

# ECE C247 HW 3 Solution

Ashish Kumar Singh (UID:105479019)

February 3, 2022

**Problem 1.** Backpropagation for autoencoders

**Sub-Problem 1a.** preserve info?

**Solution 1a.** When we minimize loss, we are minimizing the difference between  $W^T W x$  and  $x$ , which means  $Wx$  preserves information about  $x$  as we are able to recover it by multiplying it with  $W^T$ .

**Sub-Problem 1b.** Computational graph for  $\mathcal{L}$

**Solution 1b.** Fig. 1 shows the computational graph for  $\mathcal{L}$

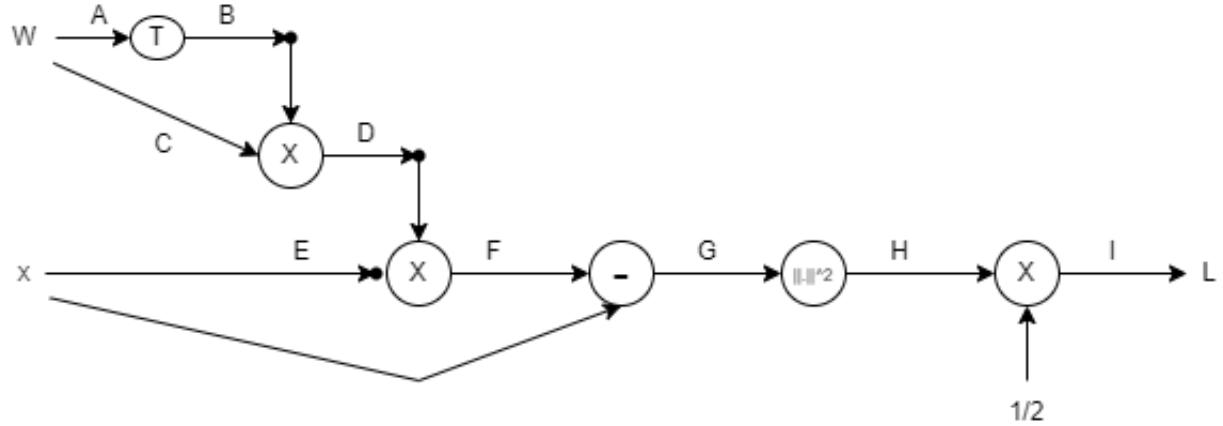


Figure 1: Computational graph

**Sub-Problem 1c.** two paths for  $W$

**Solution 1c.** We can calculate  $\nabla_W \mathcal{L}$  by adding derivative along both path, Consider following two paths,  $W -> a -> \mathcal{L}$  and  $W -> b -> \mathcal{L}$  then by law of total derivatives,

$$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial a} \frac{\partial a}{\partial W} + \frac{\partial \mathcal{L}}{\partial b} \frac{\partial b}{\partial W}$$

**Sub-Problem 1d.** Find  $\nabla_W \mathcal{L}$

**Solution 1d.** Fig. 1 shows the computational graph, in which  $A = W$ ,  $B = W^T$ ,  $C = W$ ,  $D = W^T W$ ,  $E = x$ ,  $F = W^T W x$ ,  $G = W^T W x - x$ ,  $H = \|W^T W x - x\|^2$  and  $I = 0.5 * \|W^T W x - x\|^2 = \mathcal{L}$

$$\frac{\partial \mathcal{L}}{\partial I} = 1$$

$$\frac{\partial \mathcal{L}}{\partial H} = 1/2$$

$$\frac{\partial \mathcal{L}}{\partial G} = 2G \frac{\partial \mathcal{L}}{\partial H} = G$$

$$\frac{\partial \mathcal{L}}{\partial F} = G$$

$$\frac{\partial \mathcal{L}}{\partial D} = \frac{\partial \mathcal{L}}{\partial F} E^T = G x^T$$

$$\frac{\partial \mathcal{L}}{\partial B} = \frac{\partial \mathcal{L}}{\partial D} C^T = G x^T W^T$$

$$\frac{\partial \mathcal{L}}{\partial C} = B^T \frac{\partial \mathcal{L}}{\partial D} = W G x^T$$

$$\frac{\partial \mathcal{L}}{\partial A} = (\frac{\partial \mathcal{L}}{\partial B})^T = (G x^T W^T)^T = W x G^T$$

$$\nabla_W \mathcal{L} = \frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial A} + \frac{\partial \mathcal{L}}{\partial C} = W x G^T + W G x^T$$

