

①

1) a) In the autoencoder, we have an encoding layer that outputs the low dimensionality rep^r of x .

\Rightarrow if $x \in \mathbb{R}^n$, $Wx \in \mathbb{R}^m$ is the low dimensional rep^r of x such that $W \in \mathbb{R}^{m \times n}$ & $m < n$.

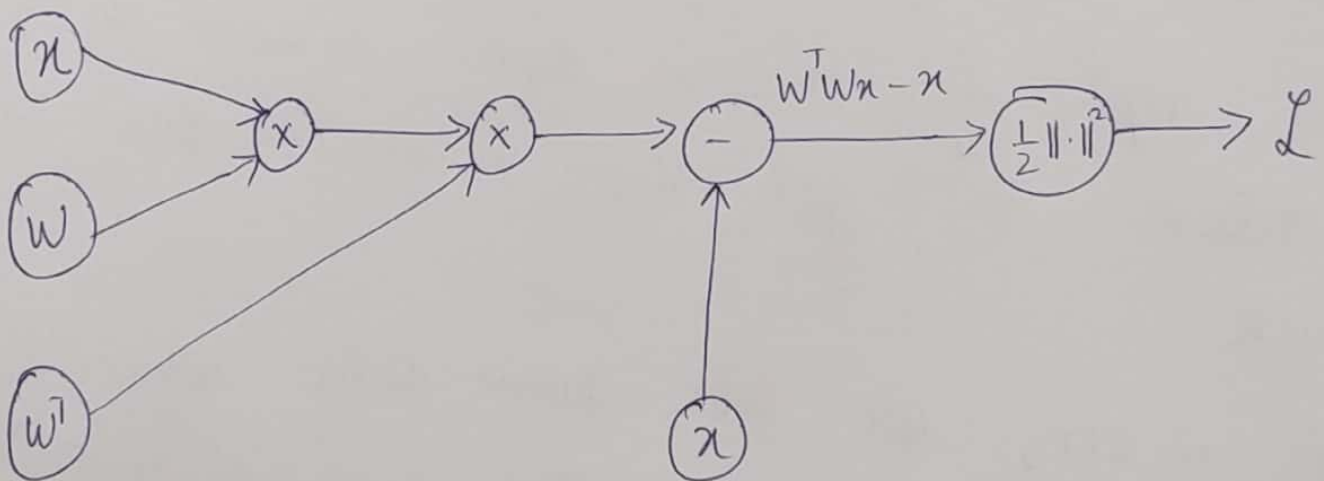
This rep^r of x is then passed to a decoder to get the reconstructed signal back.

As $m < n$, we lose some data when we perform Wx . \therefore The aim of encoding is to minimize this information loss.

To decode the rep^r back to original signal $\rightarrow \underline{W^T W x} \in \mathbb{R}^n$

\therefore by minimizing the Euclidean distance (L_2 -norm) b/n Wx & $W^T Wx$ we tune W , thus preserving most of the information.

$$(b) \quad d = \frac{1}{2} \| W^T Wx - x \|^2$$



Cp From the computational graph in b, we can see that there are two paths to WL . One corresponding to W . The other corresponding to W^T . The cp's can be added to get the final cost function.

$$d = \underset{d_1}{\text{cost of path } w} + \underset{d_2}{\text{cost of path } w^T} \quad (3)$$

$$\frac{\partial L}{\partial w} = \frac{\partial L_1}{\partial w} + \frac{\partial L_2}{\partial w}$$

$$d_1 = d_2 = d$$

$$\frac{\partial d_2}{\partial w} = \left(\frac{\partial d_2}{\partial w^T} \right)^T = \left(\frac{\partial L}{\partial w^T} \right)^T$$

\therefore total gradient =

$$\nabla_w L = \frac{\partial L}{\partial w} = \frac{\partial L}{\partial w} + \underline{\underline{\left(\frac{\partial L}{\partial w^T} \right)^T}}$$

Q7 Consider $z \in \mathbb{R}^n$

$$f(z) = \|z\|_2^2$$

$$= \left(\left(\sum_{k=1}^n z_k^2 \right)^{1/2} \right)^2 = \sum_{k=1}^n z_k^2$$

$$\frac{\partial f(z)}{\partial z_j} = \frac{\partial}{\partial z_j} \left(\sum_{k=1}^n z_k^2 \right)$$

$$\boxed{\nabla_z f(z) = 2z_j}$$

Using the above derivative,

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

thru' $\frac{1}{2} \|\cdot\|^2$

→ Back prop

$$\frac{\partial L}{\partial (W^T Wx - x)} = \frac{1}{2} \cdot 2 (W^T Wx - x)$$

$$= \underline{\underline{W^T Wx - x}}$$

→ Back prop thru subtraction, derivative is passed

$$\frac{\partial L}{\partial W^T Wx} = \frac{\partial L}{\partial (W^T Wx - x)} \quad \text{--- (1)}$$

$$\rightarrow \frac{\partial L}{\partial Wx} = (W^T)^T \left(\frac{\partial L}{\partial W^T Wx} \right) \quad \langle \text{using hint} \rangle$$

$$\therefore \frac{\partial L}{\partial Wx} = W (W^T Wx - x)$$

→ Back prop. to W ,

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial Wx} x^T$$

$$\frac{\partial L}{\partial W} = W (W^T Wx - x) x^T$$

From (1), back prop. to W^T

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

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

$$\therefore \nabla_W L = \frac{\partial L}{\partial W} + \left(\frac{\partial L}{\partial W^T} \right)^T$$

⑥

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

$$\therefore \nabla_w L = W W^T W x x^T - W x x^T + W x (x^T W^T W - x^T)$$

$$\begin{aligned} \nabla_w L = & W W^T W x x^T - W x x^T \\ & + W x x^T W^T W - W x x^T \end{aligned}$$

$$\boxed{\nabla_w L = W (W^T W x x^T + x x^T W^T W - 2 x x^T)}$$

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 [50]: 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 [51]: from nndl.neural_net import TwoLayerNet
```

```
In [52]: # 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 [53]: ## Implement the forward pass of the neural network.

