dependent on  $t$ .

In this case, the derivative of  $\mathcal{L}$  w.r.t.  $W$  must include all back propagation paths leading to it.





$$L = \frac{1}{2} \|z\|^2$$

$$z = y - x$$

$$y = Nm$$

$$m = Wx$$

$$N = W^T$$

$$\frac{\partial L}{\partial z} = z$$

$$\frac{\partial z}{\partial y} = 1$$

$$\frac{\partial y}{\partial m} = N^T, \frac{\partial y}{\partial N} = *N^T$$

$$\frac{\partial m}{\partial x} = *x^T$$

$$\frac{\partial L}{\partial W} = (*)^T$$

$$W \in \mathbb{R}^{m \times n}$$

$$x \in \mathbb{R}^n$$

$$N = W^T \in \mathbb{R}^{n \times m}$$

$$z \in \mathbb{R}^n$$

$$y \in \mathbb{R}^n$$

$$m \in \mathbb{R}^m$$

$$\frac{\partial L}{\partial z} = z$$

$$\frac{\partial L}{\partial y} = \frac{\partial z}{\partial y} \frac{\partial L}{\partial z} = z$$

$$\frac{\partial L}{\partial m} = \frac{\partial y}{\partial m} \frac{\partial L}{\partial y} = N^T z$$

$$\frac{\partial L}{\partial W_1} = \frac{\partial m}{\partial W_1} \frac{\partial L}{\partial m} = N^T z \cdot x^T$$

$$\frac{\partial L}{\partial N} = \frac{\partial y}{\partial N} \frac{\partial L}{\partial y} = z \cdot m^T$$

$$\frac{\partial L}{\partial W_2} = \frac{\partial N}{\partial W} \frac{\partial L}{\partial N} = (z m^T)^T = m z^T$$

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial W_1} + \frac{\partial L}{\partial W_2} = N^T z \cdot x^T + m z^T$$

$$= W(y-x)x^T + Wx(y-x)^T$$

$$= W(Nm-x)x^T + Wx(Nm-x)^T$$

$$\nabla_m L = W(W^T W x - x)x^T + Wx(W^T W x - x)^T$$

## This is the 2-layer neural network workbook for ECE 239AS Assignment #3

Please follow the notebook linearly to implement a two layer neural network.

Please print out the workbook entirely when completed.

We thank Serena Yeung & Justin Johnson for permission to use code written for the CS 231n class (cs231n.stanford.edu). These are the functions in the cs231n folders and code in the jupyter notebook to preprocess and show the images. The classifiers used are based off of code prepared for CS 231n as well.

The goal of this workbook is to give you experience with training a two layer neural network.

```
In [126]: import random
import numpy as np
from cs231n.data_utils import load_CIFAR10
import matplotlib.pyplot as plt

%matplotlib inline
%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))))
```

The autoreload extension is already loaded. To reload it, use:  
%reload\_ext autoreload

## Toy example

Before loading CIFAR-10, there will be a toy example to test your implementation of the forward and backward pass

```
In [127]: from nndl.neural_net import TwoLayerNet
```

```
In [158]: # Create a small net and some toy data to check your implementations.  
# Note that we set the random seed for repeatable experiments.  
  
input_size = 4  
hidden_size = 10  
num_classes = 3  
num_inputs = 5  
  
def init_toy_model():  
    np.random.seed(0)  
    return TwoLayerNet(input_size, hidden_size, num_classes, std=1e-1)  
  
def init_toy_data():  
    np.random.seed(1)  
    X = 10 * np.random.randn(num_inputs, input_size)  
    y = np.array([0, 1, 2, 2, 1])  
    return X, y  
  
net = init_toy_model()  
X, y = init_toy_data()
```

## Compute forward pass scores

```
In [159]: ## Implement the forward pass of the neural network.

# Note, there is a statement if y is None: return scores, which is wh
y
# the following call will calculate the scores.
scores = net.loss(X)
print('Your scores:')
print(scores)
print()
print('correct scores:')
correct_scores = np.asarray([
    [-1.07260209,  0.05083871, -0.87253915],
    [-2.02778743, -0.10832494, -1.52641362],
    [-0.74225908,  0.15259725, -0.39578548],
    [-0.38172726,  0.10835902, -0.17328274],
    [-0.64417314, -0.18886813, -0.41106892]])
print(correct_scores)
print()

# The difference should be very small. We get < 1e-7
print('Difference between your scores and correct scores:')
print(np.sum(np.abs(scores - correct_scores)))
```

```
Your scores:
[[-1.07260209  0.05083871 -0.87253915]
 [-2.02778743 -0.10832494 -1.52641362]
 [-0.74225908  0.15259725 -0.39578548]
 [-0.38172726  0.10835902 -0.17328274]
 [-0.64417314 -0.18886813 -0.41106892]]
()
correct scores:
[[-1.07260209  0.05083871 -0.87253915]
 [-2.02778743 -0.10832494 -1.52641362]
 [-0.74225908  0.15259725 -0.39578548]
 [-0.38172726  0.10835902 -0.17328274]
 [-0.64417314 -0.18886813 -0.41106892]]
()
Difference between your scores and correct scores:
3.381231222787662e-08
```

## Forward pass loss

```
In [160]: loss, _ = net.loss(X, y, reg=0.05)
correct_loss = 1.071696123862817

# should be very small, we get < 1e-12
print('Difference between your loss and correct loss:')
print(np.sum(np.abs(loss - correct_loss)))
```

```
Difference between your loss and correct loss:
0.0
```

```
In [161]: print(loss)
1.071696123862817
```

## Backward pass

Implements the backwards pass of the neural network. Check your gradients with the gradient check utilities provided.

```
In [162]: from cs231n.gradient_check import eval_numerical_gradient

# Use numeric gradient checking to check your implementation of the b
ackward pass.
# If your implementation is correct, the difference between the numer
ic and
# analytic gradients should be less than 1e-8 for each of W1, W2, b1,
and b2.

loss, grads = net.loss(X, y, reg=0.05)

# these should all be less than 1e-8 or so
for param_name in grads:
    f = lambda W: net.loss(X, y, reg=0.05)[0]
    param_grad_num = eval_numerical_gradient(f,
net.params[param_name], verbose=False)
    print('{} max relative error: {}'.format(param_name, rel_error(pa
ram_grad_num, grads[param_name])))

b2 max relative error: 1.83913010442e-10
b1 max relative error: 3.1726800927e-09
W1 max relative error: 1.28328233376e-09
W2 max relative error: 2.9632227682e-10
```

## Training the network

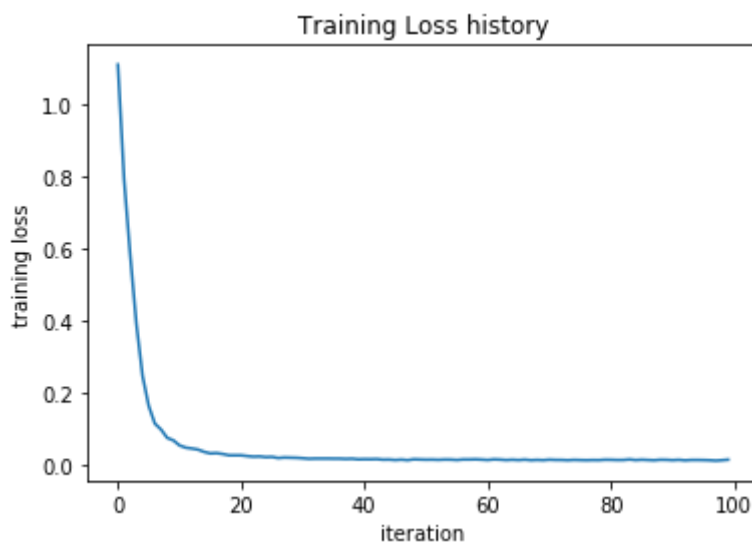
Implement `neural_net.train()` to train the network via stochastic gradient descent, much like the softmax and SVM.

```
In [163]: net = init_toy_model()
stats = net.train(X, y, X, y,
                  learning_rate=1e-1, reg=5e-6,
                  num_iters=100, verbose=False)

print('Final training loss: ', stats['loss_history'][-1])

# plot the loss history
plt.plot(stats['loss_history'])
plt.xlabel('iteration')
plt.ylabel('training loss')
plt.title('Training Loss history')
plt.show()

('Final training loss: ', 0.014497864587765875)
```



## Classify CIFAR-10

Do classification on the CIFAR-10 dataset.

```

In [164]: from cs231n.data_utils import load_CIFAR10

def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000):
    """
    Load the CIFAR-10 dataset from disk and perform preprocessing to prepare
    it for the two-layer neural net classifier. These are the same steps as
    we used for the SVM, but condensed to a single function.
    """
    # Load the raw CIFAR-10 data
    cifar10_dir = 'cifar-10-batches-py'
    X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)

    # Subsample the data
    mask = list(range(num_training, num_training + num_validation))
    X_val = X_train[mask]
    y_val = y_train[mask]
    mask = list(range(num_training))
    X_train = X_train[mask]
    y_train = y_train[mask]
    mask = list(range(num_test))
    X_test = X_test[mask]
    y_test = y_test[mask]

    # Normalize the data: subtract the mean image
    mean_image = np.mean(X_train, axis=0)
    X_train -= mean_image
    X_val -= mean_image
    X_test -= mean_image

    # Reshape data to rows
    X_train = X_train.reshape(num_training, -1)
    X_val = X_val.reshape(num_validation, -1)
    X_test = X_test.reshape(num_test, -1)

    return X_train, y_train, X_val, y_val, X_test, y_test

# Invoke the above function to get our data.
X_train, y_train, X_val, y_val, X_test, y_test = get_CIFAR10_data()
print('Train data shape: ', X_train.shape)
print('Train labels shape: ', y_train.shape)
print('Validation data shape: ', X_val.shape)
print('Validation labels shape: ', y_val.shape)
print('Test data shape: ', X_test.shape)
print('Test labels shape: ', y_test.shape)

('Train data shape: ', (49000, 3072))
('Train labels shape: ', (49000,))
('Validation data shape: ', (1000, 3072))
('Validation labels shape: ', (1000,))
('Test data shape: ', (1000, 3072))
('Test labels shape: ', (1000,))

```



## Running SGD

If your implementation is correct, you should see a validation accuracy of around 28-29%.

```
In [165]: input_size = 32 * 32 * 3
hidden_size = 50
num_classes = 10
net = TwoLayerNet(input_size, hidden_size, num_classes)

# Train the network
stats = net.train(X_train, y_train, X_val, y_val,
                  num_iters=1000, batch_size=200,
                  learning_rate=1e-4, learning_rate_decay=0.95,
                  reg=0.25, verbose=True)

# Predict on the validation set
val_acc = (net.predict(X_val) == y_val).mean()
print('Validation accuracy: ', val_acc)

# Save this net as the variable subopt_net for later comparison.
subopt_net = net

iteration 0 / 1000: loss 2.30275751861
iteration 100 / 1000: loss 2.30212015921
iteration 200 / 1000: loss 2.29561360074
iteration 300 / 1000: loss 2.25182590432
iteration 400 / 1000: loss 2.18899523505
iteration 500 / 1000: loss 2.11625277919
iteration 600 / 1000: loss 2.0646708277
iteration 700 / 1000: loss 1.99016886231
iteration 800 / 1000: loss 2.00282764012
iteration 900 / 1000: loss 1.94651768179
('Validation accuracy: ', 0.283)
```

## Questions:

The training accuracy isn't great.