where  $G = W^T W x - x$

**Problem 2.** Gaussian-process latent variable model

**Sub-Problem 2a.** Computational graph for  $\mathcal{L}_1$

**Solution 2a.** Fig. 2 shows the computational graph for  $\mathcal{L}_1$

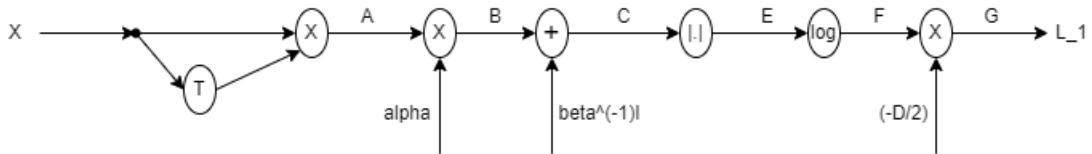


Figure 2: Computational graph

**Sub-Problem 2b.** Compute  $\frac{\partial \mathcal{L}_1}{\partial X}$

**Solution 2b.** Fig. 2 shows the computational graph for  $\mathcal{L}_1$ , where  $A = XX^T$ ,  $B = \alpha XX^T$ ,  $C = \alpha XX^T + \beta^{-1}I$ ,  $E = |\alpha XX^T + \beta^{-1}I|$ ,  $F = \log(|\alpha XX^T + \beta^{-1}I|)$  and  $G = -(D/2)\log(|\alpha XX^T + \beta^{-1}I|) = \mathcal{L}_1$

$$\frac{\partial \mathcal{L}_1}{\partial G} = 1$$

$$\frac{\partial \mathcal{L}_1}{\partial F} = (-D/2)$$

$$\frac{\partial \mathcal{L}_1}{\partial E} = (-D/2) \frac{\partial F}{\partial E} = \frac{-D}{2E}$$

$$\frac{\partial \mathcal{L}_1}{\partial C} = \frac{-D}{2E} \frac{\partial E}{\partial C} = \frac{-D|C|(C^T)^{-1}}{2E} = \frac{-D(C^T)^{-1}}{2}$$

$$\frac{\partial \mathcal{L}_1}{\partial B} = \frac{-D(C^T)^{-1}}{2}$$

$$\frac{\partial \mathcal{L}_1}{\partial A} = \frac{-\alpha D(C^T)^{-1}}{2}$$

$$\frac{\partial \mathcal{L}_1}{\partial X} = \frac{\partial \mathcal{L}_1}{\partial A} \frac{\partial A}{\partial X} = \frac{-\alpha D(C^T)^{-1}}{2} \frac{\partial XX^T}{\partial X} = -\alpha D(C^T)^{-1}X$$

Given  $K = \alpha XX^T + \beta^{-1}I = C$

$$\frac{\partial \mathcal{L}_1}{\partial X} = -\alpha D(K^T)^{-1}X$$

Since  $K^T = \alpha XX^T + \beta^{-1}I = K$

$$\frac{\partial \mathcal{L}_1}{\partial X} = -\alpha DK^{-1}X$$

**Sub-Problem 2c.** Computational graph for  $\mathcal{L}_2$

**Solution 2c.** Fig. 3 shows the computational graph for  $\mathcal{L}_2$

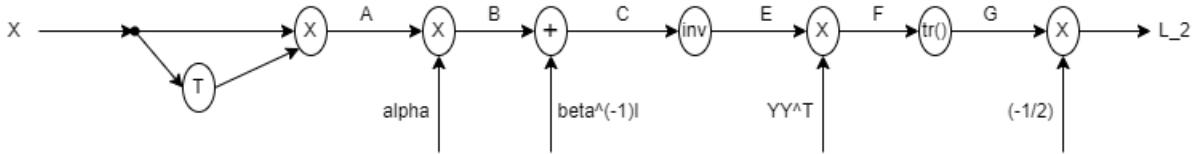


Figure 3: Computational graph

**Sub-Problem 2d.** Compute  $\frac{\partial \mathcal{L}_2}{\partial X}$

**Solution 2d.** Fig. 3 shows the computational graph for  $\mathcal{L}_2$ , where  $A = XX^T$ ,  $B = \alpha XX^T$ ,  $C = \alpha XX^T + \beta^{-1}I = K$ ,  $E = (\alpha XX^T + \beta^{-1}I)^{-1}$ ,  $F = EYY^T$  and  $G = \text{tr}(F)$