# Note, there is a statement if y is None: return scores, which is why
# 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:
3.3812311797665195e-08
```

Forward pass loss

```
In [54]: 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 [55]: print(loss)

1.071696123862817
```

Backward pass

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

```
In [56]: from cs231n.gradient_check import eval_numerical_gradient

# Use numeric gradient checking to check your implementation of the back
ward pass.
# If your implementation is correct, the difference between the numeric
and
# analytic gradients should be less than 1e-8 for each of W1, W2, b1, an
d 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(param
_grad_num, grads[param_name])))

W1 max relative error: 1.28328965625e-09
W2 max relative error: 3.42547269506e-10
b2 max relative error: 1.83916590901e-10
b1 max relative error: 3.1726802857e-09
```

Training the network

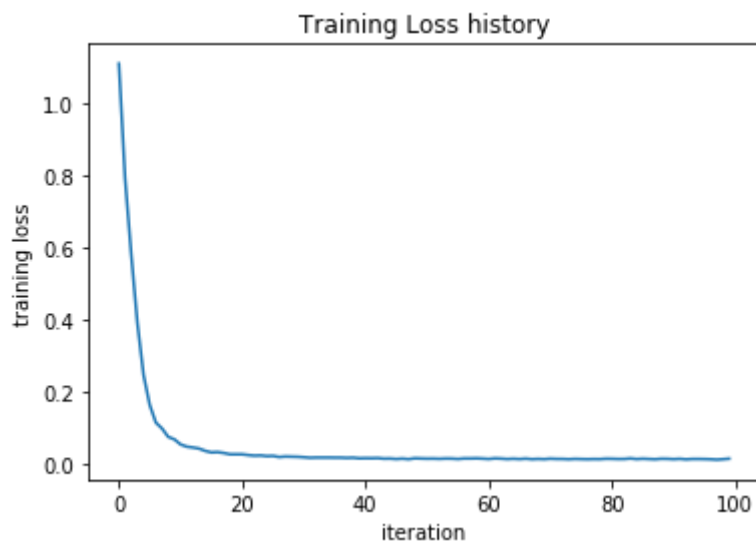
Implement `neural_net.train()` to train the network via stochastic gradient descent, much like the softmax and SVM.

```
In [57]: 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: ', 0.014497864587765906)
```



Classify CIFAR-10

Do classification on the CIFAR-10 dataset.

```

In [58]: 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 [59]: 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 [60]: stats['train_acc_history']
```

```
Out[60]: [0.095, 0.15, 0.25, 0.25, 0.315]
```

```

In [78]: # ===== #
# YOUR CODE HERE:
#   Do some debugging to gain some insight into why the optimization
#   isn't great.
# ===== #

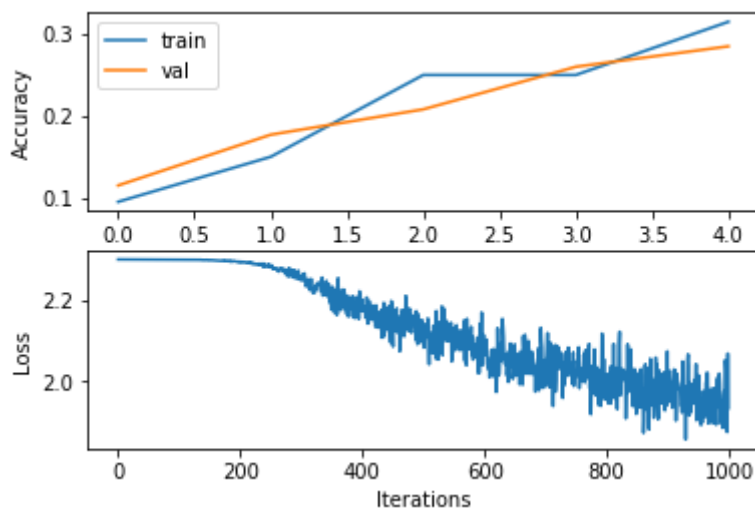
# Plot the loss function and train / validation accuracies
plt.subplot(2,1,1)

plt.ylabel('Accuracy')
plt.plot(stats['train_acc_history'],label='train')
plt.xlabel('Iterations')
plt.plot(stats['val_acc_history'],label='val')
plt.legend()
plt.subplot(2,1,2)
plt.plot(stats['loss_history'])
plt.xlabel('Iterations')
plt.ylabel('Loss')

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

```

Out[78]: <matplotlib.text.Text at 0x10b3bd810>



Answers:

(1) the loss is decreasing linearly, which seems to suggest that the learning rate may be too low. Moreover, there is no gap between the training and validation accuracy, which means that the model we are using is of low capacity.

(2) Change the learning rate (increase) and increase the size of the neural network(number of hidden layers)

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 [87]: 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 b
y:
# 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)!
# ===== #
hidden_size = 50
lr = 1e-3
regularization_strengths = [0.85, 0.89]
best_net = None # store the best model into this
best_stats = None
results = {}
best_val = -1

for reg in regularization_strengths:
    np.random.seed(0)
    print "hidden size: %d, lr: %.4f, reg: %.2f" %(hidden_size, lr, reg)
    net = TwoLayerNet(input_size, hidden_size, num_classes)
    stats = net.train(X_train, y_train, X_val, y_val,
                      num_iters=3000, batch_size=200,
                      learning_rate=lr, learning_rate_decay=0.95,
                      reg=reg, verbose=False)
    print 'train accuracy: %.4f' %stats['train_acc_history'][-1]
    print 'validation accuracy: %.4f' %stats['val_acc_history'][-1]

    # check if validation accuracy is best or not
    if best_val < stats['val_acc_history'][-1]:
        best_val = stats['val_acc_history'][-1]
        best_net = net
        best_stats = stats

print 'best validation accuracy achieved during cross-validation: %f' %
best_val

# ===== #
# END YOUR CODE HERE
# ===== #
best_net = net

hidden size: 50, lr: 0.0010, reg: 0.85
train accuracy: 0.5400
validation accuracy: 0.5040
hidden size: 50, lr: 0.0010, reg: 0.89
train accuracy: 0.5750
validation accuracy: 0.5120
best validation accuracy achieved during cross-validation: 0.512000

