

- 1) Backprop for autoencoders.

data compression

$$x \in \mathbb{R}^n, w \in \mathbb{R}^{m \times n}, m < n$$

Then Wx is of lower dim than x .

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

$$\frac{\partial}{\partial W} \quad \uparrow \\ \text{linear case}$$

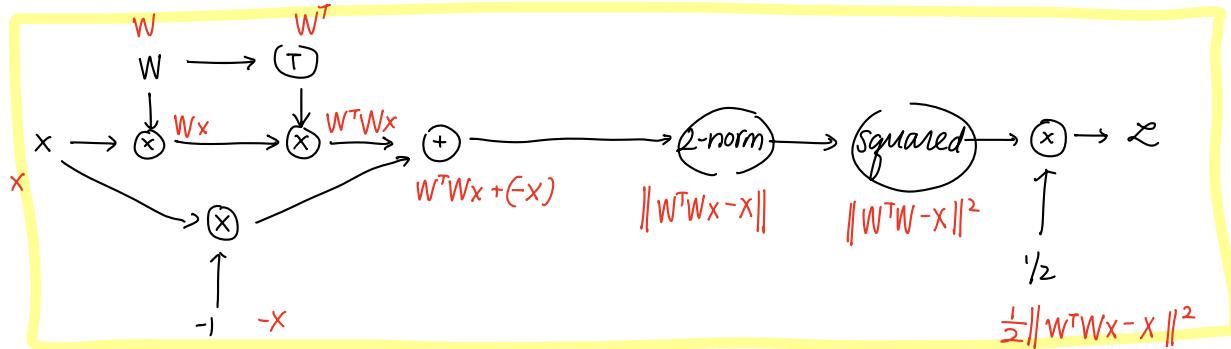
⊕ : pass through gradient

⊗ : multiply by (-1) from previous wire

- a) How does W preserve info about x ?

W reduces the dim of x , W^T increases dim of x to reconstruct it to original dim. Minimizing the difference between the reduced then reconstructed x and the original x preserves info about x , s.t. the $W^T W$ changes as little info as possible.

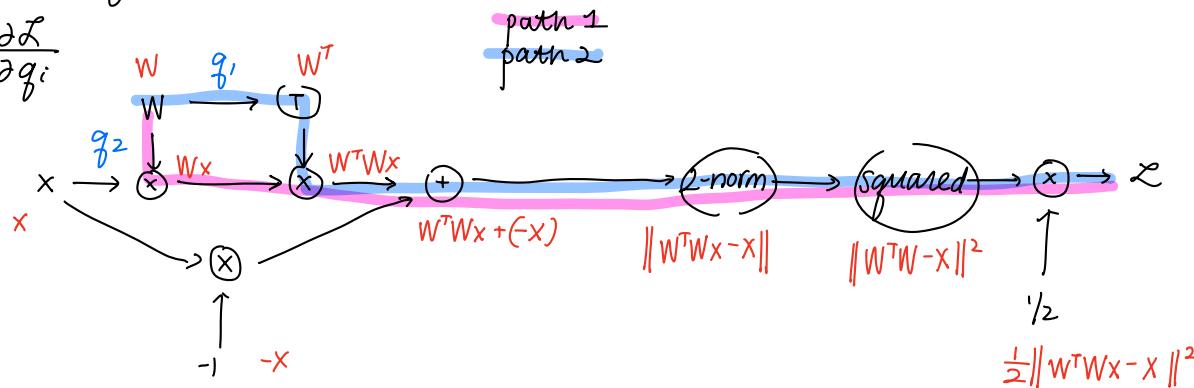
- b) Draw computational graph for \mathcal{L}



- c) Two paths to W , how to account for this?

In lecture, we saw in general,

$$\frac{\partial \mathcal{L}}{\partial x} = \sum_{i=1}^n \frac{\partial g_i}{\partial x} \frac{\partial \mathcal{L}}{\partial g_i}$$

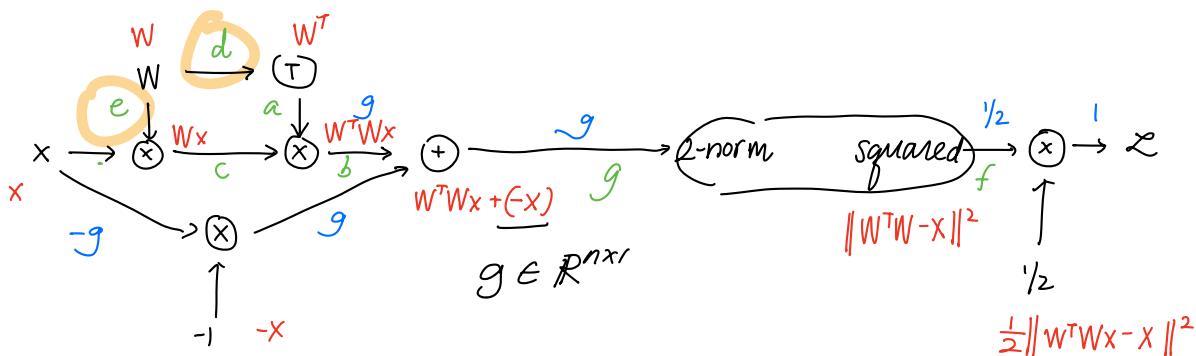


Then, $\frac{\partial \mathcal{L}}{\partial W} = (\text{path 1}) + (\text{path 2})$

$$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial g_1}{\partial W} \cdot \frac{\partial \mathcal{L}}{\partial g_1} + \frac{\partial g_2}{\partial W} \cdot \frac{\partial \mathcal{L}}{\partial g_2}$$

, where g_1, g_2 represent the gradients of the two paths to W .

d)



$$\frac{\partial f}{\partial g} = \frac{\partial \|g\|_2^2}{\partial g} = \frac{\partial g^T g}{\partial g} = 2g \cdot \frac{\partial \mathcal{L}}{\partial f} = 2g \cdot \frac{1}{2} = g$$

a)

$$\frac{\partial \mathcal{L}}{\partial a} = \frac{\partial b}{\partial a} \cdot \left(\frac{\partial \mathcal{L}}{\partial b} \right)$$

trick to switch dimensions

$$\frac{\partial b}{\partial a} = \frac{\partial \underbrace{aWx}_{\substack{n \times 1 \text{ vector} \\ n \times m \text{ matrix}}} \cdot R}{\partial a} = \underbrace{R}_{\substack{3D \text{ tensor} \\ n \times m \text{ matrix}}} \rightarrow$$

$$\frac{\partial \mathcal{L}}{\partial W^T} = \frac{\partial \mathcal{L}}{\partial XY} Y^T \quad \begin{matrix} g \\ \text{where} \\ X = a = W^T \\ Y = Wx \end{matrix} \quad \begin{matrix} \frac{\partial XY}{\partial X} \\ \dim W = m \times n \\ \dim W^T = n \times m \end{matrix}$$

$$= \frac{\partial \mathcal{L}}{\partial W^T Wx} \quad b = W^T Wx \quad \begin{matrix} (n \times m) \times (m \times n) \\ (n \times n) \times (n \times 1) \\ = (n \times 1) \end{matrix}$$

$$\frac{\partial \mathcal{L}}{\partial a} = g x^T W^T$$

$$\frac{\partial \mathcal{L}}{\partial c} = \frac{\partial b}{\partial c} \cdot \left(\frac{\partial \mathcal{L}}{\partial b} \right)$$

$$\frac{\partial \mathcal{L}}{\partial c} = Wg$$

$$\frac{\partial b}{\partial c} = \frac{\partial \underbrace{W^T Wx}_{Y} \cdot Z}{\partial Wx}$$

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

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

$$= W \cdot g$$

$$\frac{\partial \mathcal{L}}{\partial X} = \frac{\partial \mathcal{L}}{\partial Z} Y^T$$

$$\frac{\partial \mathcal{L}}{\partial Y} = X^T \frac{\partial \mathcal{L}}{\partial Z}$$

$$\frac{\partial \mathcal{L}}{\partial K} = -K^T \frac{\partial \mathcal{L}}{\partial Z} K^{-T}$$

$$\frac{\partial \mathcal{L}}{\partial d} = \left(\frac{\partial \mathcal{L}}{\partial a} \right)^T = (g x^T W^T)^T = \boxed{W x g^T} = \frac{\partial \mathcal{L}}{\partial d}$$

$$e) \frac{\partial \mathcal{L}}{\partial e} = \frac{\partial C}{\partial e} \cdot \frac{\partial \mathcal{L}}{\partial C}$$