$$\begin{aligned}\frac{\partial \mathcal{L}_2}{\partial G} &= (-1/2) \\ \frac{\partial \mathcal{L}_2}{\partial F} &= (-1/2)I \\ \frac{\partial \mathcal{L}_2}{\partial E} &= (-1/2)YY^T \\ \frac{\partial \mathcal{L}_2}{\partial C} &= -C^{-T} \frac{\partial \mathcal{L}_2}{\partial E} C^{-T} = (1/2)C^{-T}YY^TC^{-T} \\ \frac{\partial \mathcal{L}_2}{\partial B} &= (1/2)C^{-T}YY^TC^{-T} \\ \frac{\partial \mathcal{L}_2}{\partial A} &= (1/2)\alpha C^{-T}YY^TC^{-T} \\ \frac{\partial \mathcal{L}_2}{\partial X} &= \alpha C^{-T}YY^TC^{-T}X\end{aligned}$$

Given  $K = \alpha XX^T + \beta^{-1}I = C$  and  $K^T = K$

$$\frac{\partial \mathcal{L}_2}{\partial X} = \alpha K^{-1}YY^TK^{-1}X$$

**Sub-Problem 2e.** Compute  $\frac{\partial \mathcal{L}}{\partial X}$

**Solution 2e.**

$$\begin{aligned}\mathcal{L} &= -c + \mathcal{L}_1 + \mathcal{L}_2 \\ \frac{\partial \mathcal{L}}{\partial X} &= \frac{\partial \mathcal{L}_1}{\partial X} + \frac{\partial \mathcal{L}_2}{\partial X} \\ \frac{\partial \mathcal{L}}{\partial X} &= -\alpha DK^{-1}X + \alpha K^{-1}YY^TK^{-1}X\end{aligned}$$

Problem 3 and 4 code is attached below.

# 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 [90]:

```
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))))
```

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 [91]:

```
from nnndl.neural_net import TwoLayerNet
```

In [92]:

```
# 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 [93]:

```
## 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:  
3.381231204052648e-08

## Forward pass loss

In [94]:

```

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 [95]:

```

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:
    #print(param_name)
    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_

```

W2 max relative error: 2.9632216125742514e-10

```
b2 max relative error: 1.2482669498248223e-09
b1 max relative error: 3.1726804786908923e-09
W1 max relative error: 1.2832845443256344e-09
```

## Training the network

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

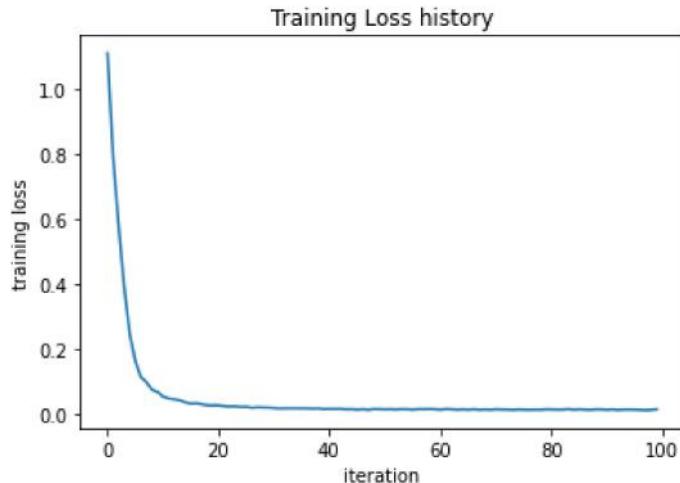
In [99]:

```
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 [101...]

```
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 = '../hw2-code/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 [130]: 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 2.3027848915701807
iteration 100 / 1000: loss 2.302335829375214
iteration 200 / 1000: loss 2.2953762275847236
iteration 300 / 1000: loss 2.254731626176731
iteration 400 / 1000: loss 2.150319523052278
iteration 500 / 1000: loss 2.219673159078906
iteration 600 / 1000: loss 2.031065845367762
iteration 700 / 1000: loss 2.0495904195319734
iteration 800 / 1000: loss 2.037775926722867
iteration 900 / 1000: loss 1.9956964919555678
Validation accuracy: 0.286