```

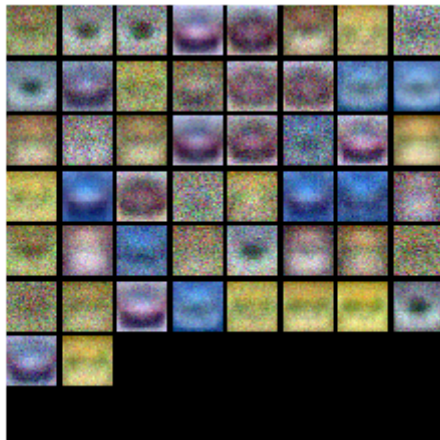


```
In [88]: from cs231n.vis_utils import visualize_grid

# 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) In most neural networks trained on visual data, the first layer weights typically show some visible structure when visualized. In suboptimal net, we are not able to visualize the structure which is clearer with the best net.

Evaluate on test set

```
In [89]: test_acc = (best_net.predict(X_test) == y_test).mean()
          print('Test accuracy: ', test_acc)

('Test accuracy: ', 0.513)
```

```

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         np.random.seed(0)
46         self.params = {}
47         self.params['W1'] = std * np.random.randn(hidden_size, input_size)
48         self.params['b1'] = np.zeros(hidden_size)
49         self.params['W2'] = std * np.random.randn(output_size, hidden_size)
50         self.params['b2'] = np.zeros(output_size)
51
52
53     def loss(self, X, y=None, reg=0.0):
54         """
55         Compute the loss and gradients for a two layer fully connected neural
56         network.
57
58         Inputs:
59         - X: Input data of shape (N, D). Each X[i] is a training sample.
60         - y: Vector of training labels. y[i] is the label for X[i], and each y[i] is
61             an integer in the range 0 <= y[i] < C. This parameter is optional; if it
62             is not passed then we only return scores, and if it is passed then we
63             instead return the loss and gradients.
64         - reg: Regularization strength.
65
66         Returns:
67         If y is None, return a matrix scores of shape (N, C) where scores[i, c] is
68         the score for class c on input X[i].
69

```

```

70 If y is not None, instead return a tuple of:
71 - loss: Loss (data loss and regularization loss) for this batch of training
72 samples.
73 - grads: Dictionary mapping parameter names to gradients of those parameters
74 with respect to the loss function; has the same keys as self.params.
75 """
76 # Unpack variables from the params dictionary
77 W1, b1 = self.params['W1'], self.params['b1']
78 W2, b2 = self.params['W2'], self.params['b2']
79 N, D = X.shape
80
81 # Compute the forward pass
82 scores = None
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 = np.maximum(0, np.dot(W1, X.T) + np.matrix(b1).T)
93 z = np.dot(W2, h1) + np.matrix(b2).T
94 scores = z.T
95 pass
96
97 # ===== #
98 # END YOUR CODE HERE
99 # ===== #
100 softmax = lambda x: np.exp(x - np.max(x)) / np.exp(x - np.max(x)).sum(axis=0)
101 # If the targets are not given then jump out, we're done
102 if y is None:
103     return scores
104
105 # Compute the loss
106 loss = None
107
108 # ===== #
109 # YOUR CODE HERE:
110 # Calculate the loss of the neural network. This includes the
111 # softmax loss and the L2 regularization for W1 and W2. Store the
112 # total loss in the variable loss. Multiply the regularization
113 # loss by 0.5 (in addition to the factor reg).
114 # ===== #
115
116 # scores is num_examples by num_classes
117 aExp = np.exp(scores)
118
119 prob = aExp / np.sum(aExp, axis = 1)
120 correctLogProb = -np.log(prob[range(N), y])
121 dataLoss = np.sum(correctLogProb) / N
122 regloss = 0.5 * reg * np.sum(W1 * W1) + 0.5 * reg * np.sum(W2 * W2)
123
124 loss = regloss + dataLoss
125 # ===== #
126 # END YOUR CODE HERE
127 # ===== #
128
129 grads = {}
130
131 # ===== #
132 # YOUR CODE HERE:
133 # Implement the backward pass. Compute the derivatives of the
134 # weights and the biases. Store the results in the grads
135 # dictionary. e.g., grads['W1'] should store the gradient for
136 # W1, and be of the same size as W1.
137 # ===== #
138 #prob = np.matrix(prob).T
139 #print y.shape[0]

```

```

140 prob = np.arange(y.shape[0], y) - 1
141 prob /= y.shape[0]
142 grads['W2'] = np.dot(h1, prob).T
143 grads['b2'] = np.sum(prob, axis=0)
144
145 dh1 = np.dot(prob, W2).T
146 dh1[h1<=0] = 0
147
148 grads['W1'] = np.dot(X.T, dh1.T).T
149 grads['b1'] = np.sum(dh1, axis=1).T
150
151
152
153 grads['W1'] += reg*W1
154 grads['W2'] += reg*W2
155
156
157 # ===== #
158 # END YOUR CODE HERE
159 # ===== #
160
161 return loss, grads
162
163 def train(self, X, y, X_val, y_val,
164           learning_rate=1e-3, learning_rate_decay=0.95,
165           reg=1e-5, num_iters=100,
166           batch_size=200, verbose=False):
167     """
168     Train this neural network using stochastic gradient descent.
169
170     Inputs:
171     - X: A numpy array of shape (N, D) giving training data.
172     - y: A numpy array of shape (N,) giving training labels; y[i] = c means that
173         X[i] has label c, where 0 <= c < C.
174     - X_val: A numpy array of shape (N_val, D) giving validation data.
175     - y_val: A numpy array of shape (N_val,) giving validation labels.
176     - learning_rate: Scalar giving learning rate for optimization.
177     - learning_rate_decay: Scalar giving factor used to decay the learning rate
178         after each epoch.
179     - reg: Scalar giving regularization strength.
180     - num_iters: Number of steps to take when optimizing.
181     - batch_size: Number of training examples to use per step.
182     - verbose: boolean; if true print progress during optimization.
183     """
184     num_train = X.shape[0]
185     iterations_per_epoch = max(num_train / batch_size, 1)
186
187     # Use SGD to optimize the parameters in self.model
188     loss_history = []
189     train_acc_history = []
190     val_acc_history = []
191
192     for it in np.arange(num_iters):
193         X_batch = None
194         y_batch = None
195
196         # ===== #
197         # YOUR CODE HERE:
198         # Create a minibatch by sampling batch_size samples randomly.
199         # ===== #
200         indic = np.random.choice(num_train, batch_size)
201         X_batch = X[indic,:]
202         y_batch = y[indic]
203
204         # ===== #
205         # END YOUR CODE HERE
206         # ===== #
207
208         # Compute loss and gradients using the current minibatch
209         loss, grads = self.loss(X_batch, y=y_batch, reg=reg)

```