(1) What are some of the reasons why this is the case? Take the following cell to do some analyses and then report your answers in the cell following the one below.

(2) How should you fix the problems you identified in (1)?

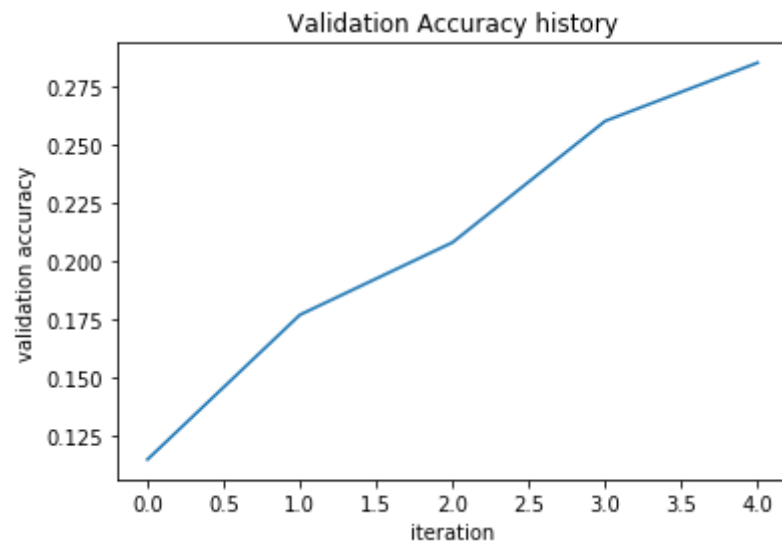
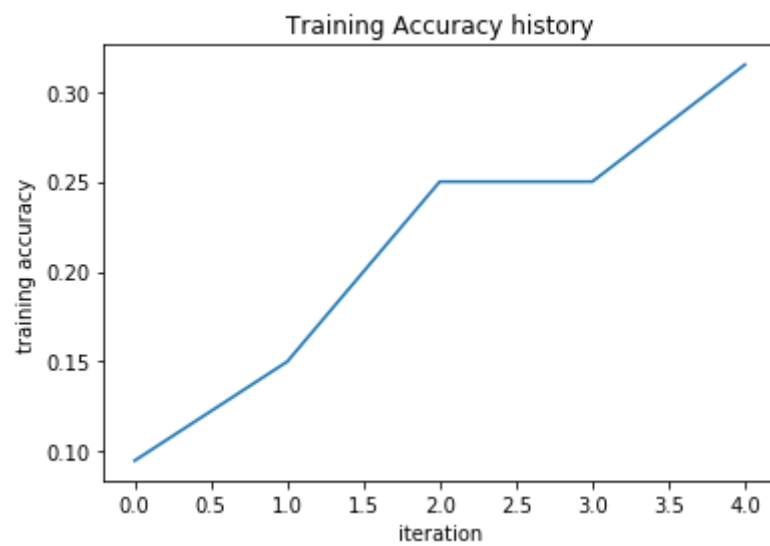
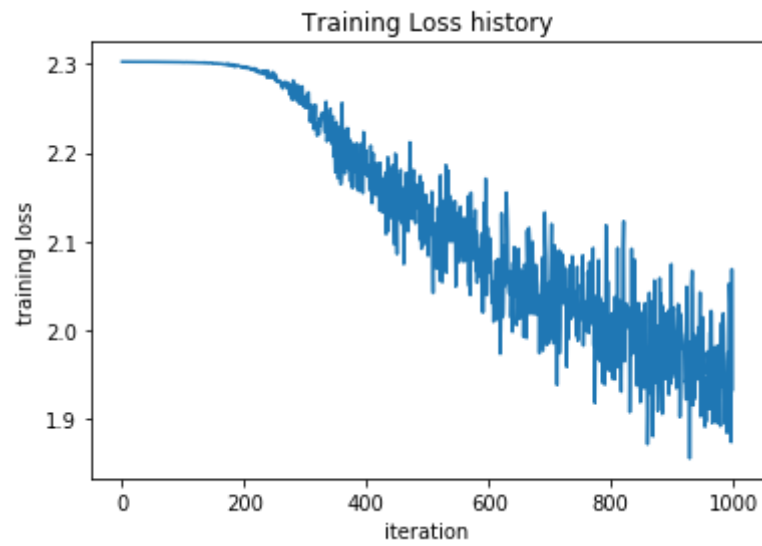
```
In [166]: stats['train_acc_history']

Out[166]: [0.095, 0.15, 0.25, 0.25, 0.315]
```

```
In [167]: # ===== #
# YOUR CODE HERE:
#   Do some debugging to gain some insight into why the optimization
#   isn't great.
# ===== #

# Plot the loss function and train / validation accuracies
plt.plot(stats['loss_history'])
plt.xlabel('iteration')
plt.ylabel('training loss')
plt.title('Training Loss history')
plt.show()
plt.plot(stats['train_acc_history'])
plt.xlabel('iteration')
plt.ylabel('training accuracy')
plt.title('Training Accuracy history')
plt.show()
plt.plot(stats['val_acc_history'])
plt.xlabel('iteration')
plt.ylabel('validation accuracy')
plt.title('Validation Accuracy history')
plt.show()

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



## Answers:

(1) All of the histories still have a slope to them, so it seems that additional accuracy could be gained with additional training. The validation accuracy and training accuracy are more or less in step, which suggest that the model is underfitting or not at full capacity.

(2) The NN needs to be trained further, which can be accomplished by either increasing the learning rates or increasing the number of iterations. The capacity could also be increased by increasing the model complexity. Since by the instructions of the HW we cannot do the third, increasing the learning rate and number of iterations will be a good first step.

## Optimize the neural network

Use the following part of the Jupyter notebook to optimize your hyperparameters on the validation set. Store your nets as `best_net`.



```

In [154]: best_net = None # store the best model into this

# ===== #
# YOUR CODE HERE:
# Optimize over your hyperparameters to arrive at the best neural
# network. You should be able to get over 50% validation accuracy.
# For this part of the notebook, we will give credit based on the
# accuracy you get. Your score on this question will be multiplied
by:
# min(floor((X - 28%)) / %22, 1)
# where if you get 50% or higher validation accuracy, you get full
# points.
#
# Note, you need to use the same network structure (keep hidden_size = 50)!
# ===== #

learning_rates = [1e-1, 5e-2, 1e-2, 5e-3, 1e-3]
learning_rates = [1e-3]
regularization_strengths = [0.1, 0.2, 0.3, 0.4, 0.5, 0.6]
best_score = -1
best_stats = None
accuracies = {}
iters = 2500

for rate in learning_rates:
    for strength in regularization_strengths:

        net = TwoLayerNet(input_size, hidden_size, num_classes)

        stats = net.train(X_train, y_train, X_val, y_val,
                           num_iters=iters, batch_size=200,
                           learning_rate=rate, learning_rate_decay=0.95,
                           reg=strength, verbose=False)

```

```
train_acc = (net.predict(X_train) == y_train).mean()
val_acc = (net.predict(X_val) == y_val).mean()

accuracies[(rate, strength)] = (train_acc, val_acc)

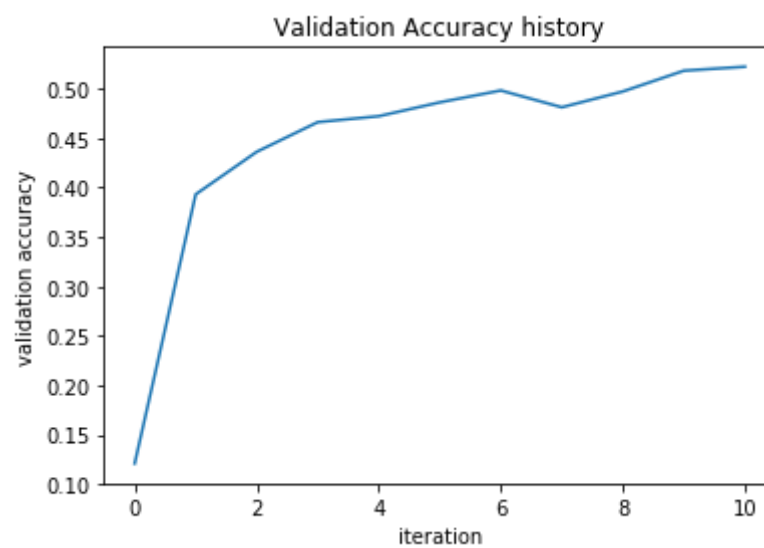
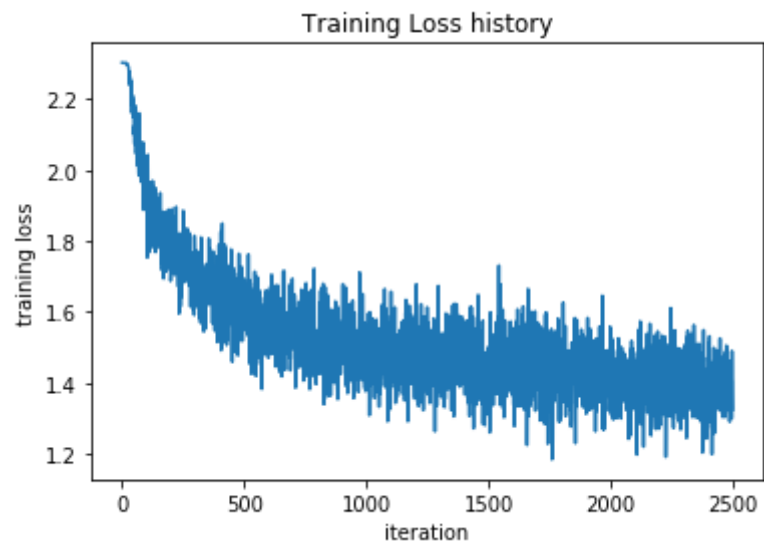
if val_acc > best_score:
    best_stats = stats
    best_score = val_acc
    best_net = net
```

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

```
In [168]: print best_score

plt.plot(best_stats['loss_history'])
plt.xlabel('iteration')
plt.ylabel('training loss')
plt.title('Training Loss history')
plt.show()
plt.plot(best_stats['train_acc_history'])
plt.xlabel('iteration')
plt.ylabel('training accuracy')
plt.title('Training Accuracy history')
plt.show()
plt.plot(best_stats['val_acc_history'])
plt.xlabel('iteration')
plt.ylabel('validation accuracy')
plt.title('Validation Accuracy history')
plt.show()
```

0.512



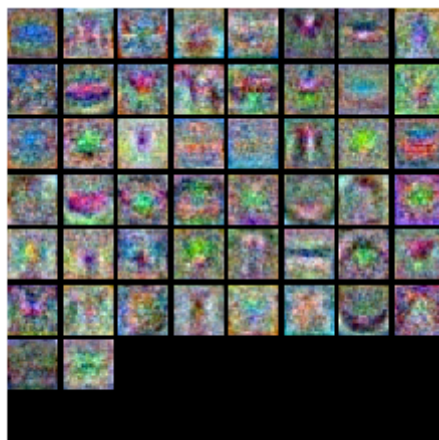
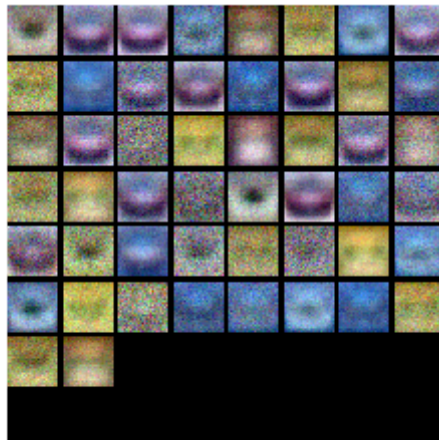