\downarrow \downarrow

$X^T Wg$
switch order

$$\frac{\partial C}{\partial e} = \frac{\partial Wx}{\partial W} = X^T$$

\overbrace{X}

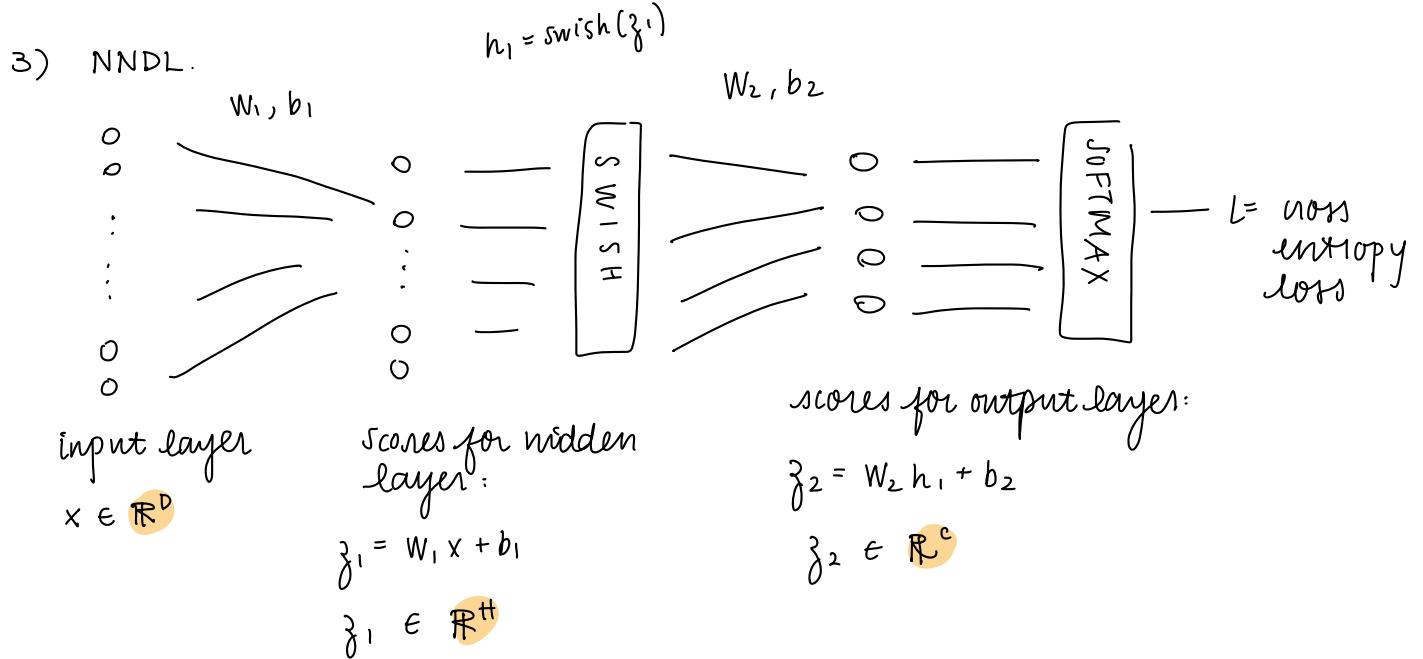
$$\frac{\partial \mathcal{L}}{\partial e} = \underbrace{Wg X^T}_{(m \times n) \times (n \times 1) \times (1 \times n)} \in \mathbb{R}^{m \times n}$$

We have $\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial d} + \frac{\partial \mathcal{L}}{\partial e}$

$$= Wxg^T + X^T Wg \quad , \text{ where } g = W^T Wx - X$$

$$\nabla_W \mathcal{L} = Wx(W^T Wx - X)^T + W(W^T Wx - X)X^T$$

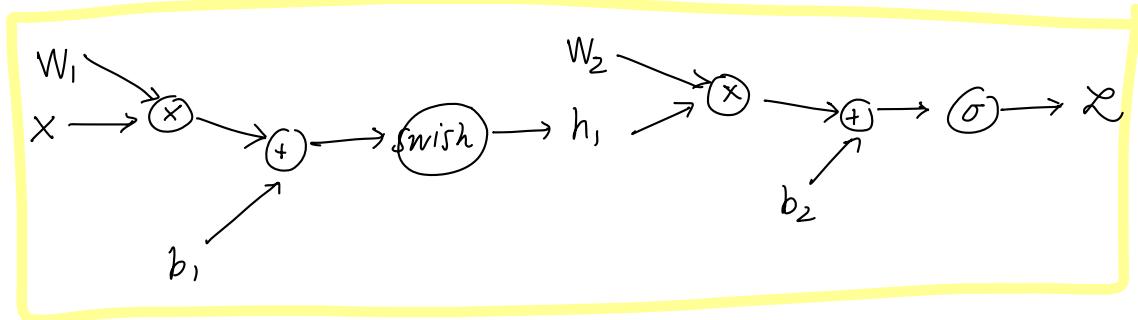
2) I am a c147 student. (Do this later)



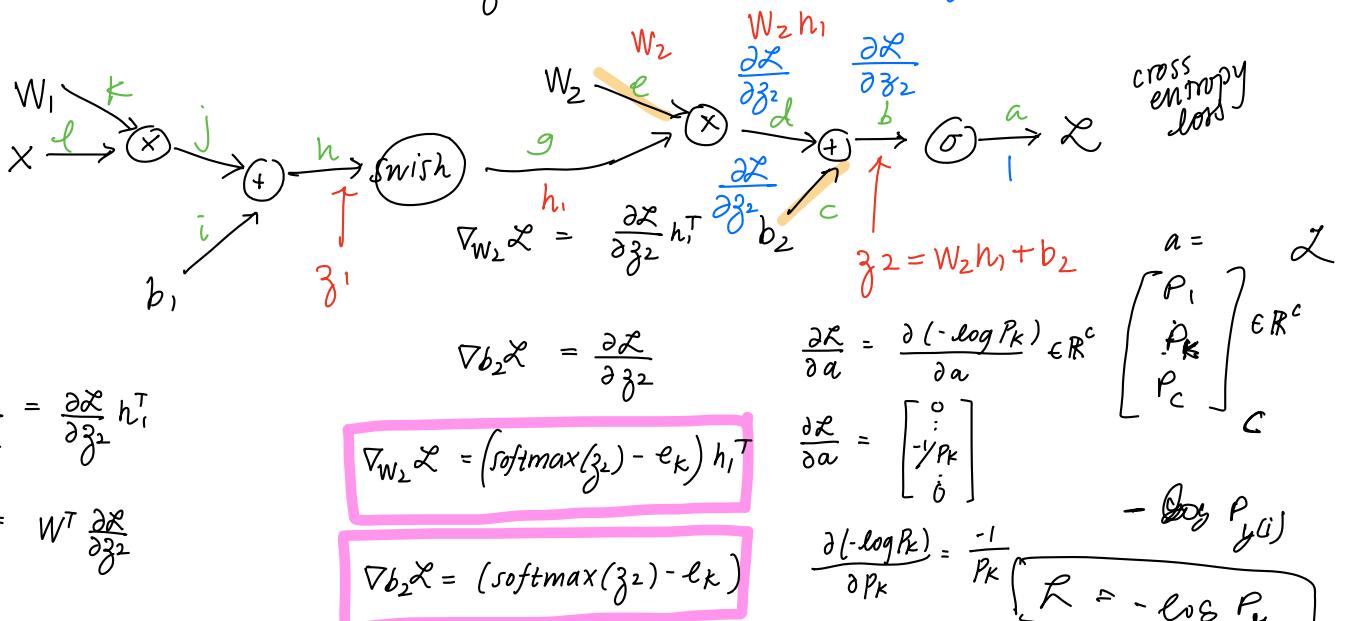
$$\text{swish}(k) = \frac{k}{1-e^{-k}} = k \delta(k), \text{ where } \delta(k) \text{ is the sigmoid activation func from lecture}$$

final answer in terms of $\frac{\partial \mathcal{L}}{\partial z_2}$

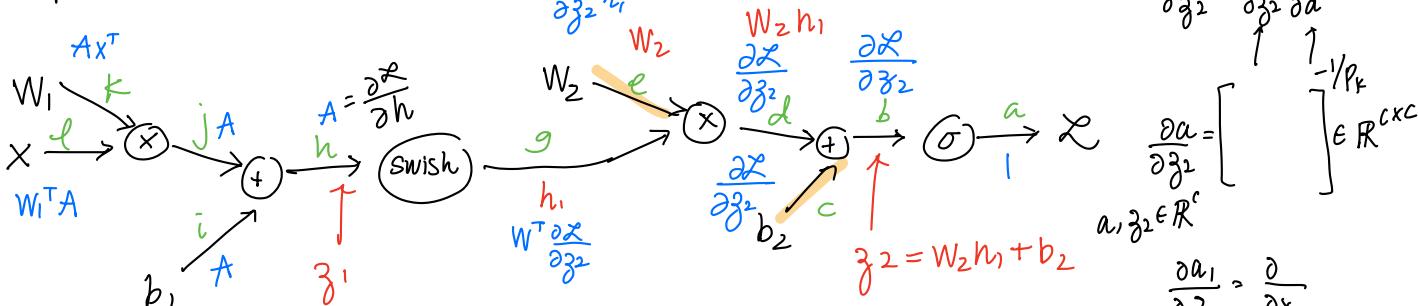
a) Draw the computational graph.



b) Compute $\nabla_{w_2} L$, $\nabla_{b_2} L$: $\frac{\partial \mathcal{L}}{\partial z_2} = \text{softmax}(z_2) - e_k$ * gradients



c) compute $\nabla_w L$, $\nabla_b L$:



$$\text{swish}(k) = \frac{k}{1 + e^{-k}}$$

$$\frac{\partial \mathcal{L}}{\partial h} = \frac{\partial g}{\partial h} \cdot \frac{\partial \mathcal{L}}{\partial g} \quad \frac{\partial g}{\partial h} = \frac{\partial}{\partial h} \left(\frac{h}{1+e^{-h}} \right) = 1 \cdot (1+e^{-h}) + h(-e^{-h})$$

$$= \frac{1 \cdot (1+e^{-h}) + h(e^{-h})}{(1+e^{-h})^2}$$

by product rule

$$\frac{gf' - fg'}{g^2}$$

$$\begin{bmatrix} D(1-\delta) & & \\ & \ddots & \\ & & D(1-\delta) \end{bmatrix} = \begin{bmatrix} 0 & & \\ & \vdots & \\ & -1/p_k & \\ & & 0 \end{bmatrix} = \begin{bmatrix} 0 & & \\ & \vdots & \\ & p_{k-1} & \\ & & 0 \end{bmatrix} = \frac{\partial \mathcal{L}}{\partial \hat{y}_2}$$

$$P_X(1-P_K)\left(\frac{1}{P_K}\right) = -1(1-P_K) = P_K - 1$$

$$\frac{\partial g}{\partial h} = \frac{1 - e^{-\hat{z}_2}}{(1 - e^{-\hat{z}_2})^2}$$

$$\frac{\partial z}{\partial h} = \left(\frac{1 - e^{-\hat{z}_2}}{(1 - e^{-\hat{z}_2})^2} \right) W \frac{\partial \hat{z}_2}{\partial z_2}$$

$$A = \frac{\partial \mathcal{L}}{\partial h}$$

$$a = \begin{bmatrix} p_1 \\ \vdots \\ p_c \end{bmatrix}_{c \times 1} \quad \hat{z}_2 = \begin{bmatrix} x_1 \\ \vdots \\ x_c \end{bmatrix}_{c \times 1}$$

$$\frac{\partial a}{\partial \hat{z}_2} = \begin{bmatrix} \frac{\partial p_1}{\partial x_1} & \cdots \\ \frac{\partial p_2}{\partial x_1} & \cdots \\ \vdots & \vdots \end{bmatrix}_{c \times c} = \begin{bmatrix} \sigma(x_1)(1 - \sigma(x_1)) & \cdots \\ \vdots & \vdots \end{bmatrix}$$

$$\nabla_{W_1} \mathcal{L} = Ax^\top =$$

$$\nabla_{W_1} \mathcal{L} = ((h_1 + \sigma(\hat{z}_1)(1 - h_1)) \odot (W^\top (\text{softmax}(\hat{z}_2) - e_k)))x^\top$$

$$\nabla_{b_1} \mathcal{L} = A =$$

$$\nabla_{b_1} \mathcal{L} = (h_1 + \sigma(\hat{z}_1)(1 - h_1)) \odot (W^\top (\text{softmax}(\hat{z}_2) - e_k))$$