```

210     loss_history.append(loss)
211
212     # ===== #
213     # YOUR CODE HERE:
214     #     Perform a gradient descent step using the minibatch to update
215     #     all parameters (i.e., W1, W2, b1, and b2).
216     # ===== #
217     self.params['W1'] -= learning_rate*grads['W1']
218     self.params['W2'] -= learning_rate*grads['W2']
219     self.params['b1'] -= learning_rate*np.asarray(grads['b1']).reshape(-1)
220     self.params['b2'] -= learning_rate*np.asarray(grads['b2']).reshape(-1)
221
222
223     # ===== #
224     # END YOUR CODE HERE
225     # ===== #
226
227     if verbose and it % 100 == 0:
228         print('iteration {} / {}: loss {}'.format(it, num_iters, loss))
229
230     # Every epoch, check train and val accuracy and decay learning rate.
231     if it % iterations_per_epoch == 0:
232         # Check accuracy
233         train_acc = (self.predict(X_batch) == y_batch).mean()
234         val_acc = (self.predict(X_val) == y_val).mean()
235         train_acc_history.append(train_acc)
236         val_acc_history.append(val_acc)
237
238         # Decay learning rate
239         learning_rate *= learning_rate_decay
240
241     return {
242         'loss_history': loss_history,
243         'train_acc_history': train_acc_history,
244         'val_acc_history': val_acc_history,
245     }
246
247 def predict(self, X):
248     """
249     Use the trained weights of this two-layer network to predict labels for
250     data points. For each data point we predict scores for each of the C
251     classes, and assign each data point to the class with the highest score.
252
253     Inputs:
254     - X: A numpy array of shape (N, D) giving N D-dimensional data points to
255         classify.
256
257     Returns:
258     - y_pred: A numpy array of shape (N,) giving predicted labels for each of
259         the elements of X. For all i, y_pred[i] = c means that X[i] is predicted
260         to have class c, where 0 <= c < C.
261     """
262     y_pred = None
263
264     # ===== #
265     # YOUR CODE HERE:
266     #     Predict the class given the input data.
267     # ===== #
268     #reluFunc = lambda x: np.multiply(x,(x>0))
269     z = np.dot(X, self.params['W1'].T) + self.params['b1']
270     h1 = np.maximum(0, z)
271     out = np.dot(h1, self.params['W2'].T) + self.params['b2']
272
273     y_pred = np.argmax(out, axis = 1)
274
275
276     # ===== #
277     # END YOUR CODE HERE
278     # ===== #

```

2/6/2018

/Users/vijayravi/Documents/UCLA/Coursework/2018Winter/neuralNetworks/homeworks/homework3/code/nndl/neural_net.py

279

280 **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).

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 [1]: ## 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))))

```

```

In [2]: # 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 [3]: # 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: 9.76985004799e-10
```

Affine layer backward pass

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

```
In [5]: # 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: 1.09375090808e-10
dw error: 1.32890806001e-10
db error: 3.27574581513e-12
```

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 [6]: # 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, 0.136
36364, ],
                        [ 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 [7]: 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: 3.27560333592e-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 [8]: 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: 5.51287985683e-10
dw error: 5.62139289988e-10
db error: 3.27559030827e-12
```

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 [9]: 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: 9.00006906888
dx error: 1.40215660067e-09

Testing softmax_loss:
loss: 2.30259244752
dx error: 9.60422880538e-09

```

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 [15]: 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, 14.5719843
4, 15.33206765, 16.09215096],
    [12.05769098, 12.74614105, 13.43459113, 14.1230412, 14.8114912
8, 15.49994135, 16.18839143],
    [12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.0509982
2, 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 l
oss'

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, gra
ds[name])))

```



```
Testing initialization ...
Testing test-time forward pass ...
Testing training loss (no regularization)
Running numeric gradient check with reg = 0.0
W1 relative error: 2.13161195546e-08
W2 relative error: 3.31027019978e-10
b1 relative error: 8.36819673248e-09
b2 relative error: 2.53077405016e-10
Running numeric gradient check with reg = 0.7
W1 relative error: 2.52791534132e-07
W2 relative error: 2.85086969908e-08
b1 relative error: 1.56468020339e-08
b2 relative error: 9.08961463813e-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 [18]: 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.
# ===== #
data = {
    'X_train': data['X_train'],
    'y_train': data['y_train'],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

std = 1e-2

model = TwoLayerNet(hidden_dims=200)
solver = Solver(model, data,
                update_rule='sgd',
                optim_config={
                    'learning_rate': 1e-3,
                },
                lr_decay=0.95,
                num_epochs=10, batch_size=100,
                print_every=100)
solver.train()

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

```

```

(Iteration 1 / 4900) loss: 2.302490
(Epoch 0 / 10) train acc: 0.159000; val_acc: 0.144000
(Iteration 101 / 4900) loss: 1.696819
(Iteration 201 / 4900) loss: 1.595609
(Iteration 301 / 4900) loss: 1.633681
(Iteration 401 / 4900) loss: 1.514735
(Epoch 1 / 10) train acc: 0.442000; val_acc: 0.438000
(Iteration 501 / 4900) loss: 1.521922
(Iteration 601 / 4900) loss: 1.481645
(Iteration 701 / 4900) loss: 1.392719
(Iteration 801 / 4900) loss: 1.262135
(Iteration 901 / 4900) loss: 1.203570
(Epoch 2 / 10) train acc: 0.515000; val_acc: 0.466000