```
In [156]: from cs231n.vis_utils import visualize_grid

# Visualize the weights of the network

def show_net_weights(net):
    W1 = net.params['W1']
    W1 = W1.T.reshape(32, 32, 3, -1).transpose(3, 0, 1, 2)
    plt.imshow(visualize_grid(W1, padding=3).astype('uint8'))
    plt.gca().axis('off')
    plt.show()

show_net_weights(subopt_net)
show_net_weights(best_net)
```



## Question:

(1) What differences do you see in the weights between the suboptimal net and the best net you arrived at?

## Answer:

(1) The weights in the suboptimized one are taking more vague approximations of the objects in the pictures. You can see in the best net that certain features are being emphasized, and these can be assumed to be unique to the classes.

## Evaluate on test set

```
In [157]: test_acc = (best_net.predict(X_test) == y_test).mean()
          print('Test accuracy: ', test_acc)
          ('Test accuracy: ', 0.5)
```

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 """
5 This code was originally written for CS 231n at Stanford University
6 (cs231n.stanford.edu). It has been modified in various areas for use in the
7 ECE 239AS class at UCLA. This includes the descriptions of what code to
8 implement as well as some slight potential changes in variable names to be
9 consistent with class nomenclature. We thank Justin Johnson & Serena Yeung for
10 permission to use this code. To see the original version, please visit
11 cs231n.stanford.edu.
12 """
13
14 class TwoLayerNet(object):
15     """
16     A two-layer fully-connected neural network. The net has an input dimension of
17     N, a hidden layer dimension of H, and performs classification over C classes.
18     We train the network with a softmax loss function and L2 regularization on the
19     weight matrices. The network uses a ReLU nonlinearity after the first fully
20     connected layer.
21
22     In other words, the network has the following architecture:
23
24     input - fully connected layer - ReLU - fully connected layer - softmax
25
26     The outputs of the second fully-connected layer are the scores for each class.
27     """
28
29     def __init__(self, input_size, hidden_size, output_size, std=1e-4):
30         """
31         Initialize the model. Weights are initialized to small random values and
32         biases are initialized to zero. Weights and biases are stored in the
33         variable self.params, which is a dictionary with the following keys:
34
35         W1: First layer weights; has shape (H, D)
36         b1: First layer biases; has shape (H,)
37         W2: Second layer weights; has shape (C, H)
38         b2: Second layer biases; has shape (C,)
39
40         Inputs:
41         - input_size: The dimension D of the input data.
42         - hidden_size: The number of neurons H in the hidden layer.
43         - output_size: The number of classes C.
44         """
45         self.params = {}
46         self.params['W1'] = std * np.random.randn(hidden_size, input_size)
47         self.params['b1'] = np.zeros(hidden_size)
48         self.params['W2'] = std * np.random.randn(output_size, hidden_size)
49         self.params['b2'] = np.zeros(output_size)
50
51     def loss(self, X, y=None, reg=0.0):
52         """
53         Compute the loss and gradients for a two layer fully connected neural
54         network.
55
56         Inputs:
57         - X: Input data of shape (N, D). Each X[i] is a training sample.
58         - y: Vector of training labels. y[i] is the label for X[i], and each y[i] is
59             an integer in the range 0 <= y[i] < C. This parameter is optional; if it
60             is not passed then we only return scores, and if it is passed then we
61             instead return the loss and gradients.
62         - reg: Regularization strength.
63
64         Returns:
65         If y is None, return a matrix scores of shape (N, C) where scores[i, c] is
66             the score for class c on input X[i].
67
68         If y is not None, instead return a tuple of:

```

```

70 - loss: Loss (data loss and regularization loss) for this batch of training
71 samples.
72 - grads: Dictionary mapping parameter names to gradients of those parameters
73 with respect to the loss function; has the same keys as self.params.
74 """
75 # Unpack variables from the params dictionary
76 W1, b1 = self.params['W1'], self.params['b1']
77 W2, b2 = self.params['W2'], self.params['b2']
78 N, D = X.shape
79
80 # Compute the forward pass
81 scores = None
82
83 # =====
84 # YOUR CODE HERE:
85 # Calculate the output scores of the neural network. The result
86 # should be (C, N). As stated in the description for this class,
87 # there should not be a ReLU layer after the second FC layer.
88 # The output of the second FC layer is the output scores. Do not
89 # use a for loop in your implementation.
90 # =====
91
92 H1 = W1.dot(X.T) + b1[:, np.newaxis]
93 H1[H1<0] = 0
94 scores = W2.dot(H1) + b2[:, np.newaxis]
95 scores = scores.T
96
97 # =====
98 # END YOUR CODE HERE
99 # =====
100
101
102 # If the targets are not given then jump out, we're done
103 if y is None:
104     return scores
105
106 # Compute the loss
107 loss = None
108
109 # =====
110 # YOUR CODE HERE:
111 # Calculate the loss of the neural network. This includes the
112 # softmax loss and the L2 regularization for W1 and W2. Store the
113 # total loss in the variable loss. Multiply the regularization
114 # loss by 0.5 (in addition to the factor reg).
115 # =====
116
117 # scores is num_examples by num_classes
118 num_examples = scores.shape[0]
119
120 max_score = np.amax(scores, axis=1)
121 scores -= max_score[:, np.newaxis]
122
123 e_scores = np.exp(scores)
124 sums = np.sum(e_scores, axis=1)
125 log_sums = np.log(sums)
126 y_terms = scores[np.arange(num_examples), y]
127 loss = np.sum(log_sums - y_terms)/num_examples + .5*reg*np.sum(W1*W1) + .5*reg*np.sum(W2*W2)
128 # =====
129 # END YOUR CODE HERE
130 # =====
131
132 grads = {}
133
134 # =====
135 # YOUR CODE HERE:
136 # Implement the backward pass. Compute the derivatives of the
137 # weights and the biases. Store the results in the grads
138 # dictionary. e.g., grads['W1'] should store the gradient for

```



```

139 # W1, and be of the same size as W1.
140 # ===== #
141
142 #print W1.shape
143 #print H1.shape
144 #print scores.shape
145
146 d_scores = e_scores/sums[:,np.newaxis]
147 d_scores[np.arange(num_examples),y] -= 1
148 d_scores = d_scores.T/num_examples
149
150 b2_grad = np.sum(d_scores,axis=1)
151 W2_grad = d_scores.dot(H1.T)
152
153 r_grad = W2.T.dot(d_scores)
154 r_grad[H1<=0] = 0
155
156 b1_grad = np.sum(r_grad,axis=1)
157 W1_grad = r_grad.dot(X)
158
159 grads['b1'] = b1_grad
160 grads['W1'] = W1_grad + reg*W1
161
162 grads['b2'] = b2_grad
163 grads['W2'] = W2_grad + reg*W2
164
165 # ===== #
166 # END YOUR CODE HERE
167 # ===== #
168
169 return loss, grads
170
171 def train(self, X, y, X_val, y_val,
172         learning_rate=1e-3, learning_rate_decay=0.95,
173         reg=1e-5, num_iters=100,
174         batch_size=200, verbose=False):
175     """
176     Train this neural network using stochastic gradient descent.
177
178     Inputs:
179     - X: A numpy array of shape (N, D) giving training data.
180     - y: A numpy array of shape (N,) giving training labels; y[i] = c means that
181         X[i] has label c, where 0 <= c < C.
182     - X_val: A numpy array of shape (N_val, D) giving validation data.
183     - y_val: A numpy array of shape (N_val,) giving validation labels.
184     - learning_rate: Scalar giving learning rate for optimization.
185     - learning_rate_decay: Scalar giving factor used to decay the learning rate
186         after each epoch.
187     - reg: Scalar giving regularization strength.
188     - num_iters: Number of steps to take when optimizing.
189     - batch_size: Number of training examples to use per step.
190     - verbose: boolean; if true print progress during optimization.
191     """
192     num_train = X.shape[0]
193     iterations_per_epoch = max(num_train / batch_size, 1)
194
195     # Use SGD to optimize the parameters in self.model
196     loss_history = []
197     train_acc_history = []
198     val_acc_history = []
199
200     for it in np.arange(num_iters):
201         X_batch = None
202         y_batch = None
203
204         # ===== #
205         # YOUR CODE HERE:
206         # Create a minibatch by sampling batch_size samples randomly.
207         # ===== #

```

```

208 mask = np.random.choice(np.arange(X.shape[0]), batch_size)
209 X_batch = X[mask]
210 y_batch = y[mask]
211
212 # ===== #
213 # END YOUR CODE HERE
214 # ===== #
215
216 # Compute loss and gradients using the current minibatch
217 loss, grads = self.loss(X_batch, y=y_batch, reg=reg)
218 loss_history.append(loss)
219
220 # ===== #
221 # YOUR CODE HERE:
222 # Perform a gradient descent step using the minibatch to update
223 # all parameters (i.e., W1, W2, b1, and b2).
224 # ===== #
225
226 self.params['W1'] -= learning_rate*grads['W1']
227 self.params['b1'] -= learning_rate*grads['b1']
228 self.params['W2'] -= learning_rate*grads['W2']
229 self.params['b2'] -= learning_rate*grads['b2']
230
231
232 # ===== #
233 # END YOUR CODE HERE
234 # ===== #
235
236 if verbose and it % 100 == 0:
237     print('iteration {} / {}: loss {}'.format(it, num_iters, loss))
238
239 # Every epoch, check train and val accuracy and decay learning rate.
240 if it % iterations_per_epoch == 0:
241     # Check accuracy
242     train_acc = (self.predict(X_batch) == y_batch).mean()
243     val_acc = (self.predict(X_val) == y_val).mean()
244     train_acc_history.append(train_acc)
245     val_acc_history.append(val_acc)
246
247     # Decay learning rate
248     learning_rate *= learning_rate_decay
249
250 return {
251     'loss_history': loss_history,
252     'train_acc_history': train_acc_history,
253     'val_acc_history': val_acc_history,
254 }
255
256 def predict(self, X):
257     """
258     Use the trained weights of this two-layer network to predict labels for
259     data points. For each data point we predict scores for each of the C
260     classes, and assign each data point to the class with the highest score.
261
262     Inputs:
263     - X: A numpy array of shape (N, D) giving N D-dimensional data points to
264       classify.
265
266     Returns:
267     - y_pred: A numpy array of shape (N,) giving predicted labels for each of
268       the elements of X. For all i, y_pred[i] = c means that X[i] is predicted
269       to have class c, where 0 <= c < C.
270     """
271     y_pred = None
272
273     # ===== #
274     # YOUR CODE HERE:
275     # Predict the class given the input data.
276     # ===== #

```

```
277 W1, b1 = self.params['W1'], self.params['b1']
278 W2, b2 = self.params['W2'], self.params['b2']
279
280 N, D = X.shape
281 H1 = W1.dot(X.T) + b1[:, np.newaxis]
282 H1[H1<0] = 0
283 scores = W2.dot(H1) + b2[:, np.newaxis]
284
285 y_pred = np.argmax(scores, axis=0)
286
287
288 # ===== #
289 # END YOUR CODE HERE
290 # ===== #
291
292 return y_pred
```

## Fully connected networks

In the previous notebook, you implemented a simple two-layer neural network class. However, this class is not modular. If you wanted to change the number of layers, you would need to write a new loss and gradient function. If you wanted to optimize the network with different optimizers, you'd need to write new training functions. If you wanted to incorporate regularizations, you'd have to modify the loss and gradient function.