```

## 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 [107... stats['train_acc_history']
```

```
Out[107... [0.11, 0.15, 0.18, 0.22, 0.315]
```

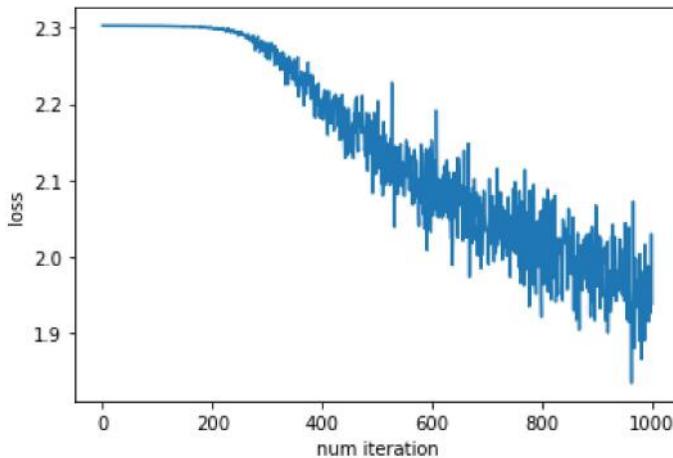
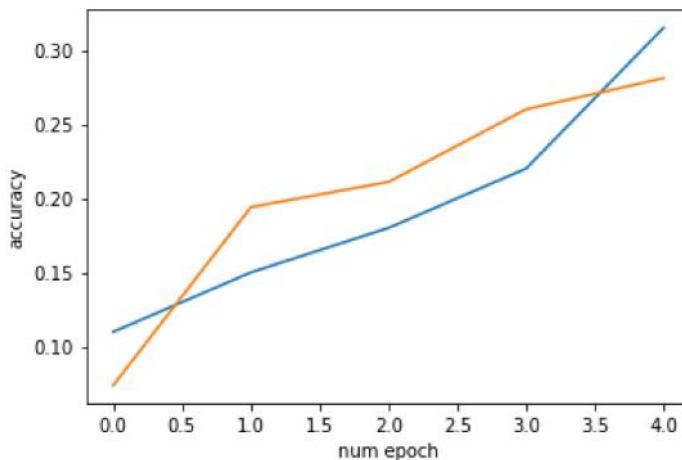
```
In [116... # ===== #
# YOUR CODE HERE:
# Do some debugging to gain some insight into why the optimization
# isn't great.
# ===== #

# Plot the loss function and train / validation accuracies

plt.plot(stats['train_acc_history'], label='train acc')
plt.plot(stats['val_acc_history'], label='val acc')
plt.xlabel('num epoch')
plt.ylabel('accuracy')
plt.show()

plt.plot(stats['loss_history'], label='train loss')
plt.xlabel('num iteration')
plt.ylabel('loss')
plt.show()

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



## Answers:

(1) We can see that the train loss is still decreasing and it is not converged yet. Also the train and val accuracy are increasing every epoch which means we need to train more or increase the learning rate.

(2) optimize hyperparameters like learning rate and training more, till the loss curve converges

## 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 [139...]

```
best_net = None # store the best model into this

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

lrs = [10**x for x in np.arange(-3.2,-2.2,0.2)]
bs = np.arange(220,250,20)
regs = np.arange(0.1,0.21,0.02)
max_iters = np.arange(1000,1600,100)
best_val = 0
print(lrs,bs, regs)

for lr in lrs:
    for batch in bs:
        for reg in regs:
            for max_iter in max_iters:
                print("lr: ",lr," bs: ",batch, " reg: ",reg, " max iter: ", max_iter, " best val")
                net = TwoLayerNet(input_size, hidden_size, num_classes)
                stats = net.train(X_train, y_train, X_val, y_val,
                                  num_iters=max_iter, batch_size=batch,
                                  learning_rate=lr, learning_rate_decay=0.95,
                                  reg=reg, verbose=False)

                val_acc = (net.predict(X_val) == y_val).mean()
                if val_acc>best_val:
                    best_val = val_acc
                    best_net = net

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

val_acc = (best_net.predict(X_val) == y_val).mean()
print('Validation accuracy: ', val_acc)
```

```
[0.000630957344480193, 0.001, 0.001584893192461114, 0.002511886431509582, 0.003981071705534978]
[220 240] [0.1 0.12 0.14 0.16 0.18 0.2 ]
lr: 0.000630957344480193 bs: 220 reg: 0.1 max iter: 1000 best val: 0
lr: 0.000630957344480193 bs: 220 reg: 0.1 max iter: 1100 best val: 0.456
lr: 0.000630957344480193 bs: 220 reg: 0.1 max iter: 1200 best val: 0.462
lr: 0.000630957344480193 bs: 220 reg: 0.1 max iter: 1300 best val: 0.462
lr: 0.000630957344480193 bs: 220 reg: 0.1 max iter: 1400 best val: 0.465
lr: 0.000630957344480193 bs: 220 reg: 0.1 max iter: 1500 best val: 0.494
lr: 0.000630957344480193 bs: 220 reg: 0.12000000000000001 max iter: 1000 best val: 0.494
lr: 0.000630957344480193 bs: 220 reg: 0.12000000000000001 max iter: 1100 best val: 0.494
lr: 0.000630957344480193 bs: 220 reg: 0.12000000000000001 max iter: 1200 best val: 0.494
```