4) 2 layer neural network - Jupyter notebook

5) General FC neural network.

$$\frac{\partial p_1}{\partial x_1} = \sigma(x_1)(1 - \sigma(x_1)) \quad \boxed{=}$$

similar to ∇_{W_2}
when it's j

$$\frac{\partial p_2}{\partial x_1} = -\sigma(x_1)\sigma(x_2)$$

$$\frac{\partial a}{\partial \hat{z}_2} \cdot \frac{\partial \mathcal{L}}{\partial a} = \begin{bmatrix} \sigma(1 - \sigma) \\ \vdots \\ \sigma(p_k) \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ 1 \end{bmatrix} \rightarrow \text{mult by the } k^{\text{th}} \text{ column}$$

$$K^{\text{th}} \text{ col} : \begin{bmatrix} -\sigma(x_1)\sigma(x_K) \\ \vdots \\ \sigma(x_K)(1 - \sigma(x_K)) \end{bmatrix}$$

$$= \begin{bmatrix} -\sigma(x_1)\sigma(x_K) \\ \vdots \\ \sigma(x_K)(1 - \sigma(x_K)) \end{bmatrix} (-1/\sigma_k)$$

$$= \begin{bmatrix} \sigma(x_1) \\ \vdots \\ \sigma(x_K) - 1 \\ \vdots \\ \sigma(x_c) \end{bmatrix} = \text{softmax}(\hat{z}_2) - \begin{bmatrix} 0 \\ \vdots \\ 1 \end{bmatrix}$$

e_k is the notation for
col vec w/ 1 @ k^{th} row,
0 otherwise

$$\frac{\partial \mathcal{L}}{\partial \hat{z}_2} = \text{softmax}(\hat{z}_2) - e_k$$

This is the 2-layer neural network notebook for ECE C147/C247 Homework #3

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

Please print out the notebook entirely when completed.

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

```
In [ ]: import random
import numpy as np
from utils.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))))
```

Toy example

Before loading CIFAR-10, there will be a toy example to test your implementation of the forward and backward pass. Make sure to read the description of TwoLayerNet class in `neural_net.py` file , understand the architecture and initializations

```
In [ ]: from nnndl.neural_net import TwoLayerNet

In [ ]: # 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 [ ]: ## Implement the forward pass of the neural network.
## See the loss() method in TwoLayerNet class for the same

# 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)))

# mine
# A D G
# J M B
# E H K
# N C F
# I L O

# correct
# A B C
# D E F
# G H I
# J K L
# M N O
```

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:
3.381231233889892e-08

Forward pass loss

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

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

```
Loss: 1.071696123862817
Difference between your loss and correct loss:
0.0
```

Backward pass

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

```
In [ ]: from utils.gradient_check import eval_numerical_gradient

# Use numeric gradient checking to check your implementation of the backward 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, 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(param_grad_num,
```

```
b2 max relative error: 1.248270530283678e-09
W2 max relative error: 2.9632227682005116e-10
b1 max relative error: 3.172680092703762e-09
W1 max relative error: 1.2832823337649917e-09
```

Training the network

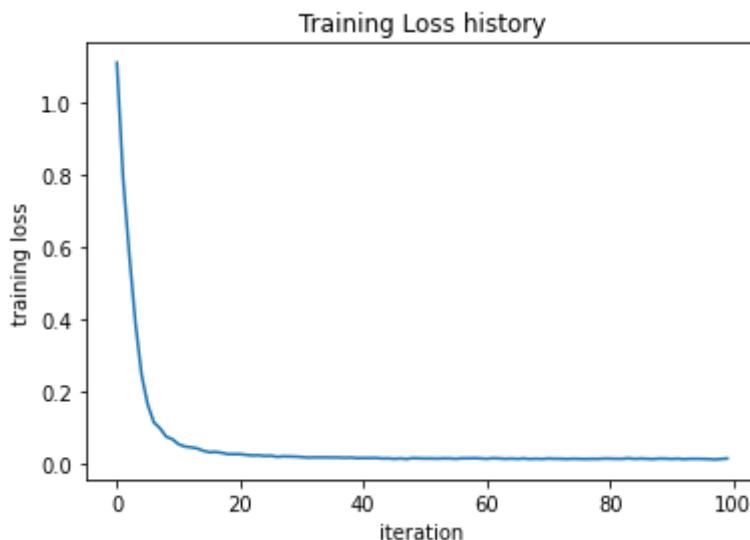
Implement `neural_net.train()` to train the network via stochastic gradient descent, much like the softmax and SVM.

```
In [ ]: 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: 0.014497864587765886
```



Classify CIFAR-10

Do classification on the CIFAR-10 dataset.

```
In [ ]: from utils.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.
    """
    # Load the raw CIFAR-10 data
    cifar10_dir = '/Users/mcapetz/Downloads/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 [ ]:
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 2.302769713193006
iteration 100 / 1000: loss 2.3024454719405805
iteration 200 / 1000: loss 2.2973223064347863
iteration 300 / 1000: loss 2.2714527896770402
iteration 400 / 1000: loss 2.1676144295711666
iteration 500 / 1000: loss 2.1187039249761
iteration 600 / 1000: loss 2.0425279262316556
iteration 700 / 1000: loss 2.0410322587782623
iteration 800 / 1000: loss 1.9946566892338784
iteration 900 / 1000: loss 1.896382745210702
Validation accuracy: 0.283
```

Questions:

The training accuracy isn't great.

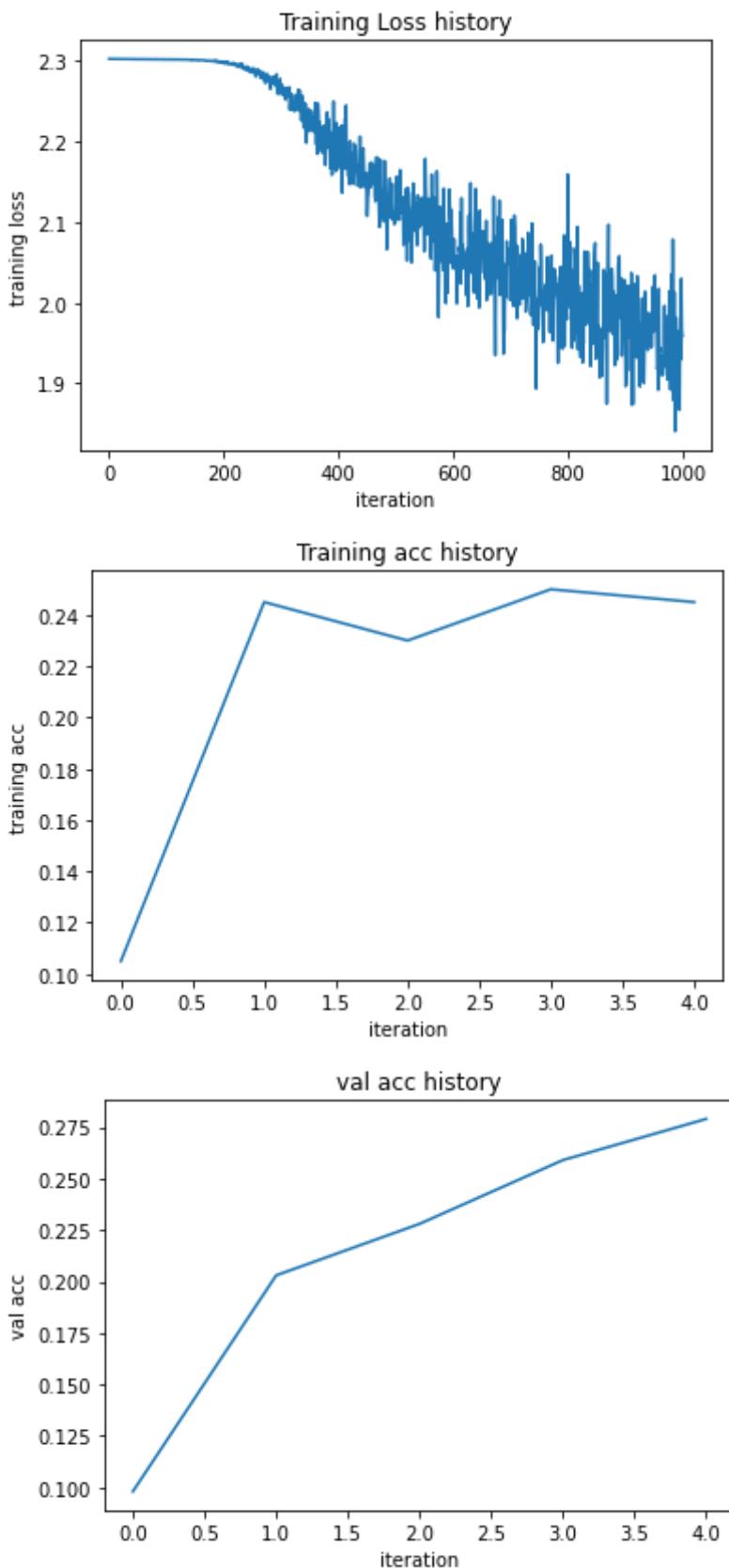
- (1) What are some of the reasons why this is the case? Take the following cell to do some analyses and then report your answers in the cell following the one below.