Instead of having to modify functions each time, for the rest of the class, we'll work in a more modular framework where we define forward and backward layers that calculate losses and gradients respectively. Since the forward and backward layers share intermediate values that are useful for calculating both the loss and the gradient, we'll also have these function return "caches" which store useful intermediate values.

The goal is that through this modular design, we can build different sized neural networks for various applications.

In this HW #3, we'll define the basic architecture, and in HW #4, we'll build on this framework to implement different optimizers and regularizations (like BatchNorm and Dropout).

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

## Modular layers

This notebook will build modular layers in the following manner. First, there will be a forward pass for a given layer with inputs ( $x$ ) and return the output of that layer ( $out$ ) as well as cached variables ( $cache$ ) that will be used to calculate the gradient in the backward pass.

```
def layer_forward(x, w):  
    """ Receive inputs  $x$  and weights  $w$  """  
    # Do some computations ...  
    z = # ... some intermediate value  
    # Do some more computations ...  
    out = # the output  
  
    cache = (x, w, z, out) # Values we need to compute gradients  
  
    return out, cache
```

The backward pass will receive upstream derivatives and the cache object, and will return gradients with respect to the inputs and weights, like this:

```
def layer_backward(dout, cache):  
    """  
    Receive derivative of loss with respect to outputs and cache,  
    and compute derivative with respect to inputs.  
    """  
    # Unpack cache values  
    x, w, z, out = cache  
  
    # Use values in cache to compute derivatives  
    dx = # Derivative of loss with respect to  $x$   
    dw = # Derivative of loss with respect to  $w$   
  
    return dx, dw
```

```
In [107]: ## Import and setups

import time
import numpy as np
import matplotlib.pyplot as plt
from nndl.fc_net import *
from cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
from cs231n.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))))
```

The autoreload extension is already loaded. To reload it, use:  
 %reload\_ext autoreload

```
In [108]: # Load the (preprocessed) CIFAR10 data.

data = get_CIFAR10_data()
for k in data.keys():
    print('{:}: {}'.format(k, data[k].shape))

X_val: (1000, 3, 32, 32)
X_train: (49000, 3, 32, 32)
X_test: (1000, 3, 32, 32)
y_val: (1000,)
y_train: (49000,)
y_test: (1000,)
```

## Linear layers

In this section, we'll implement the forward and backward pass for the linear layers.

The linear layer forward pass is the function `affine_forward` in `nndl/layers.py` and the backward pass is `affine_backward`.

After you have implemented these, test your implementation by running the cell below.

## Affine layer forward pass

Implement `affine_forward` and then test your code by running the following cell.

```
In [109]: # Test the affine_forward function

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

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

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

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

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

Testing affine_forward function:
difference: 9.76984946819e-10
```

## Affine layer backward pass

Implement `affine_backward` and then test your code by running the following cell.



```
In [110]: # Test the affine_backward function

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

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

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

# The error should be around 1e-10
print('Testing affine_backward function:')
print('dx error: {}'.format(rel_error(dx_num, dx)))
print('dw error: {}'.format(rel_error(dw_num, dw)))
print('db error: {}'.format(rel_error(db_num, db)))

Testing affine_backward function:
dx error: 8.60871575069e-11
dw error: 5.16849626491e-11
db error: 6.21942801469e-11
```

## Activation layers

In this section you'll implement the ReLU activation.

### ReLU forward pass

Implement the `relu_forward` function in `nndl/layers.py` and then test your code by running the following cell.

```
In [111]: # Test the relu_forward function

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

out, _ = relu_forward(x)
correct_out = np.array([[ 0.,          0.,          0.,          0.,
                          ],
                        [ 0.,          0.,          0.04545455,  0.13
636364, ],
                        [ 0.22727273,  0.31818182,  0.40909091,  0.5,
                          ]])

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

Testing relu_forward function:
difference: 4.99999979802e-08
```

## ReLU backward pass

Implement the `relu_backward` function in `nndl/layers.py` and then test your code by running the following cell.

```
In [112]: x = np.random.randn(10, 10)
dout = np.random.randn(*x.shape)

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

_, cache = relu_forward(x)
dx = relu_backward(dout, cache)
# The error should be around 1e-12
print('Testing relu_backward function:')
print('dx error: {}'.format(rel_error(dx_num, dx)))

Testing relu_backward function:
dx error: 3.27561396362e-12
```

## Combining the affine and ReLU layers

Often times, an affine layer will be followed by a ReLU layer. So let's make one that puts them together. Layers that are combined are stored in `nndl/layer_utils.py`.

## Affine-ReLU layers

We've implemented `affine_relu_forward()` and `affine_relu_backward` in `nndl/layer_utils.py`. Take a look at them to make sure you understand what's going on. Then run the following cell to ensure its implemented correctly.

```
In [113]: from nndl.layer_utils import affine_relu_forward, affine_relu_backward

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

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

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

print('Testing affine_relu_forward and affine_relu_backward:')
print('dx error: {}'.format(rel_error(dx_num, dx)))
print('dw error: {}'.format(rel_error(dw_num, dw)))
print('db error: {}'.format(rel_error(db_num, db)))

Testing affine_relu_forward and affine_relu_backward:
dx error: 1.47423906121e-08
dw error: 9.4703746119e-10
db error: 1.35072953084e-11
```

## Softmax and SVM losses

You've already implemented these, so we have written these in `layers.py`. The following code will ensure they are working correctly.

```

In [114]: num_classes, num_inputs = 10, 50
x = 0.001 * np.random.randn(num_inputs, num_classes)
y = np.random.randint(num_classes, size=num_inputs)

dx_num = eval_numerical_gradient(lambda x: svm_loss(x, y)[0], x, verbose=False)
loss, dx = svm_loss(x, y)

# Test svm_loss function. Loss should be around 9 and dx error should be 1e-9
print('Testing svm_loss:')
print('loss: {}'.format(loss))
print('dx error: {}'.format(rel_error(dx_num, dx)))

dx_num = eval_numerical_gradient(lambda x: softmax_loss(x, y)[0], x, verbose=False)
loss, dx = softmax_loss(x, y)

# Test softmax_loss function. Loss should be 2.3 and dx error should be 1e-8
print('\nTesting softmax_loss:')
print('loss: {}'.format(loss))
print('dx error: {}'.format(rel_error(dx_num, dx)))

Testing svm_loss:
loss: 9.00013162042
dx error: 3.0387355051e-09

Testing softmax_loss:
loss: 2.30259872626
dx error: 1.0323036143e-08

```

## Implementation of a two-layer NN

In `nndl/fc_net.py`, implement the class `TwoLayerNet` which uses the layers you made here. When you have finished, the following cell will test your implementation.

```
In [115]: N, D, H, C = 3, 5, 50, 7
X = np.random.randn(N, D)
y = np.random.randint(C, size=N)

std = 1e-2
model = TwoLayerNet(input_dim=D, hidden_dims=H, num_classes=C, weight
_scale=std)

print('Testing initialization ... ')
W1_std = abs(model.params['W1'].std() - std)
b1 = model.params['b1']
W2_std = abs(model.params['W2'].std() - std)
b2 = model.params['b2']
assert W1_std < std / 10, 'First layer weights do not seem right'
assert np.all(b1 == 0), 'First layer biases do not seem right'
assert W2_std < std / 10, 'Second layer weights do not seem right'
assert np.all(b2 == 0), 'Second layer biases do not seem right'
```

```

print('Testing test-time forward pass ... ')
model.params['W1'] = np.linspace(-0.7, 0.3, num=D*H).reshape(D, H)
model.params['b1'] = np.linspace(-0.1, 0.9, num=H)
model.params['W2'] = np.linspace(-0.3, 0.4, num=H*C).reshape(H, C)
model.params['b2'] = np.linspace(-0.9, 0.1, num=C)
X = np.linspace(-5.5, 4.5, num=N*D).reshape(D, N).T
scores = model.loss(X)
correct_scores = np.asarray(
    [[11.53165108, 12.2917344, 13.05181771, 13.81190102, 14.571984
34, 15.33206765, 16.09215096],
    [12.05769098, 12.74614105, 13.43459113, 14.1230412, 14.811491
28, 15.49994135, 16.18839143],
    [12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.050998
22, 15.66781506, 16.2846319 ]])
scores_diff = np.abs(scores - correct_scores).sum()
assert scores_diff < 1e-6, 'Problem with test-time forward pass'

print('Testing training loss (no regularization)')
y = np.asarray([0, 5, 1])
loss, grads = model.loss(X, y)
correct_loss = 3.4702243556
assert abs(loss - correct_loss) < 1e-10, 'Problem with training-time
loss'

model.reg = 1.0
loss, grads = model.loss(X, y)
correct_loss = 26.5948426952
assert abs(loss - correct_loss) < 1e-10, 'Problem with regularization
loss'

for reg in [0.0, 0.7]:
    print('Running numeric gradient check with reg = {}'.format(reg))
    model.reg = reg
    loss, grads = model.loss(X, y)

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

```

```

Testing initialization ...
Testing test-time forward pass ...
Testing training loss (no regularization)
Running numeric gradient check with reg = 0.0
W1 relative error: 1.83365627867e-08
W2 relative error: 3.11807423476e-10
b1 relative error: 9.82831520464e-09
b2 relative error: 4.32913495457e-10
Running numeric gradient check with reg = 0.7
W1 relative error: 2.52791523102e-07
W2 relative error: 7.97665280616e-08
b1 relative error: 1.34676189626e-08
b2 relative error: 7.75909428384e-10