two layer nn



## two layer nn



two layer nn

two layer nn

## two\_layer\_nn

```

lr: 0.003981071705534978 bs: 220 reg: 0.18000000000000002 max iter: 1300 best val: 0.5
08
lr: 0.003981071705534978 bs: 220 reg: 0.18000000000000002 max iter: 1400 best val: 0.5
08
lr: 0.003981071705534978 bs: 220 reg: 0.18000000000000002 max iter: 1500 best val: 0.5
08
lr: 0.003981071705534978 bs: 220 reg: 0.2 max iter: 1000 best val: 0.508
lr: 0.003981071705534978 bs: 220 reg: 0.2 max iter: 1100 best val: 0.508
lr: 0.003981071705534978 bs: 220 reg: 0.2 max iter: 1200 best val: 0.508
lr: 0.003981071705534978 bs: 220 reg: 0.2 max iter: 1300 best val: 0.508
lr: 0.003981071705534978 bs: 220 reg: 0.2 max iter: 1400 best val: 0.508
lr: 0.003981071705534978 bs: 220 reg: 0.2 max iter: 1500 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.1 max iter: 1000 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.1 max iter: 1100 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.1 max iter: 1200 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.1 max iter: 1300 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.1 max iter: 1400 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.1 max iter: 1500 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.12000000000000001 max iter: 1000 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.12000000000000001 max iter: 1100 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.12000000000000001 max iter: 1200 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.12000000000000001 max iter: 1300 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.12000000000000001 max iter: 1400 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.12000000000000001 max iter: 1500 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.14 max iter: 1000 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.14 max iter: 1100 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.14 max iter: 1200 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.14 max iter: 1300 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.14 max iter: 1400 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.14 max iter: 1500 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.16000000000000003 max iter: 1000 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.16000000000000003 max iter: 1100 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.16000000000000003 max iter: 1200 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.16000000000000003 max iter: 1300 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.16000000000000003 max iter: 1400 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.16000000000000003 max iter: 1500 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.18000000000000002 max iter: 1000 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.18000000000000002 max iter: 1100 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.18000000000000002 max iter: 1200 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.18000000000000002 max iter: 1300 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.18000000000000002 max iter: 1400 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.18000000000000002 max iter: 1500 best val: 0.5
08
lr: 0.003981071705534978 bs: 240 reg: 0.2 max iter: 1000 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.2 max iter: 1100 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.2 max iter: 1200 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.2 max iter: 1300 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.2 max iter: 1400 best val: 0.508
lr: 0.003981071705534978 bs: 240 reg: 0.2 max iter: 1500 best val: 0.508
Validation accuracy: 0.508

```

In [140...]

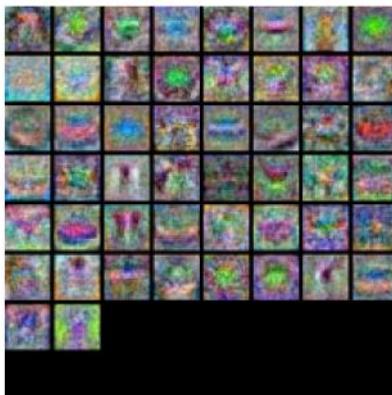
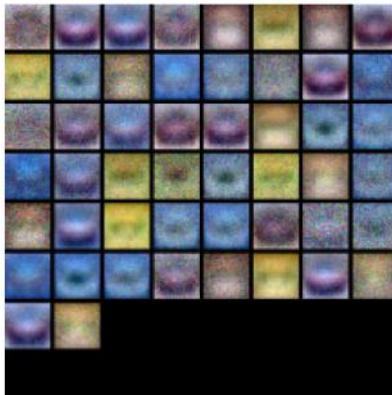
```

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?

## Answer:

- (1) Suboptimal models weight seems to be over smoothed whereas best net weight seems to have more complex features

## Evaluate on test set