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

```
In [ ]: stats['train_acc_history']  
Out[ ]: [0.105, 0.245, 0.23, 0.25, 0.245]
```

```
In [ ]:  
# ======  
# 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  
  
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()  
  
# plot the loss history  
plt.plot(stats['train_acc_history'])  
plt.xlabel('iteration')  
plt.ylabel('training acc')  
plt.title('Training acc history')  
plt.show()  
  
# plot the loss history  
plt.plot(stats['val_acc_history'])  
plt.xlabel('iteration')  
plt.ylabel('val acc')  
plt.title('val acc history')  
plt.show()  
  
# ======  
# END YOUR CODE HERE  
# ======
```

```
Final training loss: 1.9585869601827879
```



Answers:

- (1) The hyperparameters including learning rate, iteration count, batch size, etc were not tuned. In particular, the iteration count was not as high as it could have been.

(2) I fixed these problems by experimenting with different numbers and plotting the different validation accuracies to see at what point they peaked. This led me to increase the number of iterations and also alter the learning rate and regularization rate.

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 [ ]: best_net = TwoLayerNet(input_size, hidden_size, num_classes) # store the best net

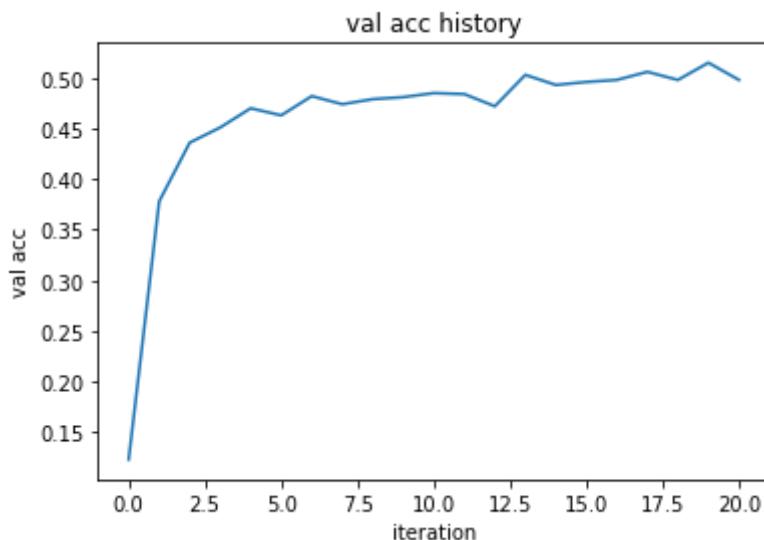
# ===== #
# YOUR CODE HERE:
#   Optimize over your hyperparameters to arrive at the best neural
#   network. You should be able to get over 50% validation accuracy.
#   For this part of the notebook, we will give credit based on the
#   accuracy you get. Your score on this question will be multiplied 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) !
# ===== #
# Train the network
stats = best_net.train(X_train, y_train, X_val, y_val,
                       num_iters=5000, batch_size=200,
                       learning_rate=1e-3, learning_rate_decay=0.95,
                       reg=0.5, verbose=True)

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

# ===== #
# END YOUR CODE HERE
# ===== #
val_acc = (best_net.predict(X_val) == y_val).mean()
print('Validation accuracy: ', val_acc)

# plot the loss history
plt.plot(stats['val_acc_history'])
plt.xlabel('iteration')
plt.ylabel('val acc')
plt.title('val acc history')
plt.show()
```

```
iteration 0 / 5000: loss 2.3029597225435428
iteration 100 / 5000: loss 1.9325933178332602
iteration 200 / 5000: loss 1.7716759550027896
iteration 300 / 5000: loss 1.7248081531071433
iteration 400 / 5000: loss 1.8443332386092726
iteration 500 / 5000: loss 1.5302909563702052
iteration 600 / 5000: loss 1.6526188579465462
iteration 700 / 5000: loss 1.5906508141462397
iteration 800 / 5000: loss 1.4855353596856218
iteration 900 / 5000: loss 1.5159556839929458
iteration 1000 / 5000: loss 1.488270908981075
iteration 1100 / 5000: loss 1.4475743724154462
iteration 1200 / 5000: loss 1.4828636891972566
iteration 1300 / 5000: loss 1.4527629724343858
iteration 1400 / 5000: loss 1.4842048768754978
iteration 1500 / 5000: loss 1.5642384638909517
iteration 1600 / 5000: loss 1.5077882504893556
iteration 1700 / 5000: loss 1.5485629034664101
iteration 1800 / 5000: loss 1.4868146985257522
iteration 1900 / 5000: loss 1.4441929644630527
iteration 2000 / 5000: loss 1.5144998528901483
iteration 2100 / 5000: loss 1.393662849428555
iteration 2200 / 5000: loss 1.5304174989897708
iteration 2300 / 5000: loss 1.445368445922107
iteration 2400 / 5000: loss 1.455080286765323
iteration 2500 / 5000: loss 1.4549381905657
iteration 2600 / 5000: loss 1.2709510874375622
iteration 2700 / 5000: loss 1.5267038798950487
iteration 2800 / 5000: loss 1.4386154710730088
iteration 2900 / 5000: loss 1.4934165464617921
iteration 3000 / 5000: loss 1.3580106333238915
iteration 3100 / 5000: loss 1.3466971155558127
iteration 3200 / 5000: loss 1.3696032551646022
iteration 3300 / 5000: loss 1.3616061593985698
iteration 3400 / 5000: loss 1.4596148696103697
iteration 3500 / 5000: loss 1.3419087222617698
iteration 3600 / 5000: loss 1.5033050723689223
iteration 3700 / 5000: loss 1.3034197628298878
iteration 3800 / 5000: loss 1.587226669889954
iteration 3900 / 5000: loss 1.3827370758432151
iteration 4000 / 5000: loss 1.5225687814985926
iteration 4100 / 5000: loss 1.4189289335169293
iteration 4200 / 5000: loss 1.353386840778016
iteration 4300 / 5000: loss 1.4960344832766483
iteration 4400 / 5000: loss 1.3407818549293629
iteration 4500 / 5000: loss 1.3933265081570054
iteration 4600 / 5000: loss 1.2662590757250827
iteration 4700 / 5000: loss 1.507514944399104
iteration 4800 / 5000: loss 1.4202320018537
iteration 4900 / 5000: loss 1.3732155789389926
Validation accuracy: 0.507
Validation accuracy: 0.507
```

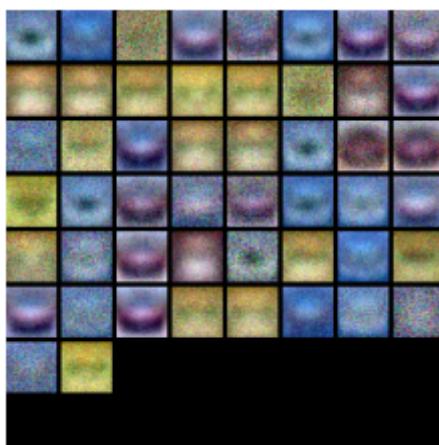


```
In [ ]: from utils.vis_utils import visualize_grid
```

```
# 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?

In []:

Answer:

- (1) The best net I arrived at is more colorful and has more contrast. There are more complex shapes and colors compared to the suboptimal net. The shapes of colors in the best net I arrived at have more distinguishable shapes and features.

Evaluate on test set

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

Test accuracy: 0.515

In []:

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).

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 []:

```
## Import and setups

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

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

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

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

In []:

```
# Load the (preprocessed) CIFAR10 data.
data = get_CIFAR10_data()
for k in data.keys():
    print('{}: {}'.format(k, data[k].shape))

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

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 [ ]: # 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: 9.769849468192957e-10

Affine layer backward pass

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

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

_, 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)))
```

```
dout (10, 5)
w (6, 5)
x shape (10, 2, 3)
dx shape (10, 2, 3) 6 5
dw shape (6, 5) 10 6
db shape (5,) 5
Testing affine_backward function:
dx error: 4.6113274414386756e-10
dw error: 5.992261138441068e-11
db error: 1.6107986410934294e-10
```

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 [ ]:
# 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.,          0.04545455,  0.13636364, ],
                       [ 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: 4.99999798022158e-08

ReLU backward pass

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

```
In [ ]:
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: 3.27560459558687e-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 []:

```
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], dx)
dw_num = eval_numerical_gradient_array(lambda w: affine_relu_forward(x, w, b)[0], dw)
db_num = eval_numerical_gradient_array(lambda b: affine_relu_forward(x, w, b)[0], db)

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: 9.272137686162983e-11
dw error: 3.880773812009996e-10
db error: 7.826730327496857e-12
```

Softmax loss

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

In []:

```
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: 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 softmax_loss:
loss: 2.302827467260296
dx error: 7.522993108189267e-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 []:

```
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, 14.57198434, 15.33,
     [12.05769098, 12.74614105, 13.43459113, 14.1230412, 14.81149128, 15.49,
     [12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.05099822, 15.66
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, grads[name])))