```

## Solver

We will now use the `cs231n Solver` class to train these networks. Familiarize yourself with the API in `cs231n/solver.py`. After you have done so, declare an instance of a `TwoLayerNet` with 200 units and then train it with the `Solver`. Choose parameters so that your validation accuracy is at least 50%.

```
In [117]: model = TwoLayerNet()
          solver = None

          # ===== #
          # YOUR CODE HERE:
          #   Declare an instance of a TwoLayerNet and then train
          #   it with the Solver. Choose hyperparameters so that your validation
          #   accuracy is at least 40%. We won't have you optimize this further
          #   since you did it in the previous notebook.
          # ===== #

          model = TwoLayerNet(hidden_dims=200, reg=.25)

          solver = Solver(model, data, print_every=490,
                          batch_size=200,
                          lr_decay=.95,
                          optim_config = {
                              'learning_rate' : 1e-3,
                          })

          solver.train()

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

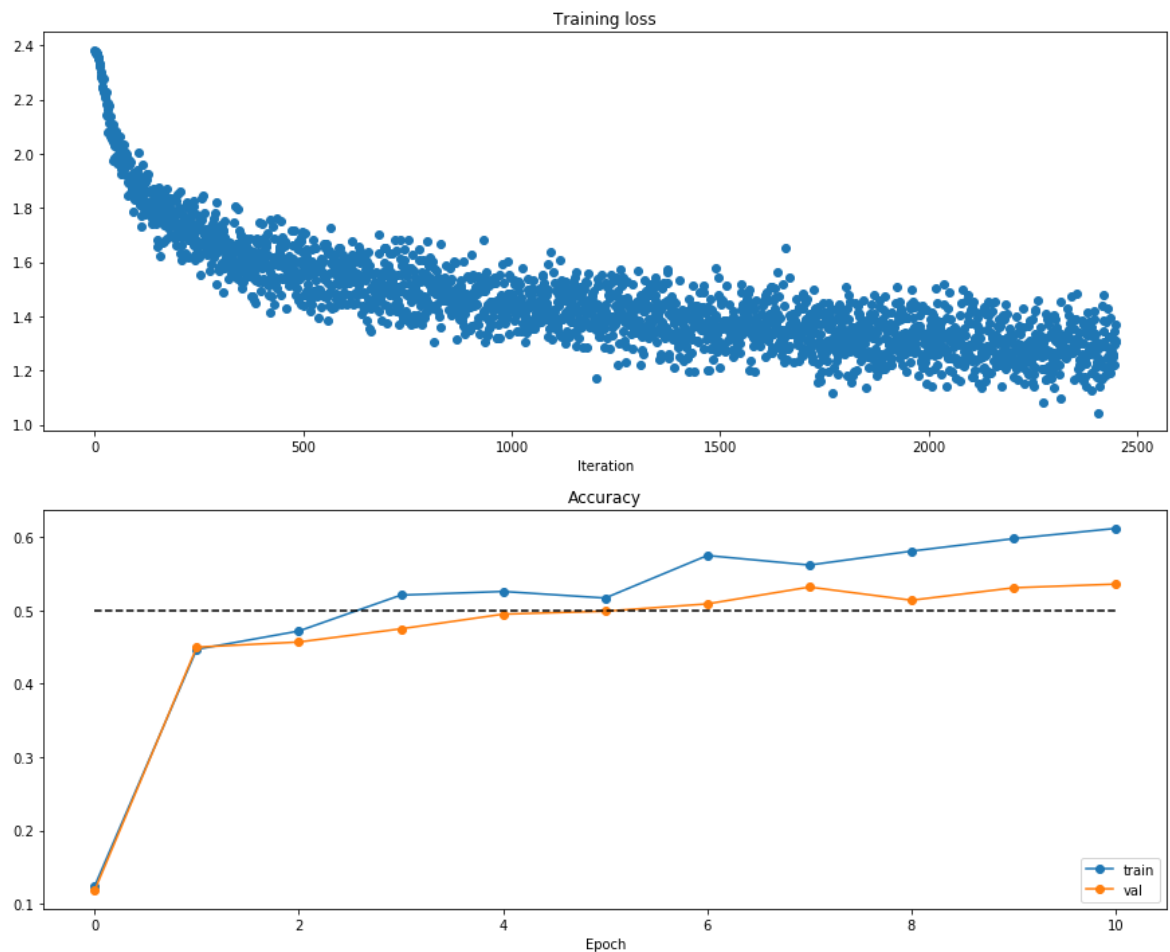
          (Iteration 1 / 2450) loss: 2.381478
          (Epoch 0 / 10) train acc: 0.124000; val_acc: 0.118000
          (Epoch 1 / 10) train acc: 0.447000; val_acc: 0.450000
          (Epoch 2 / 10) train acc: 0.472000; val_acc: 0.457000
          (Iteration 491 / 2450) loss: 1.457798
          (Epoch 3 / 10) train acc: 0.521000; val_acc: 0.475000
          (Epoch 4 / 10) train acc: 0.526000; val_acc: 0.495000
          (Iteration 981 / 2450) loss: 1.523241
          (Epoch 5 / 10) train acc: 0.517000; val_acc: 0.499000
          (Epoch 6 / 10) train acc: 0.575000; val_acc: 0.509000
          (Iteration 1471 / 2450) loss: 1.432054
          (Epoch 7 / 10) train acc: 0.562000; val_acc: 0.532000
          (Epoch 8 / 10) train acc: 0.581000; val_acc: 0.514000
          (Iteration 1961 / 2450) loss: 1.343220
          (Epoch 9 / 10) train acc: 0.598000; val_acc: 0.531000
          (Epoch 10 / 10) train acc: 0.612000; val_acc: 0.536000
```



In [118]: *# Run this cell to visualize training loss and train / val accuracy*

```
plt.subplot(2, 1, 1)
plt.title('Training loss')
plt.plot(solver.loss_history, 'o')
plt.xlabel('Iteration')

plt.subplot(2, 1, 2)
plt.title('Accuracy')
plt.plot(solver.train_acc_history, '-o', label='train')
plt.plot(solver.val_acc_history, '-o', label='val')
plt.plot([0.5] * len(solver.val_acc_history), 'k--')
plt.xlabel('Epoch')
plt.legend(loc='lower right')
plt.gcf().set_size_inches(15, 12)
plt.show()
```



## Multilayer Neural Network

Now, we implement a multi-layer neural network.

Read through the `FullyConnectedNet` class in the file `nndl/fc_net.py`.

Implement the initialization, the forward pass, and the backward pass. There will be lines for batchnorm and dropout layers and caches; ignore these all for now. That'll be in assignment #4.

```
In [148]: N, D, H1, H2, C = 2, 15, 20, 30, 10
X = np.random.randn(N, D)
y = np.random.randint(C, size=(N,))

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

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

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

```
Running check with reg = 0
Initial loss: 2.3013653843
W1 relative error: 1.35064716142e-06
W2 relative error: 4.01073941783e-07
W3 relative error: 1.04377988289e-07
b1 relative error: 8.39454578961e-08
b2 relative error: 7.1797379996e-09
b3 relative error: 8.52298577931e-11
Running check with reg = 3.14
Initial loss: 7.31589155688
W1 relative error: 6.28609875144e-09
W2 relative error: 1.44438461983e-08
W3 relative error: 1.28177508603e-08
b1 relative error: 2.11059960475e-07
b2 relative error: 5.58921132425e-09
b3 relative error: 1.99575793561e-10
```

```
In [163]: # Use the three layer neural network to overfit a small dataset.

num_train = 50
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'],
}

#### !!!!!
# Play around with the weight_scale and learning_rate so that you can
overfit a small dataset.
# Your training accuracy should be 1.0 to receive full credit on this
part.
weight_scale = 1e-2
learning_rate = 1e-3

model = FullyConnectedNet([100, 100],
                           weight_scale=weight_scale, dtype=np.float64)
solver = Solver(model, small_data,
                 print_every=400, num_epochs=200, batch_size=25,
                 update_rule='sgd',
                 optim_config={
                     'learning_rate': learning_rate,
                 })
solver.train()

plt.plot(solver.loss_history, 'o')
plt.title('Training loss history')
plt.xlabel('Iteration')
plt.ylabel('Training loss')
plt.show()
```

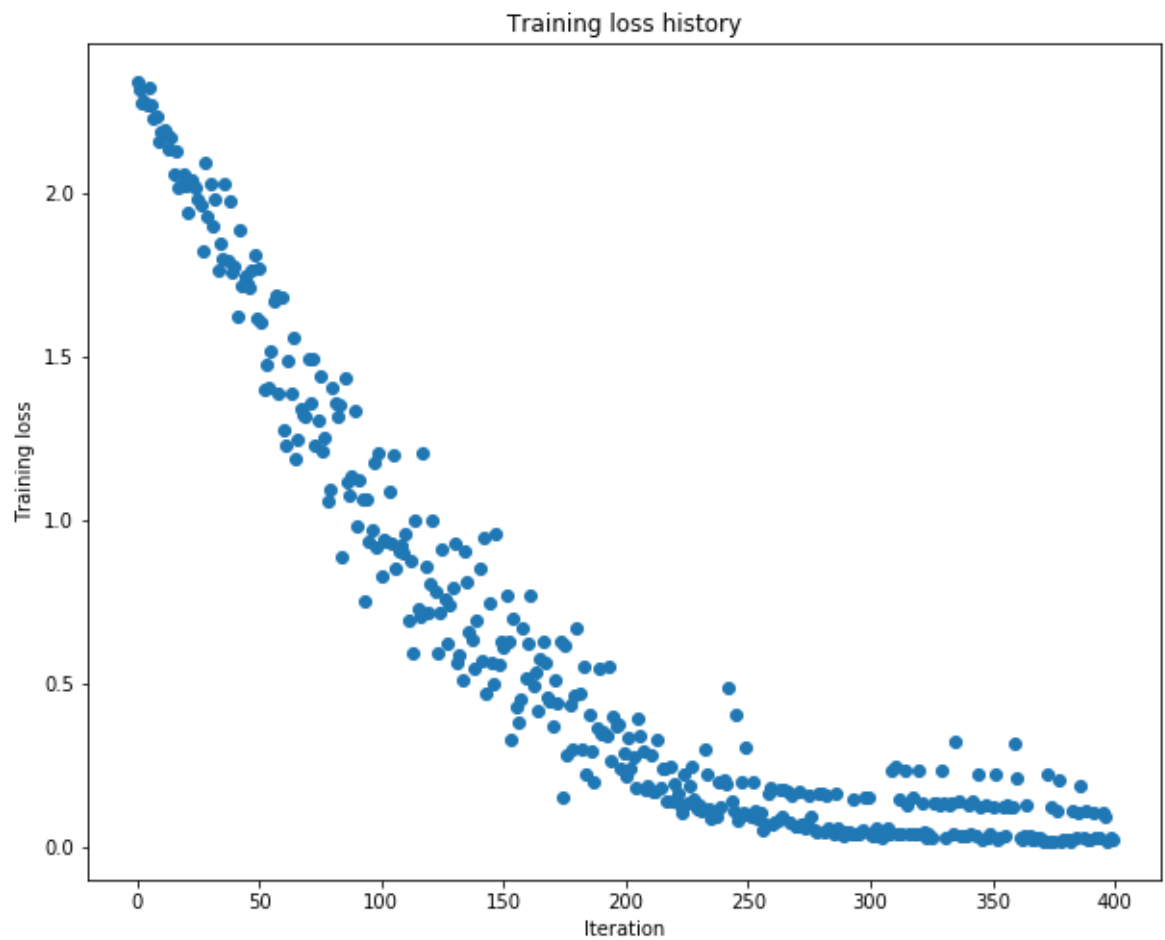
```
(Iteration 1 / 400) loss: 2.340599
(Epoch 0 / 200) train acc: 0.100000; val_acc: 0.102000
(Epoch 1 / 200) train acc: 0.120000; val_acc: 0.100000
(Epoch 2 / 200) train acc: 0.140000; val_acc: 0.107000
(Epoch 3 / 200) train acc: 0.180000; val_acc: 0.116000
(Epoch 4 / 200) train acc: 0.160000; val_acc: 0.115000
(Epoch 5 / 200) train acc: 0.200000; val_acc: 0.119000
(Epoch 6 / 200) train acc: 0.240000; val_acc: 0.125000
(Epoch 7 / 200) train acc: 0.280000; val_acc: 0.128000
(Epoch 8 / 200) train acc: 0.380000; val_acc: 0.128000
(Epoch 9 / 200) train acc: 0.360000; val_acc: 0.127000
(Epoch 10 / 200) train acc: 0.480000; val_acc: 0.131000
(Epoch 11 / 200) train acc: 0.440000; val_acc: 0.126000
(Epoch 12 / 200) train acc: 0.420000; val_acc: 0.130000
(Epoch 13 / 200) train acc: 0.460000; val_acc: 0.136000
(Epoch 14 / 200) train acc: 0.500000; val_acc: 0.147000
(Epoch 15 / 200) train acc: 0.520000; val_acc: 0.159000
(Epoch 16 / 200) train acc: 0.520000; val_acc: 0.158000
(Epoch 17 / 200) train acc: 0.540000; val_acc: 0.165000
(Epoch 18 / 200) train acc: 0.560000; val_acc: 0.157000
(Epoch 19 / 200) train acc: 0.540000; val_acc: 0.161000
(Epoch 20 / 200) train acc: 0.540000; val_acc: 0.162000
(Epoch 21 / 200) train acc: 0.560000; val_acc: 0.164000
(Epoch 22 / 200) train acc: 0.560000; val_acc: 0.161000
(Epoch 23 / 200) train acc: 0.560000; val_acc: 0.166000
(Epoch 24 / 200) train acc: 0.580000; val_acc: 0.163000
(Epoch 25 / 200) train acc: 0.600000; val_acc: 0.177000
(Epoch 26 / 200) train acc: 0.600000; val_acc: 0.175000
(Epoch 27 / 200) train acc: 0.600000; val_acc: 0.173000
(Epoch 28 / 200) train acc: 0.580000; val_acc: 0.169000
(Epoch 29 / 200) train acc: 0.580000; val_acc: 0.169000
(Epoch 30 / 200) train acc: 0.580000; val_acc: 0.165000
(Epoch 31 / 200) train acc: 0.620000; val_acc: 0.176000
(Epoch 32 / 200) train acc: 0.600000; val_acc: 0.183000
(Epoch 33 / 200) train acc: 0.620000; val_acc: 0.184000
(Epoch 34 / 200) train acc: 0.620000; val_acc: 0.176000
(Epoch 35 / 200) train acc: 0.620000; val_acc: 0.178000
(Epoch 36 / 200) train acc: 0.620000; val_acc: 0.175000
(Epoch 37 / 200) train acc: 0.620000; val_acc: 0.174000
(Epoch 38 / 200) train acc: 0.700000; val_acc: 0.188000
(Epoch 39 / 200) train acc: 0.680000; val_acc: 0.172000
(Epoch 40 / 200) train acc: 0.700000; val_acc: 0.176000
(Epoch 41 / 200) train acc: 0.700000; val_acc: 0.175000
(Epoch 42 / 200) train acc: 0.700000; val_acc: 0.173000
(Epoch 43 / 200) train acc: 0.700000; val_acc: 0.167000
(Epoch 44 / 200) train acc: 0.740000; val_acc: 0.175000
(Epoch 45 / 200) train acc: 0.740000; val_acc: 0.176000
(Epoch 46 / 200) train acc: 0.700000; val_acc: 0.178000
(Epoch 47 / 200) train acc: 0.760000; val_acc: 0.183000
(Epoch 48 / 200) train acc: 0.760000; val_acc: 0.178000
(Epoch 49 / 200) train acc: 0.780000; val_acc: 0.192000
(Epoch 50 / 200) train acc: 0.800000; val_acc: 0.193000
(Epoch 51 / 200) train acc: 0.800000; val_acc: 0.190000
(Epoch 52 / 200) train acc: 0.820000; val_acc: 0.193000
(Epoch 53 / 200) train acc: 0.820000; val_acc: 0.190000
(Epoch 54 / 200) train acc: 0.800000; val_acc: 0.189000
(Epoch 55 / 200) train acc: 0.820000; val_acc: 0.197000
```