```
In [141]: test_acc = (best_net.predict(X_test) == y_test).mean()
print('Test accuracy: ', test_acc)
```

Test accuracy: 0.49

In [ ]:

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4
5 class TwoLayerNet(object):
6     """
7         A two-layer fully-connected neural network. The net has an input dimension of
8             N, a hidden layer dimension of H, and performs classification over C classes.
9             We train the network with a softmax loss function and L2 regularization on the
10            weight matrices. The network uses a ReLU nonlinearity after the first fully
11            connected layer.
12
13     In other words, the network has the following architecture:
14
15     input - fully connected layer - ReLU - fully connected layer - softmax
16
17     The outputs of the second fully-connected layer are the scores for each class.
18     """
19
20     def __init__(self, input_size, hidden_size, output_size, std=1e-4):
21         """
22             Initialize the model. Weights are initialized to small random values and
23             biases are initialized to zero. Weights and biases are stored in the
24             variable self.params, which is a dictionary with the following keys:
25
26             W1: First layer weights; has shape (H, D)
27             b1: First layer biases; has shape (H,)
28             W2: Second layer weights; has shape (C, H)
29             b2: Second layer biases; has shape (C,)
30
31             Inputs:
32             - input_size: The dimension D of the input data.
33             - hidden_size: The number of neurons H in the hidden layer.
34             - output_size: The number of classes C.
35             """
36
37         self.params = {}
38         self.params['W1'] = std * np.random.randn(hidden_size, input_size)
39         self.params['b1'] = np.zeros(hidden_size)
40         self.params['W2'] = std * np.random.randn(output_size, hidden_size)
41         self.params['b2'] = np.zeros(output_size)
42
43     def loss(self, X, y=None, reg=0.0):
44         """
45             Compute the loss and gradients for a two layer fully connected neural
46             network.
47
48             Inputs:
49             - X: Input data of shape (N, D). Each X[i] is a training sample.
50             - y: Vector of training labels. y[i] is the label for X[i], and each y[i] is
51                 an integer in the range 0 <= y[i] < C. This parameter is optional; if it
52                 is not passed then we only return scores, and if it is passed then we
53                 instead return the loss and gradients.
54             - reg: Regularization strength.
55
56             Returns:
57             If y is None, return a matrix scores of shape (N, C) where scores[i, c] is
58             the score for class c on input X[i].
59
60             If y is not None, instead return a tuple of:
61             - loss: Loss (data loss and regularization loss) for this batch of training
62                 samples.
63             - grads: Dictionary mapping parameter names to gradients of those parameters
64                 with respect to the loss function; has the same keys as self.params.
65             """
66
67         # Unpack variables from the params dictionary
68         W1, b1 = self.params['W1'], self.params['b1']
69         W2, b2 = self.params['W2'], self.params['b2']
70         N, D = X.shape
71
72         # Compute the forward pass
73         scores = None

```

```

73      # ===== #
74      # YOUR CODE HERE:
75      #   Calculate the output scores of the neural network. The result
76      #   should be (N, C). As stated in the description for this class,
77      #   there should not be a ReLU layer after the second FC layer.
78      #   The output of the second FC layer is the output scores. Do not
79      #   use a for loop in your implementation.
80      # ===== #
81
82
83      layer_1 = X.dot(W1.T)+b1
84      relu_1 = np.maximum(layer_1,0)
85      layer_2 = relu_1.dot(W2.T)+b2
86      scores = layer_2
87
88      # ===== #
89      # END YOUR CODE HERE
90      # ===== #
91
92
93      # If the targets are not given then jump out, we're done
94      if y is None:
95          return scores
96
97      # Compute the loss
98      loss = None
99
100     # ===== #
101     # YOUR CODE HERE:
102     #   Calculate the loss of the neural network. This includes the
103     #   softmax loss and the L2 regularization for W1 and W2. Store the
104     #   total loss in teh variable loss. Multiply the regularization
105     #   loss by 0.5 (in addition to the factor reg).
106     # ===== #
107
108     # scores is num_examples by num_classes
109     scores = scores.T
110     ma_scores = scores.max(axis=0, keepdims=True)
111     scores -= ma_scores
112     exp_scores = np.exp(scores) #(c,n)
113     scores_sum = np.sum(exp_scores, axis=0) #(1,n)
114     loglike =
115     np.log(np.squeeze(np.take_along_axis(exp_scores,y[np.newaxis,:,0],0))/scores_sum)
116     softmax_loss = -np.average(loglike)
117
118     l2_reg_loss = 0.5*(np.sum(np.square(W1)) + np.sum(np.square(W2)))
119     loss = softmax_loss + reg*l2_reg_loss
120
121     # ===== #
122     # END YOUR CODE HERE
123     # ===== #
124
125     grads = {}
126
127     # ===== #
128     # YOUR CODE HERE:
129     #   Implement the backward pass. Compute the derivatives of the
130     #   weights and the biases. Store the results in the grads
131     #   dictionary. e.g., grads['W1'] should store the gradient for
132     #   W1, and be of the same size as W1.
133     # ===== #
134
135     scores_temp = (exp_scores/scores_sum)
136     scores_temp[y, np.arange(X.shape[0])] -= 1
137     scores_temp/= X.shape[0]
138     grads['W2'] = (scores_temp).dot(relu_1)
139     grads['b2'] = np.sum(scores_temp, axis=1)
140
141     dld_relu = scores_temp.T.dot(W2)
142     dld_relu[layer_1<0] = 0
143     grads['b1'] = np.sum(dld_relu, axis=0)
144     grads['W1'] = dld_relu.T.dot(X)