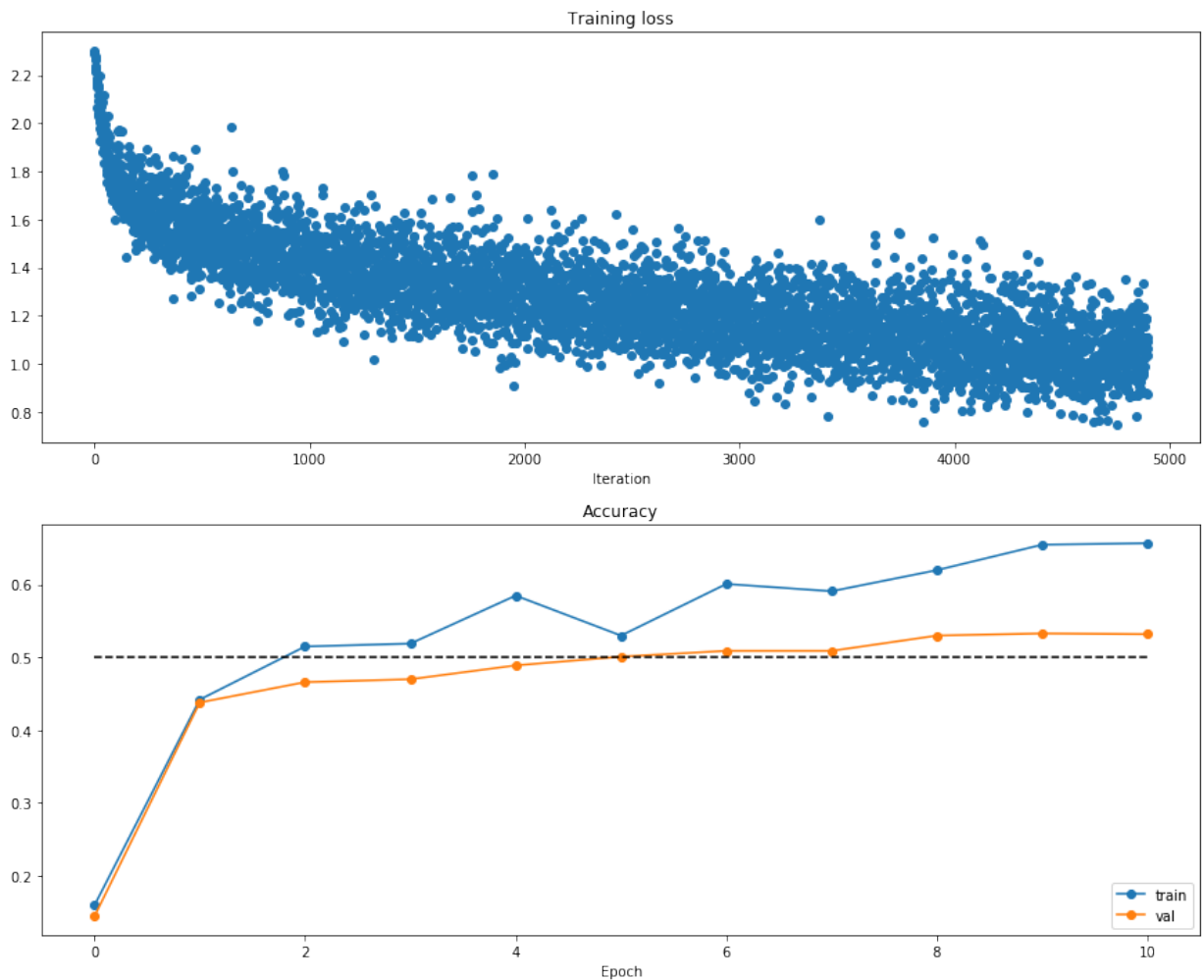
```

```
(Iteration 1001 / 4900) loss: 1.631719
(Iteration 1101 / 4900) loss: 1.220347
(Iteration 1201 / 4900) loss: 1.423363
(Iteration 1301 / 4900) loss: 1.348886
(Iteration 1401 / 4900) loss: 1.409443
(Epoch 3 / 10) train acc: 0.519000; val_acc: 0.470000
(Iteration 1501 / 4900) loss: 1.372770
(Iteration 1601 / 4900) loss: 1.479888
(Iteration 1701 / 4900) loss: 1.153154
(Iteration 1801 / 4900) loss: 1.314581
(Iteration 1901 / 4900) loss: 1.251971
(Epoch 4 / 10) train acc: 0.585000; val_acc: 0.489000
(Iteration 2001 / 4900) loss: 1.228331
(Iteration 2101 / 4900) loss: 1.302921
(Iteration 2201 / 4900) loss: 1.297642
(Iteration 2301 / 4900) loss: 1.312953
(Iteration 2401 / 4900) loss: 1.155149
(Epoch 5 / 10) train acc: 0.530000; val_acc: 0.501000
(Iteration 2501 / 4900) loss: 1.259114
(Iteration 2601 / 4900) loss: 1.156750
(Iteration 2701 / 4900) loss: 1.253140
(Iteration 2801 / 4900) loss: 1.232998
(Iteration 2901 / 4900) loss: 1.286295
(Epoch 6 / 10) train acc: 0.601000; val_acc: 0.509000
(Iteration 3001 / 4900) loss: 1.241539
(Iteration 3101 / 4900) loss: 1.217068
(Iteration 3201 / 4900) loss: 1.039278
(Iteration 3301 / 4900) loss: 1.075537
(Iteration 3401 / 4900) loss: 1.066222
(Epoch 7 / 10) train acc: 0.591000; val_acc: 0.509000
(Iteration 3501 / 4900) loss: 1.090363
(Iteration 3601 / 4900) loss: 1.225429
(Iteration 3701 / 4900) loss: 1.193425
(Iteration 3801 / 4900) loss: 0.995088
(Iteration 3901 / 4900) loss: 1.033185
(Epoch 8 / 10) train acc: 0.620000; val_acc: 0.530000
(Iteration 4001 / 4900) loss: 1.188090
(Iteration 4101 / 4900) loss: 0.974699
(Iteration 4201 / 4900) loss: 1.239576
(Iteration 4301 / 4900) loss: 1.151068
(Iteration 4401 / 4900) loss: 1.005735
(Epoch 9 / 10) train acc: 0.655000; val_acc: 0.533000
(Iteration 4501 / 4900) loss: 1.040810
(Iteration 4601 / 4900) loss: 0.896869
(Iteration 4701 / 4900) loss: 0.804847
(Iteration 4801 / 4900) loss: 0.936875
(Epoch 10 / 10) train acc: 0.657000; val_acc: 0.532000
```

In [19]: *# 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 [38]: 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.float
64)

    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, gra
ds[name])))
```

```
Running check with reg = 0
Initial loss: 2.29782803292
W1 relative error: 1.51411949888e-07
W2 relative error: 3.14968357984e-06
W3 relative error: 2.01298510858e-07
b1 relative error: 1.04128304704e-08
b2 relative error: 9.76663118471e-08
b3 relative error: 1.72535524188e-10
Running check with reg = 3.14
Initial loss: 6.75808296138
W1 relative error: 8.76162821278e-09
W2 relative error: 4.47343936753e-08
W3 relative error: 9.75892478667e-07
b1 relative error: 3.90889119541e-08
b2 relative error: 8.56314343458e-09
b3 relative error: 2.94139688503e-10
```

```
In [43]: # Use the three layer neural network to overfit a small dataset.

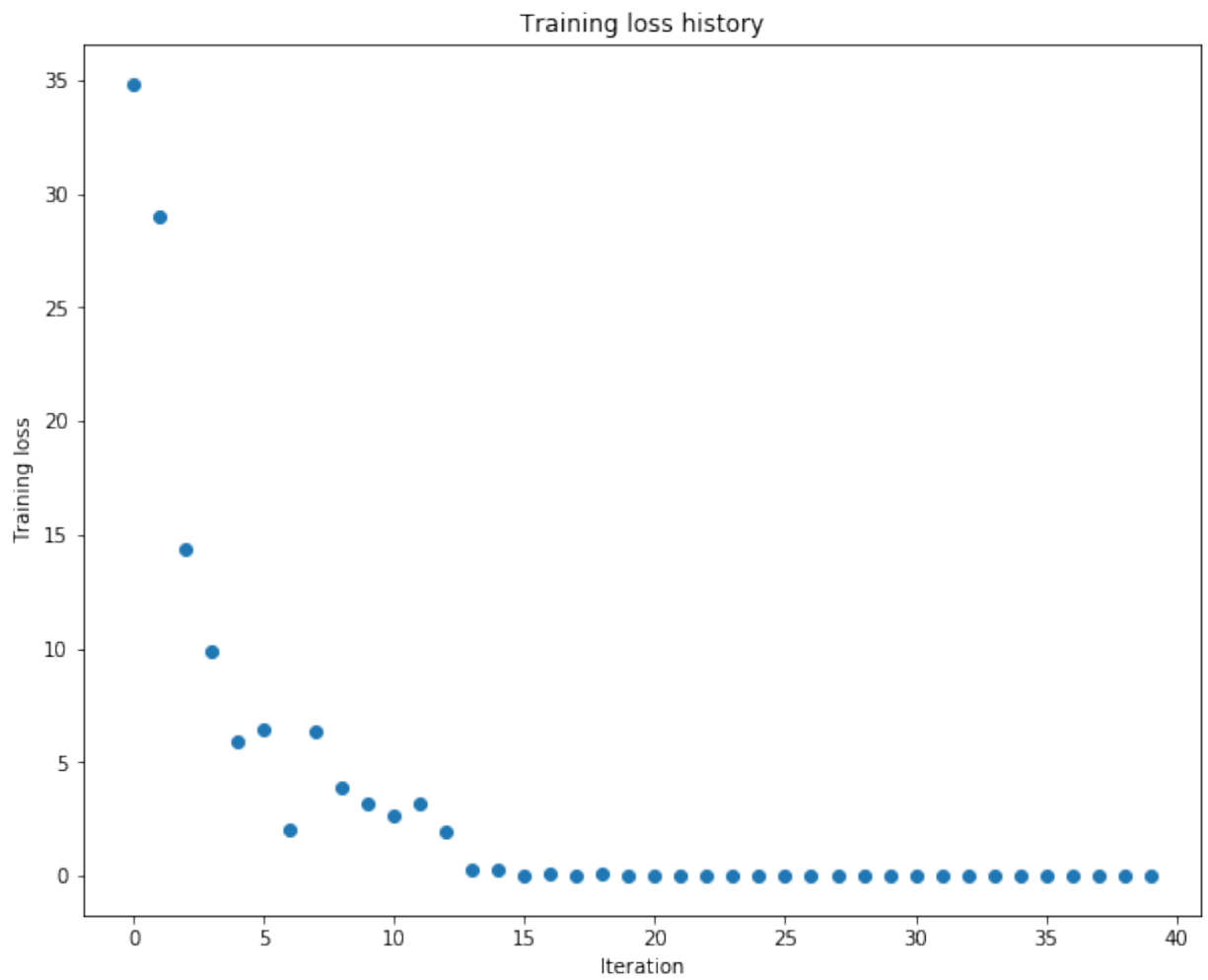
num_train = 50
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