Testing initialization ...
Testing test-time forward pass ...
Testing training loss (no regularization)
Running numeric gradient check with reg = 0.0
W1 relative error: 1.5215703686475096e-08
W2 relative error: 3.2068321167375225e-10
b1 relative error: 8.368195737354163e-09
b2 relative error: 4.3291360264321544e-10
Running numeric gradient check with reg = 0.7
W1 relative error: 2.527915175868136e-07
W2 relative error: 2.8508510893102143e-08
b1 relative error: 1.5646801536371197e-08
b2 relative error: 7.759095355706557e-10

```

Solver

We will now use the utils Solver class to train these networks. Familiarize yourself with the API in `utils/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 []:

```

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 50%. We won't have you optimize this further
# since you did it in the previous notebook.
#
# ===== #
model = TwoLayerNet(input_dim=3*32*32, hidden_dims=200, num_classes=10)
data = get_CIFAR10_data()
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, num_train_samples=3000) # how to declare

solver.train()

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

```

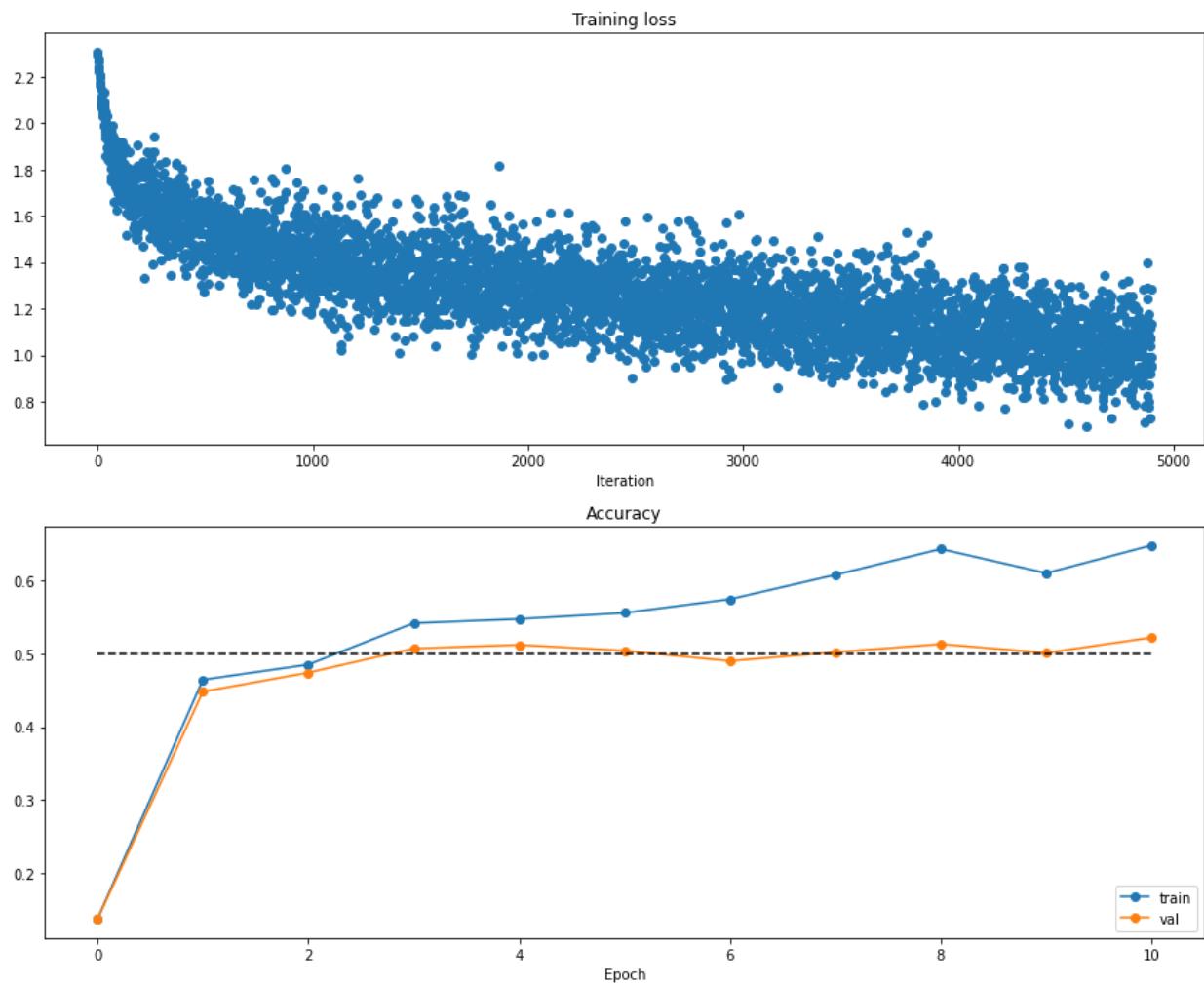
```
(Iteration 1 / 4900) loss: 2.306589
(Epoch 0 / 10) train acc: 0.137333; val_acc: 0.137000
(Iteration 101 / 4900) loss: 1.851719
(Iteration 201 / 4900) loss: 1.697929
(Iteration 301 / 4900) loss: 1.678193
(Iteration 401 / 4900) loss: 1.637067
(Epoch 1 / 10) train acc: 0.464333; val_acc: 0.448000
(Iteration 501 / 4900) loss: 1.449944
(Iteration 601 / 4900) loss: 1.458583
(Iteration 701 / 4900) loss: 1.513891
(Iteration 801 / 4900) loss: 1.288612
(Iteration 901 / 4900) loss: 1.388139
(Epoch 2 / 10) train acc: 0.485000; val_acc: 0.474000
(Iteration 1001 / 4900) loss: 1.270315
(Iteration 1101 / 4900) loss: 1.115091
(Iteration 1201 / 4900) loss: 1.211733
(Iteration 1301 / 4900) loss: 1.155765
(Iteration 1401 / 4900) loss: 1.240034
(Epoch 3 / 10) train acc: 0.541667; val_acc: 0.507000
(Iteration 1501 / 4900) loss: 1.346661
(Iteration 1601 / 4900) loss: 1.391279
(Iteration 1701 / 4900) loss: 1.287578
(Iteration 1801 / 4900) loss: 1.231270
(Iteration 1901 / 4900) loss: 1.326436
(Epoch 4 / 10) train acc: 0.547333; val_acc: 0.512000
(Iteration 2001 / 4900) loss: 1.387629
(Iteration 2101 / 4900) loss: 1.340292
(Iteration 2201 / 4900) loss: 1.231554
(Iteration 2301 / 4900) loss: 1.257639
(Iteration 2401 / 4900) loss: 1.393826
(Epoch 5 / 10) train acc: 0.555667; val_acc: 0.504000
(Iteration 2501 / 4900) loss: 1.396578
(Iteration 2601 / 4900) loss: 1.217887
(Iteration 2701 / 4900) loss: 1.190587
(Iteration 2801 / 4900) loss: 1.391394
(Iteration 2901 / 4900) loss: 1.265140
(Epoch 6 / 10) train acc: 0.574333; val_acc: 0.490000
(Iteration 3001 / 4900) loss: 1.163238
(Iteration 3101 / 4900) loss: 1.353440
(Iteration 3201 / 4900) loss: 1.012828
(Iteration 3301 / 4900) loss: 1.111804
(Iteration 3401 / 4900) loss: 1.107361
(Epoch 7 / 10) train acc: 0.607667; val_acc: 0.502000
(Iteration 3501 / 4900) loss: 1.064390
(Iteration 3601 / 4900) loss: 1.077404
(Iteration 3701 / 4900) loss: 1.041872
(Iteration 3801 / 4900) loss: 1.019779
(Iteration 3901 / 4900) loss: 1.071310
(Epoch 8 / 10) train acc: 0.643000; val_acc: 0.513000
(Iteration 4001 / 4900) loss: 1.097171
(Iteration 4101 / 4900) loss: 1.041861
(Iteration 4201 / 4900) loss: 1.144432
(Iteration 4301 / 4900) loss: 1.062828
(Iteration 4401 / 4900) loss: 1.066107
(Epoch 9 / 10) train acc: 0.610000; val_acc: 0.501000
(Iteration 4501 / 4900) loss: 1.046569
(Iteration 4601 / 4900) loss: 0.967775
(Iteration 4701 / 4900) loss: 1.135792
(Iteration 4801 / 4900) loss: 0.922348
(Epoch 10 / 10) train acc: 0.648000; val_acc: 0.522000
```

In []:

```
# 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 HW #4.

In []:

```
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-6)
        print('{} relative error: {}'.format(name, rel_error(grad_num, grads[name])))
```

```
Running check with reg = 0
Initial loss: 2.304351726737881
W1 relative error: 2.7493890788638557e-07
W2 relative error: 3.7279521977638624e-07
W3 relative error: 3.8546101120557736e-07
b1 relative error: 1.7666727565980365e-08
b2 relative error: 1.2110105191568281e-08
b3 relative error: 8.914590960924575e-11
Running check with reg = 3.14
Initial loss: 7.032554247856135
W1 relative error: 7.76263593984685e-09
W2 relative error: 7.606066279049599e-08
W3 relative error: 2.0078346520220995e-08
b1 relative error: 1.5856469842866925e-08
b2 relative error: 8.096749833757487e-09
b3 relative error: 1.7697076134496875e-10
```