```
(Epoch 56 / 200) train acc: 0.840000; val_acc: 0.191000
(Epoch 57 / 200) train acc: 0.840000; val_acc: 0.196000
(Epoch 58 / 200) train acc: 0.840000; val_acc: 0.188000
(Epoch 59 / 200) train acc: 0.860000; val_acc: 0.187000
(Epoch 60 / 200) train acc: 0.840000; val_acc: 0.184000
(Epoch 61 / 200) train acc: 0.820000; val_acc: 0.183000
(Epoch 62 / 200) train acc: 0.860000; val_acc: 0.189000
(Epoch 63 / 200) train acc: 0.860000; val_acc: 0.189000
(Epoch 64 / 200) train acc: 0.860000; val_acc: 0.195000
(Epoch 65 / 200) train acc: 0.840000; val_acc: 0.199000
(Epoch 66 / 200) train acc: 0.840000; val_acc: 0.195000
(Epoch 67 / 200) train acc: 0.860000; val_acc: 0.191000
(Epoch 68 / 200) train acc: 0.880000; val_acc: 0.190000
(Epoch 69 / 200) train acc: 0.880000; val_acc: 0.188000
(Epoch 70 / 200) train acc: 0.880000; val_acc: 0.192000
(Epoch 71 / 200) train acc: 0.880000; val_acc: 0.182000
(Epoch 72 / 200) train acc: 0.880000; val_acc: 0.179000
(Epoch 73 / 200) train acc: 0.880000; val_acc: 0.191000
(Epoch 74 / 200) train acc: 0.880000; val_acc: 0.196000
(Epoch 75 / 200) train acc: 0.880000; val_acc: 0.186000
(Epoch 76 / 200) train acc: 0.880000; val_acc: 0.193000
(Epoch 77 / 200) train acc: 0.880000; val_acc: 0.189000
(Epoch 78 / 200) train acc: 0.880000; val_acc: 0.186000
(Epoch 79 / 200) train acc: 0.880000; val_acc: 0.182000
(Epoch 80 / 200) train acc: 0.880000; val_acc: 0.188000
(Epoch 81 / 200) train acc: 0.880000; val_acc: 0.183000
(Epoch 82 / 200) train acc: 0.880000; val_acc: 0.187000
(Epoch 83 / 200) train acc: 0.880000; val_acc: 0.189000
(Epoch 84 / 200) train acc: 0.900000; val_acc: 0.183000
(Epoch 85 / 200) train acc: 0.900000; val_acc: 0.186000
(Epoch 86 / 200) train acc: 0.900000; val_acc: 0.184000
(Epoch 87 / 200) train acc: 0.900000; val_acc: 0.186000
(Epoch 88 / 200) train acc: 0.940000; val_acc: 0.187000
(Epoch 89 / 200) train acc: 0.920000; val_acc: 0.190000
(Epoch 90 / 200) train acc: 0.940000; val_acc: 0.177000
(Epoch 91 / 200) train acc: 0.960000; val_acc: 0.184000
(Epoch 92 / 200) train acc: 0.960000; val_acc: 0.182000
(Epoch 93 / 200) train acc: 0.960000; val_acc: 0.181000
(Epoch 94 / 200) train acc: 0.980000; val_acc: 0.183000
(Epoch 95 / 200) train acc: 0.980000; val_acc: 0.187000
(Epoch 96 / 200) train acc: 0.980000; val_acc: 0.188000
(Epoch 97 / 200) train acc: 0.980000; val_acc: 0.189000
(Epoch 98 / 200) train acc: 0.980000; val_acc: 0.188000
(Epoch 99 / 200) train acc: 0.980000; val_acc: 0.185000
(Epoch 100 / 200) train acc: 0.980000; val_acc: 0.188000
(Epoch 101 / 200) train acc: 0.980000; val_acc: 0.187000
(Epoch 102 / 200) train acc: 0.980000; val_acc: 0.187000
(Epoch 103 / 200) train acc: 0.980000; val_acc: 0.187000
(Epoch 104 / 200) train acc: 0.980000; val_acc: 0.192000
(Epoch 105 / 200) train acc: 0.980000; val_acc: 0.190000
(Epoch 106 / 200) train acc: 0.980000; val_acc: 0.190000
(Epoch 107 / 200) train acc: 0.980000; val_acc: 0.189000
(Epoch 108 / 200) train acc: 0.980000; val_acc: 0.188000
(Epoch 109 / 200) train acc: 0.980000; val_acc: 0.183000
(Epoch 110 / 200) train acc: 0.980000; val_acc: 0.188000
(Epoch 111 / 200) train acc: 0.980000; val_acc: 0.187000
(Epoch 112 / 200) train acc: 0.980000; val_acc: 0.190000
```

```
(Epoch 113 / 200) train acc: 0.980000; val_acc: 0.186000
(Epoch 114 / 200) train acc: 0.980000; val_acc: 0.190000
(Epoch 115 / 200) train acc: 0.980000; val_acc: 0.200000
(Epoch 116 / 200) train acc: 0.980000; val_acc: 0.194000
(Epoch 117 / 200) train acc: 0.980000; val_acc: 0.196000
(Epoch 118 / 200) train acc: 0.980000; val_acc: 0.190000
(Epoch 119 / 200) train acc: 0.980000; val_acc: 0.195000
(Epoch 120 / 200) train acc: 0.980000; val_acc: 0.191000
(Epoch 121 / 200) train acc: 0.980000; val_acc: 0.198000
(Epoch 122 / 200) train acc: 0.980000; val_acc: 0.191000
(Epoch 123 / 200) train acc: 0.980000; val_acc: 0.193000
(Epoch 124 / 200) train acc: 0.980000; val_acc: 0.193000
(Epoch 125 / 200) train acc: 0.980000; val_acc: 0.197000
(Epoch 126 / 200) train acc: 0.980000; val_acc: 0.192000
(Epoch 127 / 200) train acc: 0.980000; val_acc: 0.188000
(Epoch 128 / 200) train acc: 0.980000; val_acc: 0.191000
(Epoch 129 / 200) train acc: 0.980000; val_acc: 0.190000
(Epoch 130 / 200) train acc: 0.980000; val_acc: 0.190000
(Epoch 131 / 200) train acc: 0.980000; val_acc: 0.191000
(Epoch 132 / 200) train acc: 0.980000; val_acc: 0.194000
(Epoch 133 / 200) train acc: 0.980000; val_acc: 0.191000
(Epoch 134 / 200) train acc: 0.980000; val_acc: 0.192000
(Epoch 135 / 200) train acc: 0.980000; val_acc: 0.191000
(Epoch 136 / 200) train acc: 0.980000; val_acc: 0.191000
(Epoch 137 / 200) train acc: 0.980000; val_acc: 0.191000
(Epoch 138 / 200) train acc: 0.980000; val_acc: 0.194000
(Epoch 139 / 200) train acc: 0.980000; val_acc: 0.193000
(Epoch 140 / 200) train acc: 0.980000; val_acc: 0.194000
(Epoch 141 / 200) train acc: 0.980000; val_acc: 0.194000
(Epoch 142 / 200) train acc: 0.980000; val_acc: 0.193000
(Epoch 143 / 200) train acc: 0.980000; val_acc: 0.195000
(Epoch 144 / 200) train acc: 0.980000; val_acc: 0.194000
(Epoch 145 / 200) train acc: 0.980000; val_acc: 0.192000
(Epoch 146 / 200) train acc: 0.980000; val_acc: 0.192000
(Epoch 147 / 200) train acc: 0.980000; val_acc: 0.192000
(Epoch 148 / 200) train acc: 0.980000; val_acc: 0.195000
(Epoch 149 / 200) train acc: 0.980000; val_acc: 0.193000
(Epoch 150 / 200) train acc: 0.980000; val_acc: 0.196000
(Epoch 151 / 200) train acc: 0.980000; val_acc: 0.196000
(Epoch 152 / 200) train acc: 0.980000; val_acc: 0.195000
(Epoch 153 / 200) train acc: 0.980000; val_acc: 0.197000
(Epoch 154 / 200) train acc: 0.980000; val_acc: 0.197000
(Epoch 155 / 200) train acc: 0.980000; val_acc: 0.195000
(Epoch 156 / 200) train acc: 0.980000; val_acc: 0.196000
(Epoch 157 / 200) train acc: 0.980000; val_acc: 0.197000
(Epoch 158 / 200) train acc: 0.980000; val_acc: 0.198000
(Epoch 159 / 200) train acc: 0.980000; val_acc: 0.199000
(Epoch 160 / 200) train acc: 0.980000; val_acc: 0.198000
(Epoch 161 / 200) train acc: 0.980000; val_acc: 0.198000
(Epoch 162 / 200) train acc: 0.980000; val_acc: 0.197000
(Epoch 163 / 200) train acc: 0.980000; val_acc: 0.197000
(Epoch 164 / 200) train acc: 0.980000; val_acc: 0.198000
(Epoch 165 / 200) train acc: 0.980000; val_acc: 0.197000
(Epoch 166 / 200) train acc: 0.980000; val_acc: 0.198000
(Epoch 167 / 200) train acc: 0.980000; val_acc: 0.196000
(Epoch 168 / 200) train acc: 0.980000; val_acc: 0.195000
(Epoch 169 / 200) train acc: 0.980000; val_acc: 0.196000
```