#### !!!!!
# 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 = 5e-2
learning_rate = 5e-4

model = FullyConnectedNet([100, 100],
                           weight_scale=weight_scale, dtype=np.float64)
solver = Solver(model, small_data,
                 print_every=10, num_epochs=20, batch_size=25,
                 update_rule='sgd',
                 optim_config={
                     'learning_rate': learning_rate,
                 })
solver.train()
```

```
(Iteration 1 / 40) loss: 34.848935
(Epoch 0 / 20) train acc: 0.180000; val_acc: 0.113000
(Epoch 1 / 20) train acc: 0.220000; val_acc: 0.144000
(Epoch 2 / 20) train acc: 0.360000; val_acc: 0.139000
(Epoch 3 / 20) train acc: 0.580000; val_acc: 0.136000
(Epoch 4 / 20) train acc: 0.560000; val_acc: 0.143000
(Epoch 5 / 20) train acc: 0.740000; val_acc: 0.149000
(Iteration 11 / 40) loss: 2.633929
(Epoch 6 / 20) train acc: 0.700000; val_acc: 0.162000
(Epoch 7 / 20) train acc: 0.940000; val_acc: 0.145000
(Epoch 8 / 20) train acc: 1.000000; val_acc: 0.151000
(Epoch 9 / 20) train acc: 1.000000; val_acc: 0.153000
(Epoch 10 / 20) train acc: 1.000000; val_acc: 0.152000
(Iteration 21 / 40) loss: 0.006856
(Epoch 11 / 20) train acc: 1.000000; val_acc: 0.153000
(Epoch 12 / 20) train acc: 1.000000; val_acc: 0.153000
(Epoch 13 / 20) train acc: 1.000000; val_acc: 0.153000
(Epoch 14 / 20) train acc: 1.000000; val_acc: 0.153000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.151000
(Iteration 31 / 40) loss: 0.010032
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.151000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.151000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.150000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.150000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.151000
```

```
In [44]: plt.plot(solver.loss_history, 'o')
plt.title('Training loss history')
plt.xlabel('Iteration')
plt.ylabel('Training loss')
plt.show()
```




```

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         self.cache = {}
49
50         # ===== #
51         # YOUR CODE HERE:
52         # Initialize W1, W2, b1, and b2. Store these as self.params['W1'],
53         # self.params['W2'], self.params['b1'] and self.params['b2']. The
54         # biases are initialized to zero and the weights are initialized
55         # so that each parameter has mean 0 and standard deviation weight_scale.
56         # The dimensions of W1 should be (input_dim, hidden_dim) and the
57         # dimensions of W2 should be (hidden_dims, num_classes)
58         # ===== #
59
60         np.random.seed(0)
61         self.params = {}
62         self.params['W1'] = weight_scale * np.random.randn(input_dim, hidden_dims)
63         self.params['b1'] = np.zeros(hidden_dims)
64         self.params['W2'] = weight_scale * np.random.randn(hidden_dims, num_classes)
65         self.params['b2'] = np.zeros(num_classes)
66
67         # ===== #
68         # END YOUR CODE HERE
69         # ===== #
70