```

```

144
145     grads['W1'] += reg*W1
146     grads['W2'] += reg*W2
147
148     # ===== #
149     # END YOUR CODE HERE
150     # ===== #
151
152     return loss, grads
153
154 def train(self, X, y, X_val, y_val,
155             learning_rate=1e-3, learning_rate_decay=0.95,
156             reg=1e-5, num_iters=100,
157             batch_size=200, verbose=False):
158     """
159     Train this neural network using stochastic gradient descent.
160
161     Inputs:
162     - X: A numpy array of shape (N, D) giving training data.
163     - y: A numpy array f shape (N,) giving training labels; y[i] = c means that
164       X[i] has label c, where 0 <= c < C.
165     - X_val: A numpy array of shape (N_val, D) giving validation data.
166     - y_val: A numpy array of shape (N_val,) giving validation labels.
167     - learning_rate: Scalar giving learning rate for optimization.
168     - learning_rate_decay: Scalar giving factor used to decay the learning rate
169       after each epoch.
170     - reg: Scalar giving regularization strength.
171     - num_iters: Number of steps to take when optimizing.
172     - batch_size: Number of training examples to use per step.
173     - verbose: boolean; if true print progress during optimization.
174     """
175     num_train = X.shape[0]
176     iterations_per_epoch = max(num_train / batch_size, 1)
177
178     # Use SGD to optimize the parameters in self.model
179     loss_history = []
180     train_acc_history = []
181     val_acc_history = []
182
183     for it in np.arange(num_iters):
184         X_batch = None
185         y_batch = None
186
187         # ===== #
188         # YOUR CODE HERE:
189         #   Create a minibatch by sampling batch_size samples randomly.
190         # ===== #
191         mask = np.random.choice(num_train, batch_size, replace=True)
192         X_batch = X[mask]
193         y_batch = y[mask]
194
195         # ===== #
196         # END YOUR CODE HERE
197         # ===== #
198
199         # Compute loss and gradients using the current minibatch
200         loss, grads = self.loss(X_batch, y=y_batch, reg=reg)
201         loss_history.append(loss)
202
203         # ===== #
204         # YOUR CODE HERE:
205         #   Perform a gradient descent step using the minibatch to update
206         #   all parameters (i.e., W1, W2, b1, and b2).
207         # ===== #
208
209         self.params['W1'] -= learning_rate * grads['W1']
210         self.params['b1'] -= learning_rate * grads['b1']
211         self.params['W2'] -= learning_rate * grads['W2']
212         self.params['b2'] -= learning_rate * grads['b2']
213
214         # ===== #
215         # END YOUR CODE HERE
216         # ===== #

```

```

216
217     if verbose and it % 100 == 0:
218         print('iteration {} / {}: loss {}'.format(it, num_iters, loss))
219
220     # Every epoch, check train and val accuracy and decay learning rate.
221     if it % iterations_per_epoch == 0:
222         # Check accuracy
223         train_acc = (self.predict(X_batch) == y_batch).mean()
224         val_acc = (self.predict(X_val) == y_val).mean()
225         train_acc_history.append(train_acc)
226         val_acc_history.append(val_acc)
227
228         # Decay learning rate
229         learning_rate *= learning_rate_decay
230
231     return {
232         'loss_history': loss_history,
233         'train_acc_history': train_acc_history,
234         'val_acc_history': val_acc_history,
235     }
236
237 def predict(self, X):
238     """
239     Use the trained weights of this two-layer network to predict labels for
240     data points. For each data point we predict scores for each of the C
241     classes, and assign each data point to the class with the highest score.
242
243     Inputs:
244     - X: A numpy array of shape (N, D) giving N D-dimensional data points to
245       classify.
246
247     Returns:
248     - y_pred: A numpy array of shape (N,) giving predicted labels for each of
249       the elements of X. For all i, y_pred[i] = c means that X[i] is predicted
250       to have class c, where 0 <= c < C.
251     """
252     y_pred = None
253
254     # ===== #
255     # YOUR CODE HERE:
256     #   Predict the class given the input data.
257     # ===== #
258     scores = self.loss(X)
259     softmax = np.exp(scores)/np.sum(np.exp(scores), axis=1, keepdims=True)
# print("sc", scores.shape, softmax.shape, X.shape)
261     y_pred = np.zeros(X.shape[0])
262     for i in range(X.shape[0]):
263         y_pred[i] = np.argmax(softmax[i])
264
265     # ===== #
266     # END YOUR CODE HERE
267     # ===== #
268
269     return y_pred
270