```
(Epoch 170 / 200) train acc: 0.980000; val_acc: 0.194000
(Epoch 171 / 200) train acc: 0.980000; val_acc: 0.195000
(Epoch 172 / 200) train acc: 0.980000; val_acc: 0.195000
(Epoch 173 / 200) train acc: 0.980000; val_acc: 0.195000
(Epoch 174 / 200) train acc: 0.980000; val_acc: 0.194000
(Epoch 175 / 200) train acc: 0.980000; val_acc: 0.193000
(Epoch 176 / 200) train acc: 0.980000; val_acc: 0.196000
(Epoch 177 / 200) train acc: 0.980000; val_acc: 0.195000
(Epoch 178 / 200) train acc: 0.980000; val_acc: 0.194000
(Epoch 179 / 200) train acc: 0.980000; val_acc: 0.194000
(Epoch 180 / 200) train acc: 0.980000; val_acc: 0.194000
(Epoch 181 / 200) train acc: 0.980000; val_acc: 0.194000
(Epoch 182 / 200) train acc: 0.980000; val_acc: 0.193000
(Epoch 183 / 200) train acc: 0.980000; val_acc: 0.198000
(Epoch 184 / 200) train acc: 0.980000; val_acc: 0.196000
(Epoch 185 / 200) train acc: 0.980000; val_acc: 0.196000
(Epoch 186 / 200) train acc: 0.980000; val_acc: 0.197000
(Epoch 187 / 200) train acc: 0.980000; val_acc: 0.196000
(Epoch 188 / 200) train acc: 0.980000; val_acc: 0.196000
(Epoch 189 / 200) train acc: 0.980000; val_acc: 0.195000
(Epoch 190 / 200) train acc: 0.980000; val_acc: 0.198000
(Epoch 191 / 200) train acc: 0.980000; val_acc: 0.198000
(Epoch 192 / 200) train acc: 0.980000; val_acc: 0.195000
(Epoch 193 / 200) train acc: 0.980000; val_acc: 0.196000
(Epoch 194 / 200) train acc: 1.000000; val_acc: 0.192000
(Epoch 195 / 200) train acc: 1.000000; val_acc: 0.192000
(Epoch 196 / 200) train acc: 1.000000; val_acc: 0.192000
(Epoch 197 / 200) train acc: 1.000000; val_acc: 0.194000
(Epoch 198 / 200) train acc: 1.000000; val_acc: 0.198000
(Epoch 199 / 200) train acc: 1.000000; val_acc: 0.198000
(Epoch 200 / 200) train acc: 1.000000; val_acc: 0.196000
```





```
In [167]: print solver.train_acc_history[-1]
```

```
1.0
```

```

1 import numpy as np
2 import pdb
3
4 """
5 This code was originally written for CS 231n at Stanford University
6 (cs231n.stanford.edu). It has been modified in various areas for use in the
7 ECE 239AS class at UCLA. This includes the descriptions of what code to
8 implement as well as some slight potential changes in variable names to be
9 consistent with class nomenclature. We thank Justin Johnson & Serena Yeung for
10 permission to use this code. To see the original version, please visit
11 cs231n.stanford.edu.
12 """
13
14
15 def affine_forward(x, w, b):
16     """
17     Computes the forward pass for an affine (fully-connected) layer.
18
19     The input x has shape (N, d_1, ..., d_k) and contains a minibatch of N
20     examples, where each example x[i] has shape (d_1, ..., d_k). We will
21     reshape each input into a vector of dimension D = d_1 * ... * d_k, and
22     then transform it to an output vector of dimension M.
23
24     Inputs:
25     - x: A numpy array containing input data, of shape (N, d_1, ..., d_k)
26     - w: A numpy array of weights, of shape (D, M)
27     - b: A numpy array of biases, of shape (M,)
28
29     Returns a tuple of:
30     - out: output, of shape (N, M)
31     - cache: (x, w, b)
32     """
33
34     # ===== #
35     # YOUR CODE HERE:
36     # Calculate the output of the forward pass. Notice the dimensions
37     # of w are D x M, which is the transpose of what we did in earlier
38     # assignments.
39     # ===== #
40     shape = x.shape
41     N = shape[0]
42     D = np.prod(shape[1:])
43     reshaped_x = np.reshape(x, (N,D))
44
45     out = reshaped_x.dot(w) + b[:, np.newaxis].T
46
47     # ===== #
48     # END YOUR CODE HERE
49     # ===== #
50
51     cache = (x, w, b)
52     return out, cache
53
54
55 def affine_backward(dout, cache):
56     """
57     Computes the backward pass for an affine layer.
58
59     Inputs:
60     - dout: Upstream derivative, of shape (N, M)
61     - cache: Tuple of:
62       - x: Input data, of shape (N, d_1, ... d_k)
63       - w: Weights, of shape (D, M)
64
65     Returns a tuple of:
66     - dx: Gradient with respect to x, of shape (N, d_1, ..., d_k)
67     - dw: Gradient with respect to w, of shape (D, M)
68     - db: Gradient with respect to b, of shape (M,)
69     """

```

```

70 x, w, b = cache
71 dx, dw, db = None, None, None
72
73 # ===== #
74 # YOUR CODE HERE:
75 # Calculate the gradients for the backward pass.
76 # ===== #
77 N, M = dout.shape
78 D = w.shape[0]
79
80 reshaped_x = np.reshape(x, (N,D))
81
82 db = np.sum(dout, axis=0)
83 dw = reshaped_x.T.dot(dout)
84 dx = np.reshape(dout.dot(w.T), x.shape)
85
86
87 # ===== #
88 # END YOUR CODE HERE
89 # ===== #
90
91 return dx, dw, db
92
93 def relu_forward(x):
94     """
95     Computes the forward pass for a layer of rectified linear units (ReLU).
96
97     Input:
98     - x: Inputs, of any shape
99
100    Returns a tuple of:
101    - out: Output, of the same shape as x
102    - cache: x
103    """
104    # ===== #
105    # YOUR CODE HERE:
106    # Implement the ReLU forward pass.
107    # ===== #
108    out = np.empty_like(x)
109    out[:] = x
110    out[out<0] = 0
111    # ===== #
112    # END YOUR CODE HERE
113    # ===== #
114
115    cache = x
116    return out, cache
117
118
119 def relu_backward(dout, cache):
120     """
121     Computes the backward pass for a layer of rectified linear units (ReLU).
122
123     Input:
124     - dout: Upstream derivatives, of any shape
125     - cache: Input x, of same shape as dout
126
127     Returns:
128     - dx: Gradient with respect to x
129     """
130    x = cache
131
132    # ===== #
133    # YOUR CODE HERE:
134    # Implement the ReLU backward pass
135    # ===== #
136    dx = np.empty_like(dout)
137    dx[:] = dout
138    dx[x<0] = 0

```

```

139
140 # ===== #
141 # END YOUR CODE HERE
142 # ===== #
143
144 return dx
145
146 def svm_loss(x, y):
147     """
148     Computes the loss and gradient using for multiclass SVM classification.
149
150     Inputs:
151     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
152       for the ith input.
153     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
154       0 <= y[i] < C
155
156     Returns a tuple of:
157     - loss: Scalar giving the loss
158     - dx: Gradient of the loss with respect to x
159     """
160     N = x.shape[0]
161     correct_class_scores = x[np.arange(N), y]
162     margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
163     margins[np.arange(N), y] = 0
164     loss = np.sum(margins) / N
165     num_pos = np.sum(margins > 0, axis=1)
166     dx = np.zeros_like(x)
167     dx[margins > 0] = 1
168     dx[np.arange(N), y] -= num_pos
169     dx /= N
170     return loss, dx
171
172
173 def softmax_loss(x, y):
174     """
175     Computes the loss and gradient for softmax classification.
176
177     Inputs:
178     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
179       for the ith input.
180     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
181       0 <= y[i] < C
182
183     Returns a tuple of:
184     - loss: Scalar giving the loss
185     - dx: Gradient of the loss with respect to x
186     """
187
188     probs = np.exp(x - np.max(x, axis=1, keepdims=True))
189     probs /= np.sum(probs, axis=1, keepdims=True)
190     N = x.shape[0]
191     loss = -np.sum(np.log(probs[np.arange(N), y])) / N
192     dx = probs.copy()
193     dx[np.arange(N), y] -= 1
194     dx /= N
195     return loss, dx

```