```

```

71 def loss(self, X, y=None):
72     """
73     Compute loss and gradient for a minibatch of data.
74
75     Inputs:
76     - X: Array of input data of shape (N, d_1, ..., d_k)
77     - y: Array of labels, of shape (N,). y[i] gives the label for X[i].
78
79     Returns:
80     If y is None, then run a test-time forward pass of the model and return:
81     - scores: Array of shape (N, C) giving classification scores, where
82       scores[i, c] is the classification score for X[i] and class c.
83
84     If y is not None, then run a training-time forward and backward pass and
85     return a tuple of:
86     - loss: Scalar value giving the loss
87     - grads: Dictionary with the same keys as self.params, mapping parameter
88       names to gradients of the loss with respect to those parameters.
89     """
90     scores = None
91     W1, b1 = self.params['W1'], self.params['b1']
92     W2, b2 = self.params['W2'], self.params['b2']
93     N = X.shape[0]
94     D = np.prod(X.shape[1:])
95     # ===== #
96     # YOUR CODE HERE:
97     # Implement the forward pass of the two-layer neural network. Store
98     # the class scores as the variable 'scores'. Be sure to use the layers
99     # you prior implemented.
100    # ===== #
101
102    out1, cache1 = affine_relu_forward(X, W1, b1)
103
104    out2, cache2 = affine_forward(out1, W2, b2)
105    scores = out2
106
107    # ===== #
108    # END YOUR CODE HERE
109    # ===== #
110
111    # If y is None then we are in test mode so just return scores
112    if y is None:
113        return scores
114
115    loss, grads = 0, {}
116    # ===== #
117    # YOUR CODE HERE:
118    # Implement the backward pass of the two-layer neural net. Store
119    # the loss as the variable 'loss' and store the gradients in the
120    # 'grads' dictionary. For the grads dictionary, grads['W1'] holds
121    # the gradient for W1, grads['b1'] holds the gradient for b1, etc.
122    # i.e., grads[k] holds the gradient for self.params[k].
123    #
124    # Add L2 regularization, where there is an added cost 0.5*self.reg*W^2
125    # for each W. Be sure to include the 0.5 multiplying factor to
126    # match our implementation.
127    #
128    # And be sure to use the layers you prior implemented.
129    # ===== #
130
131    loss, dscore = softmax_loss(out2, y)
132    loss += 0.5 * self.reg * np.sum(W1*W1) + 0.5*self.reg * np.sum(W2*W2)
133
134    dx1, grads['W2'], grads['b2'] = affine_backward(dscore, cache2)
135    _, grads['W1'], grads['b1'] = affine_relu_backward(dx1, cache1)
136
137
138    grads['W2'] += self.reg * W2
139    grads['W1'] += self.reg * W1
140    # ===== #
141    # END YOUR CODE HERE