In []:

```
# 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
# Your training accuracy should be 1.0 to receive full credit on this part.

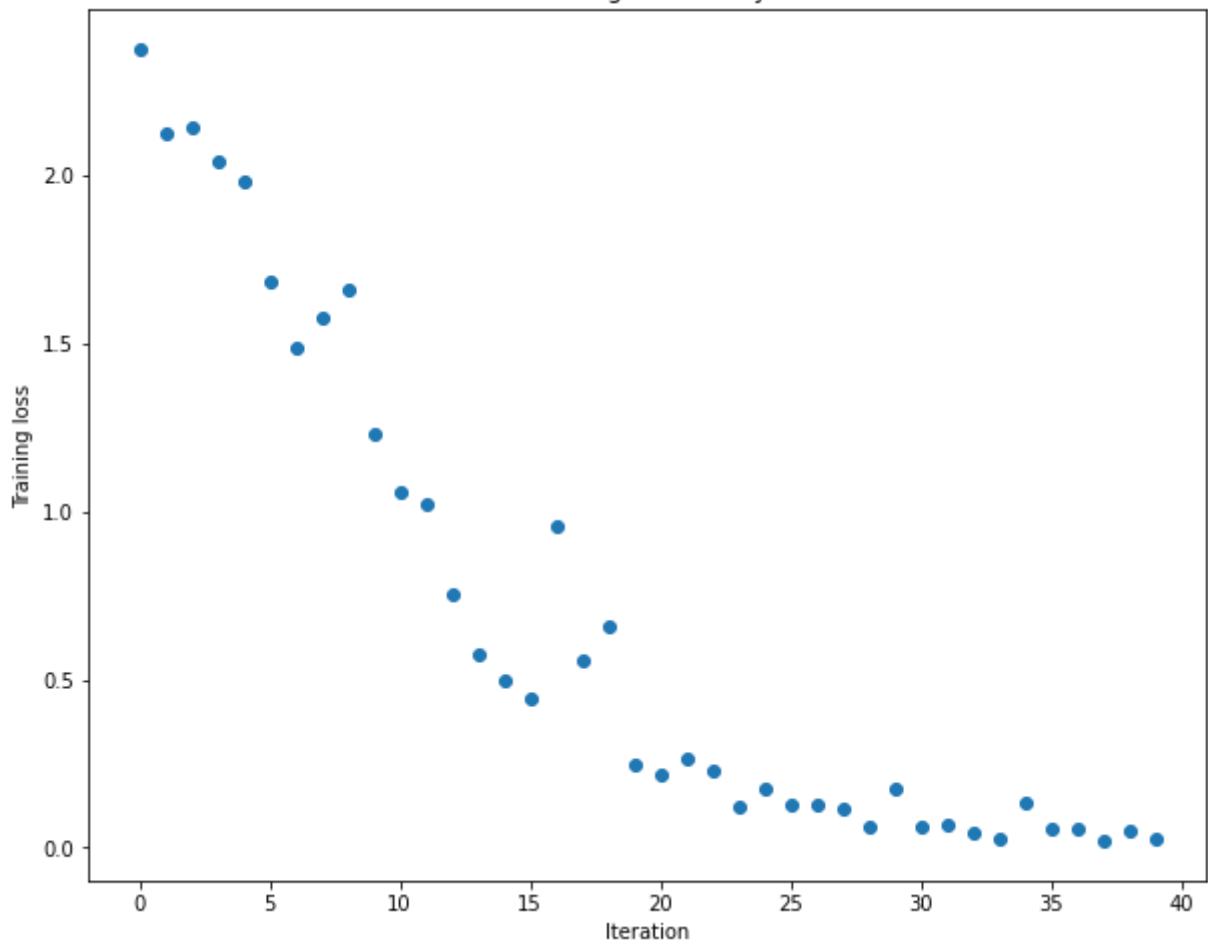
weight_scale = 1e-2
learning_rate = 1e-2

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()

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

```
(Iteration 1 / 40) loss: 2.377337
(Epoch 0 / 20) train acc: 0.220000; val_acc: 0.115000
(Epoch 1 / 20) train acc: 0.300000; val_acc: 0.115000
(Epoch 2 / 20) train acc: 0.480000; val_acc: 0.159000
(Epoch 3 / 20) train acc: 0.560000; val_acc: 0.120000
(Epoch 4 / 20) train acc: 0.540000; val_acc: 0.166000
(Epoch 5 / 20) train acc: 0.700000; val_acc: 0.160000
(Iteration 11 / 40) loss: 1.057305
(Epoch 6 / 20) train acc: 0.840000; val_acc: 0.203000
(Epoch 7 / 20) train acc: 0.800000; val_acc: 0.191000
(Epoch 8 / 20) train acc: 0.860000; val_acc: 0.193000
(Epoch 9 / 20) train acc: 0.920000; val_acc: 0.188000
(Epoch 10 / 20) train acc: 0.900000; val_acc: 0.201000
(Iteration 21 / 40) loss: 0.216215
(Epoch 11 / 20) train acc: 0.960000; val_acc: 0.197000
(Epoch 12 / 20) train acc: 0.980000; val_acc: 0.187000
(Epoch 13 / 20) train acc: 0.960000; val_acc: 0.194000
(Epoch 14 / 20) train acc: 0.940000; val_acc: 0.165000
(Epoch 15 / 20) train acc: 0.940000; val_acc: 0.176000
(Iteration 31 / 40) loss: 0.059425
(Epoch 16 / 20) train acc: 0.960000; val_acc: 0.195000
(Epoch 17 / 20) train acc: 0.980000; val_acc: 0.192000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.206000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.207000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.197000
```

Training loss history

In []:

In []:

```

import numpy as np
import matplotlib.pyplot as plt
from .layer_utils import *

class TwoLayerNet(object):
    """
    A two-layer fully-connected neural network. The net has an input
    dimension of
    D, a hidden layer dimension of H, and performs classification over C
    classes.
    We train the network with a softmax loss function and L2
    regularization on the
    weight matrices. The network uses a ReLU nonlinearity after the
    first fully
    connected layer.

    In other words, the network has the following architecture:

    input - fully connected layer - ReLU - fully connected layer -
    softmax

    The outputs of the second fully-connected layer are the scores for
    each class.
    """

    def __init__(self, input_size, hidden_size, output_size, std=1e-4):
        """
        Initialize the model. Weights are initialized to small random
        values and
        biases are initialized to zero. Weights and biases are stored in
        the
        variable self.params, which is a dictionary with the following
        keys:

        W1: First layer weights; has shape (H, D)
        b1: First layer biases; has shape (H,)
        W2: Second layer weights; has shape (C, H)
        b2: Second layer biases; has shape (C,)

        Inputs:
        - input_size: The dimension D of the input data.
        - hidden_size: The number of neurons H in the hidden layer.
        - output_size: The number of classes C.
        """
        self.params = {}
        self.params['W1'] = std * np.random.randn(hidden_size, input_size)
        self.params['b1'] = np.zeros(hidden_size)
        self.params['W2'] = std * np.random.randn(output_size,
hidden_size)

```

```

        self.params['b2'] = np.zeros(output_size)

    def loss(self, X, y=None, reg=0.0):
        """
        Compute the loss and gradients for a two layer fully connected
        neural
        network.

        Inputs:
        - X: Input data of shape (N, D). Each X[i] is a training sample.
        - y: Vector of training labels. y[i] is the label for X[i], and
        each y[i] is
            an integer in the range  $0 \leq y[i] < C$ . This parameter is
        optional; if it
            is not passed then we only return scores, and if it is passed
        then we
            instead return the loss and gradients.
        - reg: Regularization strength.

        Returns:
        If y is None, return a matrix scores of shape (N, C) where
        scores[i, c] is
            the score for class c on input X[i].
        If y is not None, instead return a tuple of:
        - loss: Loss (data loss and regularization loss) for this batch of
        training
            samples.
        - grads: Dictionary mapping parameter names to gradients of those
        parameters
            with respect to the loss function; has the same keys as
        self.params.
        """
        # Unpack variables from the params dictionary
        W1, b1 = self.params['W1'], self.params['b1']
        W2, b2 = self.params['W2'], self.params['b2']
        N, D = X.shape

        # Compute the forward pass
        scores = None

        # =====#
        # YOUR CODE HERE:
        # Calculate the output scores of the neural network. The result
        # should be (N, C). As stated in the description for this class,
        # there should not be a ReLU layer after the second FC layer.
        # The output of the second FC layer is the output scores. Do not
        # use a for loop in your implementation.

```

```

# =====#
#
f = lambda x: x * (x > 0) # relu

h1 = W1.dot(X.T) + b1.reshape(b1.shape[0], 1) # (H, D) x (D, N) +
(H, ) --> (H, N)
relu_1 = f(h1)
h2 = W2.dot(relu_1) + b2.reshape(b2.shape[0], 1) # (C, H) x (H, N)
--> (C, N)
scores = h2.T # (C, N) --> (N, C)

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



# If the targets are not given then jump out, we're done
if y is None:
    return scores

# Compute the loss
loss = None

# =====#
#
# YOUR CODE HERE:
#   Calculate the loss of the neural network. This includes the
#   softmax loss and the L2 regularization for W1 and W2. Store
the
#   total loss in the variable loss. Multiply the regularization
#   loss by 0.5 (in addition to the factor reg).
# =====#
#



# scores is num_examples by num_classes
s_loss = softmax_loss(h2.T, y)[0]
L2_reg = 0.5 * reg * (np.sum(W1 ** 2) + np.sum(W2 ** 2))
loss = s_loss + L2_reg
# =====#
#
# END YOUR CODE HERE
# =====#
#



grads = {}

```

```

# =====
#
# YOUR CODE HERE:
#   Implement the backward pass. Compute the derivatives of the
#   weights and the biases. Store the results in the grads
#   dictionary. e.g., grads['W1'] should store the gradient for
#   W1, and be of the same size as W1.
# =====
#
d1 = softmax_loss(h2.T, y)[1]

# print("w1 dim", W1.shape)
# print("w2 dim", W2.shape)
# print("relu_1 dim", relu_1.shape)
# print("d1 dim", d1.shape)

grads['b2'] = sum(d1)
d2 = d1 @ W2
# print("i got here")
grads['W2'] = d1.T @ relu_1.T + reg * W2

d3 = (h1 > 0) * d2.T
grads['b1'] = sum(d3.T)

# print("w1 dim", W1.shape)
# print("x dim", X.shape)
# print("d3 dim", d3.shape)
grads['W1'] = d3 @ X + reg * W1

# =====
#
# END YOUR CODE HERE
# =====
#
return loss, grads

def train(self, X, y, X_val, y_val,
          learning_rate=1e-3, learning_rate_decay=0.95,
          reg=1e-5, num_iters=100,
          batch_size=200, verbose=False):
    """
    Train this neural network using stochastic gradient descent.