```

1 import numpy as np
2
3 from .layers import *
4 from .layer_utils import *
5
6 """
7 This code was originally written for CS 231n at Stanford University
8 (cs231n.stanford.edu). It has been modified in various areas for use in the
9 ECE 239AS class at UCLA. This includes the descriptions of what code to
10 implement as well as some slight potential changes in variable names to be
11 consistent with class nomenclature. We thank Justin Johnson & Serena Yeung for
12 permission to use this code. To see the original version, please visit
13 cs231n.stanford.edu.
14 """
15
16 class TwoLayerNet(object):
17     """
18     A two-layer fully-connected neural network with ReLU nonlinearity and
19     softmax loss that uses a modular layer design. We assume an input dimension
20     of D, a hidden dimension of H, and perform classification over C classes.
21
22     The architecture should be affine - relu - affine - softmax.
23
24     Note that this class does not implement gradient descent; instead, it
25     will interact with a separate Solver object that is responsible for running
26     optimization.
27
28     The learnable parameters of the model are stored in the dictionary
29     self.params that maps parameter names to numpy arrays.
30     """
31
32     def __init__(self, input_dim=3*32*32, hidden_dims=100, num_classes=10,
33                 dropout=0, weight_scale=1e-3, reg=0.0):
34         """
35         Initialize a new network.
36
37         Inputs:
38         - input_dim: An integer giving the size of the input
39         - hidden_dims: An integer giving the size of the hidden layer
40         - num_classes: An integer giving the number of classes to classify
41         - dropout: Scalar between 0 and 1 giving dropout strength.
42         - weight_scale: Scalar giving the standard deviation for random
43           initialization of the weights.
44         - reg: Scalar giving L2 regularization strength.
45         """
46         self.params = {}
47         self.reg = reg
48
49         # ===== #
50         # YOUR CODE HERE:
51         # Initialize W1, W2, b1, and b2. Store these as self.params['W1'],
52         # self.params['W2'], self.params['b1'] and self.params['b2']. The
53         # biases are initialized to zero and the weights are initialized
54         # so that each parameter has mean 0 and standard deviation weight_scale.
55         # The dimensions of W1 should be (input_dim, hidden_dim) and the
56         # dimensions of W2 should be (hidden_dims, num_classes)
57         # ===== #
58
59         self.params['W1'] = weight_scale * np.random.randn(input_dim, hidden_dims)
60         self.params['b1'] = np.zeros(hidden_dims)
61         self.params['W2'] = weight_scale * np.random.randn(hidden_dims, num_classes)
62         self.params['b2'] = np.zeros(num_classes)
63
64         # ===== #
65         # END YOUR CODE HERE
66         # ===== #
67
68     def loss(self, X, y=None):
69         """
70         Compute loss and gradient for a minibatch of data.
71
72         Inputs:
73         - X: Array of input data of shape (N, d_1, ..., d_k)
74         - y: Array of labels, of shape (N,). y[i] gives the label for X[i].
75
76         Returns:
77         If y is None, then run a test-time forward pass of the model and return:
78         - scores: Array of shape (N, C) giving classification scores, where
79           scores[i, c] is the classification score for X[i] and class c.
80
81         If y is not None, then run a training-time forward and backward pass and
82         return a tuple of:
83         - loss: Scalar value giving the loss
84         - grads: Dictionary with the same keys as self.params, mapping parameter
85           names to gradients of the loss with respect to those parameters.
86         """
87         scores = None
88
89         # ===== #
90         # YOUR CODE HERE:
91         # Implement the forward pass of the two-layer neural network. Store
92         # the class scores as the variable 'scores'. Be sure to use the layers

```

```

93     # you prior implemented.
94     # ===== #
95
96     h1, cache1 = affine_relu_forward(X, self.params['W1'], self.params['b1'])
97     scores, cache2 = affine_relu_forward(h1, self.params['W2'], self.params['b2'])
98     # ===== #
99     # END YOUR CODE HERE
100    # ===== #
101
102    # If y is None then we are in test mode so just return scores
103    if y is None:
104        return scores
105
106    loss, grads = 0, {}
107    # ===== #
108    # YOUR CODE HERE:
109    # Implement the backward pass of the two-layer neural net. Store
110    # the loss as the variable 'loss' and store the gradients in the
111    # 'grads' dictionary. For the grads dictionary, grads['W1'] holds
112    # the gradient for W1, grads['b1'] holds the gradient for b1, etc.
113    # i.e., grads[k] holds the gradient for self.params[k].
114    #
115    # Add L2 regularization, where there is an added cost 0.5*self.reg*W^2
116    # for each W. Be sure to include the 0.5 multiplying factor to
117    # match our implementation.
118    #
119    # And be sure to use the layers you prior implemented.
120    # ===== #
121    W1 = self.params['W1']
122    W2 = self.params['W2']
123
124    num_examples = scores.shape[0]
125
126    max_score = np.amax(scores, axis=1)
127    scores -= max_score[:, np.newaxis]
128
129    e_scores = np.exp(scores)
130    sums = np.sum(e_scores, axis=1)
131    log_sums = np.log(sums)
132    y_terms = scores[np.arange(num_examples), y]
133    loss = np.sum(log_sums - y_terms)/num_examples + .5*self.reg*np.sum(W1*W1) + .5*self.reg*np.sum(W2*W2)
134
135    d_scores = e_scores/sums[:, np.newaxis]
136    d_scores[np.arange(num_examples), y] -= 1
137    d_scores = d_scores.T/num_examples
138
139    dx2, dw2, db2 = affine_relu_backward(d_scores.T, cache2)
140    dx1, dw1, db1 = affine_relu_backward(dx2, cache1)
141
142
143    grads['W1'] = dw1 + self.reg*W1
144    grads['b1'] = db1
145    grads['W2'] = dw2 + self.reg*W2
146    grads['b2'] = db2
147
148
149
150
151
152
153    # ===== #
154    # END YOUR CODE HERE
155    # ===== #
156
157    return loss, grads
158
159
160 class FullyConnectedNet(object):
161     """
162     A fully-connected neural network with an arbitrary number of hidden layers,
163     ReLU nonlinearities, and a softmax loss function. This will also implement
164     dropout and batch normalization as options. For a network with L layers,
165     the architecture will be
166
167     {affine - [batch norm] - relu - [dropout]} x (L - 1) - affine - softmax
168
169     where batch normalization and dropout are optional, and the {...} block is
170     repeated L - 1 times.
171
172     Similar to the TwoLayerNet above, learnable parameters are stored in the
173     self.params dictionary and will be learned using the Solver class.
174     """
175
176     def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,
177                  dropout=0, use_batchnorm=False, reg=0.0,
178                  weight_scale=1e-2, dtype=np.float32, seed=None):
179         """
180         Initialize a new FullyConnectedNet.
181
182         Inputs:
183         - hidden_dims: A list of integers giving the size of each hidden layer.
184         - input_dim: An integer giving the size of the input.

```

```

185 - num_classes: An integer giving the number of classes to classify.
186 - dropout: Scalar between 0 and 1 giving dropout strength. If dropout=0 then
187   the network should not use dropout at all.
188 - use_batchnorm: Whether or not the network should use batch normalization.
189 - reg: Scalar giving L2 regularization strength.
190 - weight_scale: Scalar giving the standard deviation for random
191   initialization of the weights.
192 - dtype: A numpy datatype object; all computations will be performed using
193   this datatype. float32 is faster but less accurate, so you should use
194   float64 for numeric gradient checking.
195 - seed: If not None, then pass this random seed to the dropout layers. This
196   will make the dropout layers deterministic so we can gradient check the
197   model.
198 """
199 self.use_batchnorm = use_batchnorm
200 self.use_dropout = dropout > 0
201 self.reg = reg
202 self.num_layers = 1 + len(hidden_dims)
203 self.dtype = dtype
204 self.params = {}
205
206 # ===== #
207 # YOUR CODE HERE:
208 # Initialize all parameters of the network in the self.params dictionary.
209 # The weights and biases of layer 1 are W1 and b1; and in general the
210 # weights and biases of layer i are Wi and bi. The
211 # biases are initialized to zero and the weights are initialized
212 # so that each parameter has mean 0 and standard deviation weight_scale.
213 # ===== #
214
215 dimensions = [input_dim] + hidden_dims + [num_classes]
216
217 for i in np.arange(self.num_layers):
218     self.params['W{}'.format(i+1)] = weight_scale * np.random.randn(dimensions[i], dimensions[i+1])
219     self.params['b{}'.format(i+1)] = np.zeros(dimensions[i+1])
220
221
222 # ===== #
223 # END YOUR CODE HERE
224 # ===== #
225
226
227 # When using dropout we need to pass a dropout_param dictionary to each
228 # dropout layer so that the layer knows the dropout probability and the mode
229 # (train / test). You can pass the same dropout_param to each dropout layer.
230 self.dropout_param = {}
231 if self.use_dropout:
232     self.dropout_param = {'mode': 'train', 'p': dropout}
233     if seed is not None:
234         self.dropout_param['seed'] = seed
235
236 # With batch normalization we need to keep track of running means and
237 # variances, so we need to pass a special bn_param object to each batch
238 # normalization layer. You should pass self.bn_params[0] to the forward pass
239 # of the first batch normalization layer, self.bn_params[1] to the forward
240 # pass of the second batch normalization layer, etc.
241 self.bn_params = []
242 if self.use_batchnorm:
243     self.bn_params = [{'mode': 'train'} for i in np.arange(self.num_layers - 1)]
244
245 # Cast all parameters to the correct datatype
246 for k, v in self.params.items():
247     self.params[k] = v.astype(dtype)
248
249
250 def loss(self, X, y=None):
251     """
252     Compute loss and gradient for the fully-connected net.
253
254     Input / output: Same as TwoLayerNet above.
255     """
256     X = X.astype(self.dtype)
257     mode = 'test' if y is None else 'train'
258
259     # Set train/test mode for batchnorm params and dropout param since they
260     # behave differently during training and testing.
261     if self.dropout_param is not None:
262         self.dropout_param['mode'] = mode
263     if self.use_batchnorm:
264         for bn_param in self.bn_params:
265             bn_param[mode] = mode
266
267     scores = None
268
269     # ===== #
270     # YOUR CODE HERE:
271     # Implement the forward pass of the FC net and store the output
272     # scores as the variable "scores".
273     # ===== #
274
275     caches = []
276     layer_scores = []

```



```

277 layer_scores.append(X)
278
279 for i in np.arange(self.num_layers):
280     temp_score, temp_cache = affine_relu_forward(layer_scores[i], self.params['W{}'.format(i+1)], self.params['b{}'.format(i+1)])
281     caches.append(temp_cache)
282     layer_scores.append(temp_score)
283
284 scores = layer_scores[-1]
285
286 # ===== #
287 # END YOUR CODE HERE
288 # ===== #
289
290 # If test mode return early
291 if mode == 'test':
292     return scores
293
294 loss, grads = 0.0, {}
295 # ===== #
296 # YOUR CODE HERE:
297 # Implement the backwards pass of the FC net and store the gradients
298 # in the grads dict, so that grads[k] is the gradient of self.params[k]
299 # Be sure your L2 regularization includes a 0.5 factor.
300 # ===== #
301 num_examples = scores.shape[0]
302
303 max_score = np.amax(scores, axis=1)
304 scores -= max_score[:, np.newaxis]
305
306 e_scores = np.exp(scores)
307 sums = np.sum(e_scores, axis=1)
308 log_sums = np.log(sums)
309 y_terms = scores[np.arange(num_examples), y]
310
311 reg_loss = 0
312 for i in np.arange(self.num_layers):
313     W = self.params['W{}'.format(i+1)]
314     reg_loss += .5*self.reg*np.sum(W*W)
315
316 loss = np.sum(log_sums - y_terms)/num_examples + reg_loss
317
318
319
320 d_scores = e_scores/sums[:, np.newaxis]
321 d_scores[np.arange(num_examples), y] -= 1
322 d_scores = d_scores/num_examples
323
324 #print len(caches)
325 #print self.num_layers
326
327 for i in np.arange(self.num_layers-1, -1, -1):
328     dx, dw, db = affine_relu_backward(d_scores, caches[i])
329     grads['W{}'.format(i+1)] = dw + self.reg*self.params['W{}'.format(i+1)]
330     grads['b{}'.format(i+1)] = db
331     d_scores = dx
332
333 #dx2, dw2, db2 = affine_relu_backward(d_scores, cache[i])
334 # dx1, dw1, db1 = affine_relu_backward(dx2, cache1)
335
336
337 # ===== #
338 # END YOUR CODE HERE
339 # ===== #
340 return loss, grads

```