```

```

142 # ===== #
143
144 return loss, grads
145
146
147 class FullyConnectedNet(object):
148     """
149     A fully-connected neural network with an arbitrary number of hidden layers,
150     ReLU nonlinearities, and a softmax loss function. This will also implement
151     dropout and batch normalization as options. For a network with L layers,
152     the architecture will be
153
154     {affine - [batch norm] - relu - [dropout]} x (L - 1) - affine - softmax
155
156     where batch normalization and dropout are optional, and the {...} block is
157     repeated L - 1 times.
158
159     Similar to the TwoLayerNet above, learnable parameters are stored in the
160     self.params dictionary and will be learned using the Solver class.
161     """
162
163     def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,
164                 dropout=0, use_batchnorm=False, reg=0.0,
165                 weight_scale=1e-2, dtype=np.float32, seed=None):
166         """
167         Initialize a new FullyConnectedNet.
168
169         Inputs:
170         - hidden_dims: A list of integers giving the size of each hidden layer.
171         - input_dim: An integer giving the size of the input.
172         - num_classes: An integer giving the number of classes to classify.
173         - dropout: Scalar between 0 and 1 giving dropout strength. If dropout=0 then
174           the network should not use dropout at all.
175         - use_batchnorm: Whether or not the network should use batch normalization.
176         - reg: Scalar giving L2 regularization strength.
177         - weight_scale: Scalar giving the standard deviation for random
178           initialization of the weights.
179         - dtype: A numpy datatype object; all computations will be performed using
180           this datatype. float32 is faster but less accurate, so you should use
181           float64 for numeric gradient checking.
182         - seed: If not None, then pass this random seed to the dropout layers. This
183           will make the dropout layers deterministic so we can gradient check the
184           model.
185         """
186         self.use_batchnorm = use_batchnorm
187         self.use_dropout = dropout > 0
188         self.reg = reg
189         self.num_layers = 1 + len(hidden_dims)
190         self.dtype = dtype
191         self.params = {}
192
193         # ===== #
194         # YOUR CODE HERE:
195         # Initialize all parameters of the network in the self.params dictionary.
196         # The weights and biases of layer 1 are W1 and b1; and in general the
197         # weights and biases of layer i are Wi and bi. The
198         # biases are initialized to zero and the weights are initialized
199         # so that each parameter has mean 0 and standard deviation weight_scale.
200         # ===== #
201         for i in range(self.num_layers):
202             if i == 0:
203                 # weights between input and hidden layer
204                 self.params['W' + str(i+1)] = weight_scale * np.random.randn(input_dim, hidden_dims[i])
205                 self.params['b' + str(i+1)] = np.zeros(hidden_dims[i])
206
207                 # weights between hidden and hidden layer
208             elif i < self.num_layers - 1:
209                 self.params['W' + str(i+1)] = weight_scale * np.random.randn(hidden_dims[i-1], hidden_dims[i])
210                 self.params['b' + str(i+1)] = np.zeros(hidden_dims[i])
211
212             # weights between hidden and output layer

```

```

213     else:
214         self.params['W' + str(i+1)] = weight_scale * np.random.randn(hidden_dims[i-1], num_classes)
215         self.params['b' + str(i+1)] = np.zeros(num_classes)
216
217
218     # ===== #
219     # END YOUR CODE HERE
220     # ===== #
221
222     # When using dropout we need to pass a dropout_param dictionary to each
223     # dropout layer so that the layer knows the dropout probability and the mode
224     # (train / test). You can pass the same dropout_param to each dropout layer.
225     self.dropout_param = {}
226     if self.use_dropout:
227         self.dropout_param = {'mode': 'train', 'p': dropout}
228         if seed is not None:
229             self.dropout_param['seed'] = seed
230
231     # With batch normalization we need to keep track of running means and
232     # variances, so we need to pass a special bn_param object to each batch
233     # normalization layer. You should pass self.bn_params[0] to the forward pass
234     # of the first batch normalization layer, self.bn_params[1] to the forward
235     # pass of the second batch normalization layer, etc.
236     self.bn_params = []
237     if self.use_batchnorm:
238         self.bn_params = [{'mode': 'train'} for i in np.arange(self.num_layers - 1)]
239
240     # Cast all parameters to the correct datatype
241     for k, v in self.params.items():
242         self.params[k] = v.astype(dtype)
243
244
245 def loss(self, X, y=None):
246     """
247     Compute loss and gradient for the fully-connected net.
248
249     Input / output: Same as TwoLayerNet above.
250     """
251     X = X.astype(self.dtype)
252     mode = 'test' if y is None else 'train'
253
254     # Set train/test mode for batchnorm params and dropout param since they
255     # behave differently during training and testing.
256     if self.dropout_param is not None:
257         self.dropout_param['mode'] = mode
258     if self.use_batchnorm:
259         for bn_param in self.bn_params:
260             bn_param['mode'] = mode
261
262     scores = None
263     caches = []
264     # ===== #
265     # YOUR CODE HERE:
266     # Implement the forward pass of the FC net and store the output
267     # scores as the variable "scores".
268     # ===== #
269
270     data_in = X
271     for i in range(self.num_layers):
272
273         if i < self.num_layers - 1:
274             out, cache = affine_relu_forward(data_in, self.params['W'+str(i+1)], self.params['b'+str(i+1)])
275             data_in = out
276             caches.append(cache)
277         else:
278             out2, cache2 = affine_forward(data_in, self.params['W'+str(i+1)], self.params['b'+str(i+1)])
279
280     scores = out2
281
282
283     # ===== #

```

```

284 # END YOUR CODE HERE
285 # ===== #
286
287 # If test mode return early
288 if mode == 'test':
289     return scores
290
291 loss, grads = 0.0, {}
292 # ===== #
293 # YOUR CODE HERE:
294 # Implement the backwards pass of the FC net and store the gradients
295 # in the grads dict, so that grads[k] is the gradient of self.params[k]
296 # Be sure your L2 regularization includes a 0.5 factor.
297 # ===== #
298
299 loss, dscore = softmax_loss(out2, y)
300 regularization = 0
301 for i in range(self.num_layers):
302     #regularization += 0.5 * self.reg * np.sum(self.params['W'+str(i+1)])
303     loss += 0.5 * self.reg * np.sum(self.params['W'+str(i+1)] * self.params['W'+str(i+1)])
304
305 for i in reversed(range(self.num_layers)):
306     if i == self.num_layers - 1:
307         dx1, grads['W'+str(i+1)], grads['b'+str(i+1)] = affine_backward(dscore, cache2)
308     else:
309         dx, grads['W'+str(i+1)], grads['b'+str(i+1)] = affine_relu_backward(dx1, caches[i])
310         dx1 = dx
311
312 for i in reversed(range(self.num_layers)):
313     grads['W'+str(i+1)] += self.reg * self.params['W'+str(i+1)]
314
315
316 # ===== #
317 # END YOUR CODE HERE
318 # ===== #
319 return loss, grads

```