    Inputs:
    - X: A numpy array of shape (N, D) giving training data.

```

```

    - y: A numpy array of shape (N,) giving training labels; y[i] = c
means that
    X[i] has label c, where 0 <= c < C.
    - X_val: A numpy array of shape (N_val, D) giving validation data.
    - y_val: A numpy array of shape (N_val,) giving validation labels.
    - learning_rate: Scalar giving learning rate for optimization.
    - learning_rate_decay: Scalar giving factor used to decay the
learning rate
        after each epoch.
    - reg: Scalar giving regularization strength.
    - num_iters: Number of steps to take when optimizing.
    - batch_size: Number of training examples to use per step.
    - verbose: boolean; if true print progress during optimization.
"""

num_train = X.shape[0]
iterations_per_epoch = max(num_train / batch_size, 1)

# Use SGD to optimize the parameters in self.model
loss_history = []
train_acc_history = []
val_acc_history = []

for it in np.arange(num_iters):
    X_batch = None
    y_batch = None

    #
===== #
    # YOUR CODE HERE:
    #   Create a minibatch by sampling batch_size samples randomly.
    #
===== #
    b_samples = []
    b_labels = []

    for i in range(batch_size):
        index = np.random.choice(num_train)
        b_samples.append(X[index])
        b_labels.append(y[index])

    X_batch = np.array(b_samples)
    y_batch = np.array(b_labels)

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

    # Compute loss and gradients using the current minibatch

```

```

        loss, grads = self.loss(X_batch, y=y_batch, reg=reg)
        loss_history.append(loss)

        #
===== #
# YOUR CODE HERE:
#   Perform a gradient descent step using the minibatch to
update
#   all parameters (i.e., W1, W2, b1, and b2).
#
===== #

# Unpack variables from the params dictionary
W1, b1 = self.params['W1'], self.params['b1']
W2, b2 = self.params['W2'], self.params['b2']

# update
self.params['W2'] = W2 - grads['W2'] * learning_rate
self.params['W1'] = W1 - grads['W1'] * learning_rate
self.params['b2'] = b2 - grads['b2'] * learning_rate
self.params['b1'] = b1 - grads['b1'] * learning_rate

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

if verbose and it % 100 == 0:
    print('iteration {} / {}: loss {}'.format(it, num_iters,
loss))

    # Every epoch, check train and val accuracy and decay learning
rate.
    if it % iterations_per_epoch == 0:
        # Check accuracy
        train_acc = (self.predict(X_batch) == y_batch).mean()
        val_acc = (self.predict(X_val) == y_val).mean()
        train_acc_history.append(train_acc)
        val_acc_history.append(val_acc)

        # Decay learning rate
        learning_rate *= learning_rate_decay

return {
    'loss_history': loss_history,
    'train_acc_history': train_acc_history,
    'val_acc_history': val_acc_history,
}

```

```

def predict(self, X):
    """
    Use the trained weights of this two-layer network to predict
    labels for
    data points. For each data point we predict scores for each of the
    C
    classes, and assign each data point to the class with the highest
    score.

    Inputs:
    - X: A numpy array of shape (N, D) giving N D-dimensional data
    points to
        classify.

    Returns:
    - y_pred: A numpy array of shape (N,) giving predicted labels for
    each of
        the elements of X. For all i, y_pred[i] = c means that X[i] is
    predicted
        to have class c, where 0 <= c < C.
    """
    y_pred = None

    # =====#
    # YOUR CODE HERE:
    # Predict the class given the input data.
    # =====#
    # p = (self.W @ X.T).T

    y_pred = np.zeros(X.shape[0])

    # Unpack variables from the params dictionary
    W1, b1 = self.params['W1'], self.params['b1']
    W2, b2 = self.params['W2'], self.params['b2']

    f = lambda x: x * (x > 0) # relu

    h1 = W1.dot(X.T) + b1.reshape(b1.shape[0], 1) # (H, D) x (D, N) +
    (H, ) --> (H, N)
    relu_1 = f(h1)
    h2 = W2.dot(relu_1) + b2.reshape(b2.shape[0], 1) # (C, H) x (H, N)
    --> (C, N)
    p = h2.T # (C, N) --> (N, C)

    for i in range(y_pred.shape[0]):

```

```
y_pred[i] = np.argmax(p[i])

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

return y_pred
```

```

import numpy as np
import pdb

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

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

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

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

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

# print("w", w.shape)
# print("xr", x.shape)
# print("b", b.shape)
# print("x shape", x.shape)
xr = x.reshape(x.shape[0], -1)
# print("x reshape", xr.shape)

out = xr.dot(w) + b
# print("wxr", xr.dot(w).shape) # check this shape to see if h2
# shape is even correct???


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

```

```

cache = (x, w, b)
return out, cache

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

    Inputs:
    - dout: Upstream derivative, of shape (N, M)
    - cache: Tuple of:
        - x: Input data, of shape (N, d_1, ..., d_k)
        - w: Weights, of shape (D, M)

    Returns a tuple of:
    - dx: Gradient with respect to x, of shape (N, d1, ..., d_k)
    - dw: Gradient with respect to w, of shape (D, M)
    - db: Gradient with respect to b, of shape (M, )
    """
    x, w, b = cache
    dx, dw, db = None, None, None

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

    # dout is N x M
    # dx should be N x d1 x ... x dk; it relates to dout through
    multiplication with w, which is D x M
    # dw should be D x M; it relates to dout through multiplication with
    x, which is N x D after reshaping
    # db should be M; it is just the sum over dout examples

    # dout: 5 x 3
    # w: 3 x 10
    # x: 10 x 5 which is M x N
    # print("x shape", x.shape)

    # n = 5
    # m =
    # d = 3
    # m = THERE IS SMTH WRONG HERE!!! WITH THE NUMBERS
    # print("dout", dout.shape)
    # print("w", w.shape)
    dx = np.dot(dout, w.T) # (N, M) x (M, D) = (N, D) should be D,M
    # print("dout shape", dout.shape)
    # print("w.T shape", w.T.shape)
    # print("x shape", x.shape)

```

```

# dx = dx.reshape(x.shape)

# dw = np.dot(x.reshape(x.shape[0], -1).T, dout) # (N, M) x (1 x N)
# should be (D, M) original
xr = x.reshape(x.shape[0], -1)
dw = np.dot(xr.T, dout) #not sure if this works
db = np.sum(dout, axis=0) # M i think this one is fine
# print("x shape", x.shape)
# print("dx shape", dx.shape, w.shape[0], dout.shape[1])
# print("dw shape", dw.shape, dout.shape[0],w.shape[0])
# print("db shape", db.shape, dout.shape[1])

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

return dx, dw, db

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

    Input:
    - x: Inputs, of any shape

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

    f = lambda x: x * (x > 0)
    out = f(x)
    # print("out", out.shape)
    # print("x", x.shape)
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    cache = x
    return out, cache

def relu_backward(dout, cache):
    """

```

Computes the backward pass for a layer of rectified linear units (ReLUs).

Input:

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

Returns:

- dx: Gradient with respect to x
- """

x = cache

```
# ===== #
# YOUR CODE HERE:
# Implement the ReLU backward pass
# ===== #

# ReLU directs linearly to those > 0
# xr = x.reshape(x.shape[0], -1)
dx = dout * (x > 0) # hadamard product is element wise operation

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

return dx
```

def softmax_loss(x, y):

"""

Computes the loss and gradient for softmax classification.

Inputs:

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

Returns a tuple of:

- loss: Scalar giving the loss
- dx: Gradient of the loss with respect to x
- """

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

```
dx[np.arange(N), y] -= 1  
dx /= N  
return loss, dx
```

```

import numpy as np

from .layers import *
from .layer_utils import *


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

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

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

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

    def __init__(self, input_dim=3*32*32, hidden_dims=100,
                 num_classes=10,
                 dropout=0, weight_scale=1e-3, reg=0.0):
        """
        Initialize a new network.

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

        # =====#
        #
        # YOUR CODE HERE:
        #   Initialize W1, W2, b1, and b2. Store these as
        self.params['W1'],

```

```

#      self.params['W2'], self.params['b1'] and self.params['b2'].
The
#      biases are initialized to zero and the weights are initialized
#      so that each parameter has mean 0 and standard deviation
weight_scale.
#      The dimensions of W1 should be (input_dim, hidden_dim) and the
#      dimensions of W2 should be (hidden_dims, num_classes)
# =====
#
#      self.params = {}
#      self.params['W1'] = weight_scale * np.random.randn(input_dim,
hidden_dims)
#      self.params['b1'] = np.zeros(hidden_dims)
#      self.params['W2'] = weight_scale * np.random.randn(hidden_dims,
num_classes)
#      self.params['b2'] = np.zeros(num_classes)

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

    Inputs:
    - X: Array of input data of shape (N, d_1, ..., d_k)
    - y: Array of labels, of shape (N,). y[i] gives the label for
X[i].
    Returns:
    If y is None, then run a test-time forward pass of the model and
return:
        - scores: Array of shape (N, C) giving classification scores,
where
            scores[i, c] is the classification score for X[i] and class c.

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

```

```

# =====#
#
# YOUR CODE HERE:
# Implement the forward pass of the two-layer neural network.
Store
# the class scores as the variable 'scores'. Be sure to use the
layers
# you prior implemented.
# =====#
#
# Unpack variables from the params dictionary
W1, b1 = self.params['W1'], self.params['b1']
W2, b2 = self.params['W2'], self.params['b2']

h1 = affine_forward(X, W1, b1)
relu_1 = relu_forward(h1[0])
h2 = affine_forward(relu_1[0], W2, b2)
scores = h2[0]
# =====#
#
# END YOUR CODE HERE
# =====#
#
# If y is None then we are in test mode so just return scores
if y is None:
    return scores

loss, grads = 0, {}
# =====#
#
# YOUR CODE HERE:
# Implement the backward pass of the two-layer neural net.
Store
# the loss as the variable 'loss' and store the gradients in the
# 'grads' dictionary. For the grads dictionary, grads['W1']
holds
# the gradient for W1, grads['b1'] holds the gradient for b1,
etc.
# i.e., grads[k] holds the gradient for self.params[k].
#
# Add L2 regularization, where there is an added cost
0.5*self.reg*W^2
# for each W. Be sure to include the 0.5 multiplying factor to
# match our implementation.
#
# And be sure to use the layers you prior implemented.
# =====#

```

```

#
# grads['W2'], grads['b2'], grads['W1'], grads['b1'] = None, None,
None, None
# W1, b1 = self.params['W1'], self.params['b1']
# W2, b2 = self.params['W2'], self.params['b2']

softmax = softmax_loss(h2[0], y)
loss = softmax[0] + 0.5 * self.reg * (np.sum(W1 ** 2) + np.sum(W2
** 2))

# h1_grad, h2_grad, relu_grad = None, None, None

# since w2 has problems need to check if h2[0] is correct shape
but i think it is?
# need to draw out and check dims
# print("h2[0] shape", h2[0].shape)
# print("w2 shape", W2.shape)
# print("b2 shape", b2.shape)
# print("W1 shape", W1.shape)
# print("b1 shape", b1.shape)
# print("relu shape", relu_1[0].shape)
# print("x shape", X.shape)
(h2_grad, grads['W2'], grads['b2']) = affine_backward(softmax[1],
(relu_1[0], W2, b2))

relu_grad = relu_backward(h2_grad, h1[0])

# print("W2 shape", W2.shape)
# print("grads['W2'] shape", grads['W2'].shape)

(h1_grad, grads['W1'], grads['b1']) = affine_backward(relu_grad,
(X, W1, b1))

# print("grads w2", grads['W2'].shape)
# print("w2 shape", W2.shape)

# add regularization
# print("w1 shape", W1.shape)
# print("w1 grad shape", grads['W1'].shape)

grads['W2'] += self.reg*W2
grads['W1'] += self.reg*W1

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

```

```

#  

    return loss, grads  

class FullyConnectedNet(object):  

    """  

        A fully-connected neural network with an arbitrary number of hidden  

        layers,  

        ReLU nonlinearities, and a softmax loss function. This will also  

        implement  

        dropout and batch normalization as options. For a network with L  

        layers,  

        the architecture will be  

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

        softmax  

        where batch normalization and dropout are optional, and the {...}  

        block is  

        repeated L - 1 times.  

        Similar to the TwoLayerNet above, learnable parameters are stored in  

        the  

        self.params dictionary and will be learned using the Solver class.  

    """  

    def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,  

                 dropout=0, use_batchnorm=False, reg=0.0,  

                 weight_scale=1e-2, dtype=np.float32, seed=None):  

        """  

            Initialize a new FullyConnectedNet.  

            Inputs:  

            - hidden_dims: A list of integers giving the size of each hidden  

            layer.  

            - input_dim: An integer giving the size of the input.  

            - num_classes: An integer giving the number of classes to  

            classify.  

            - dropout: Scalar between 0 and 1 giving dropout strength. If  

            dropout=0 then  

                the network should not use dropout at all.  

            - use_batchnorm: Whether or not the network should use batch  

            normalization.  

            - reg: Scalar giving L2 regularization strength.  

            - weight_scale: Scalar giving the standard deviation for random  

            initialization of the weights.  

            - dtype: A numpy datatype object; all computations will be  

            performed using  

                this datatype. float32 is faster but less accurate, so you

```

```

should use
    float64 for numeric gradient checking.
    - seed: If not None, then pass this random seed to the dropout
layers. This
    will make the dropout layers deteriministic so we can gradient
check the
    model.
=====
self.use_batchnorm = use_batchnorm
self.use_dropout = dropout > 0
self.reg = reg
self.num_layers = 1 + len(hidden_dims)
self.dtype = dtype
self.params = {}

# =====
#
# YOUR CODE HERE:
#   Initialize all parameters of the network in the self.params
dictionary.
#   The weights and biases of layer 1 are W1 and b1; and in
general the
#   weights and biases of layer i are Wi and bi. The
#   biases are initialized to zero and the weights are initialized
#   so that each parameter has mean 0 and standard deviation
weight_scale.
# =====
#
        self.params['W1'] = weight_scale * np.random.randn(input_dim,
hidden_dims[0])
        self.params['b1'] = np.zeros(hidden_dims[0])

for i in range(1, self.num_layers - 1):
    w = str("W" + str(i + 1))
    b = str("b" + str(i + 1))

    self.params[w] = weight_scale *
np.random.randn(hidden_dims[i-1], hidden_dims[i])
    self.params[b] = np.zeros(hidden_dims[i])

    w = str("W" + str(self.num_layers))
    b = str("b" + str(self.num_layers))
    self.params[w] = weight_scale * np.random.randn(hidden_dims[-1],
num_classes)
    self.params[b] = np.zeros(num_classes)

# for key in self.params:
#     print(key, self.params[key].shape)

```

```

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

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

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

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

    Input / output: Same as TwoLayerNet above.
    """
    X = X.astype(self.dtype)
    mode = 'test' if y is None else 'train'

    # Set train/test mode for batchnorm params and dropout param since
they
    # behave differently during training and testing.
    if self.dropout_param is not None:

```

```

        self.dropout_param['mode'] = mode
    if self.use_batchnorm:
        for bn_param in self.bn_params:
            bn_param[mode] = mode

    scores = None

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

# forward prop
Hs = [X]
Zs = [X]
Ws = []
bs = []

W = self.params['W1']
b = self.params['b1']
aff_fwd = affine_forward(X, W, b)
Z = aff_fwd[0]

for i in range(1, self.num_layers):
    relu_h = relu_forward(Z)
    H = relu_h[0]
    Hs.append(H)
    Ws.append(W)
    Zs.append(Z)
    bs.append(b)

    H = Hs[-1]
    W = self.params[str('W' + str(i+1))]
    b = self.params[str('b' + str(i+1))]

    aff_fwd = affine_forward(H, W, b)
    Z = aff_fwd[0]
    scores = Z

    Zs.append(Z)
    Ws.append(W)
    bs.append(b)
    # print("done with fwd pass")

# =====#
#

```

```

# END YOUR CODE HERE
# =====
#
# If test mode return early
if mode == 'test':
    return scores

loss, grads = 0.0, {}
# =====
#
# YOUR CODE HERE:
# Implement the backwards pass of the FC net and store the
gradients
# in the grads dict, so that grads[k] is the gradient of
self.params[k]
# Be sure your L2 regularization includes a 0.5 factor.
# =====
#
softmax = softmax_loss(Zs[-1], y)

# regularization
agg_sum = 0
for W in Ws:
    agg_sum += np.sum(W ** 2)

loss = softmax[0] + 0.5 * self.reg * agg_sum

# backprop
loss_grads = [softmax[1]]

for i in range(self.num_layers, 1, -1):
    loss_grad = loss_grads[-1]
    (h_grad, grads[str('W' + str(i))], grads[str('b' + str(i))]) =
affine_backward(loss_grad, (Hs[i-1], Ws[i-1], bs[i-1]))
    relu_grad = relu_backward(h_grad, Zs[i-1])
    loss_grads.append(relu_grad)

    loss_grad = loss_grads[-1]
    (h_grad, grads['W1'], grads['b1']) = affine_backward(loss_grad,
(X, Ws[0], bs[0]))

# regularization
for i in range(self.num_layers):
    grads[str('W' + str(i+1))] += self.reg * self.params[str('W' +
str(i+1))]

# for key in grads.keys():

```

```

#     print(key, grads[key].shape)

# loss_grad = loss_grads[-1]
# (h_grad, grads['W3'], grads['b3']) = affine_backward(loss_grad,
(Zs[-1], Ws[-1], bs[-1]))
# print("loss grad shape", loss_grad[0].shape)
# print("h grad shape", h_grad.shape)

# # print("w grad size", w_grad.shape)

# for i in range(self.num_layers -1, 0, -1):
#     print(i)
#     loss_grad = loss_grads[-1]
#     print("h_grad size", h_grad.shape)
#     print("Hs[i] shape", Hs[i-1].shape)
#     relu_grad = relu_backward(h_grad, Hs[i-1])
#     print("relu grad shape", relu_grad.shape)
#     print("h shape", Hs[i-1].shape)
#     print("w shape", Ws[i-1].shape)
#     print("b shape", bs[i-1].shape)
#     (h_grad, w_grad, b_grad) = affine_backward(relu_grad,
(Hs[i-1], Ws[i-1], bs[i-1]))
#     # print("new h shape", h_grad.shape)
#     loss_grads.append(h_grad)
#     print(str('W' + str(i)))
#     grads[str('W' + str(i))] = w_grad
#     grads[str('b' + str(i))] = b_grad

# regularization

# print("i", i)
# print(str('W'+str(i+1)))
# # need to figure this part out
# print(grads['W1'])
# #self.reg*self.params['W1']
# print("i got here 2")

# j = 0
# for W in Ws:
#     ^
#     j = j+1

```

```
# (h2_grad, grads['W2'], grads['b2']) =
affine_backward(softmax[1], (relu_1[0], W2, b2))

# relu_grad = relu_backward(h2_grad, h1[0])

# # print("W2 shape", W2.shape)
# # print("grads['W2'] shape", grads['W2'].shape)

# (h1_grad, grads['W1'], grads['b1']) = affine_backward(relu_grad,
(X, W1, b1))

# # print("grads w2", grads['W2'].shape)
# # print("w2 shape", W2.shape)

# # add regularization
# # print("w1 shape", W1.shape)
# # print("w1 grad shape", grads['W1'].shape)

# grads['W2'] += self.reg*W2
# grads['W1'] += self.reg*W1

# =====#
# END YOUR CODE HERE
# =====#
# return loss, grads
```