```

# 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 functions 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 [24]:

```
## Import and setups

import time
import numpy as np
import matplotlib.pyplot as plt
from nn1.fc_net import *
from utils.data_utils import get_CIFAR10_data
```

```

from utils.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
from utils.solver import Solver

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

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

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

```

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

```
%reload_ext autoreload
```

In [3]:

```

# 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 [6]:

```

# 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: 9.769847728806635e-10
```

## Affine layer backward pass

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

In [14]:

```
# 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: 3.282907944638132e-10
dw error: 7.219169322822716e-11
db error: 7.385661415666503e-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 [16]:

```
# 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,  0.13636364,],
                      [ 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 [18]:

```
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: 3.2756153540684154e-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 [19]:

```
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: 2.4075713315106393e-10  
dw error: 6.442206192819315e-10  
db error: 3.2756008250058654e-12

## Softmax losses

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

In [20]:

```
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: 2.3026990443765962
dx error: 8.984473846912682e-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 [30]:

```
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, 15.33206765, 16.09215,
     [12.05769098, 12.74614105, 13.43459113, 14.1230412, 14.81149128, 15.49994135, 16.18839
     [12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.05099822, 15.66781506, 16.28463
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: 1.521570279286004e-08
W2 relative error: 3.4803693682531243e-10
b1 relative error: 6.5485474139109215e-09
```

```
b2 relative error: 4.3291413857436005e-10
Running numeric gradient check with reg = 0.7
W1 relative error: 8.175466200078585e-07
W2 relative error: 2.8508696990815807e-08
b1 relative error: 1.0895946645012713e-09
b2 relative error: 9.089615724390711e-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 [33]:

```
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.
#
# ===== #
lr = 0.001584893192461114
batch_size = 220
reg = 0.12
model = TwoLayerNet(hidden_dims=2000, reg = reg)
solver = Solver(model,data,batch_size=batch_size,
                optim_config={
                    'learning_rate': lr,
                },
                lr_decay=0.95,
                num_epochs=10,print_every=100)
solver.train()
# ===== #
# END YOUR CODE HERE
# ===== #
```

```
(Iteration 1 / 2220) loss: 2.664031
(Epoch 0 / 10) train acc: 0.153000; val_acc: 0.172000
(Iteration 101 / 2220) loss: 1.962112
(Iteration 201 / 2220) loss: 1.839125
(Epoch 1 / 10) train acc: 0.441000; val_acc: 0.426000
(Iteration 301 / 2220) loss: 1.772910
(Iteration 401 / 2220) loss: 1.806921
(Epoch 2 / 10) train acc: 0.514000; val_acc: 0.481000
(Iteration 501 / 2220) loss: 1.717610
(Iteration 601 / 2220) loss: 1.792207
(Epoch 3 / 10) train acc: 0.504000; val_acc: 0.459000
(Iteration 701 / 2220) loss: 1.617884
(Iteration 801 / 2220) loss: 1.798771
(Epoch 4 / 10) train acc: 0.559000; val_acc: 0.516000
(Iteration 901 / 2220) loss: 1.645398
(Iteration 1001 / 2220) loss: 1.483347
(Iteration 1101 / 2220) loss: 1.450971
(Epoch 5 / 10) train acc: 0.585000; val_acc: 0.505000
(Iteration 1201 / 2220) loss: 1.393434
(Iteration 1301 / 2220) loss: 1.533003
(Epoch 6 / 10) train acc: 0.551000; val_acc: 0.481000
(Iteration 1401 / 2220) loss: 1.479555
(Iteration 1501 / 2220) loss: 1.446862
(Epoch 7 / 10) train acc: 0.642000; val_acc: 0.513000
(Iteration 1601 / 2220) loss: 1.397423
(Iteration 1701 / 2220) loss: 1.381091
(Epoch 8 / 10) train acc: 0.605000; val_acc: 0.526000
(Iteration 1801 / 2220) loss: 1.416650
(Iteration 1901 / 2220) loss: 1.246269
```

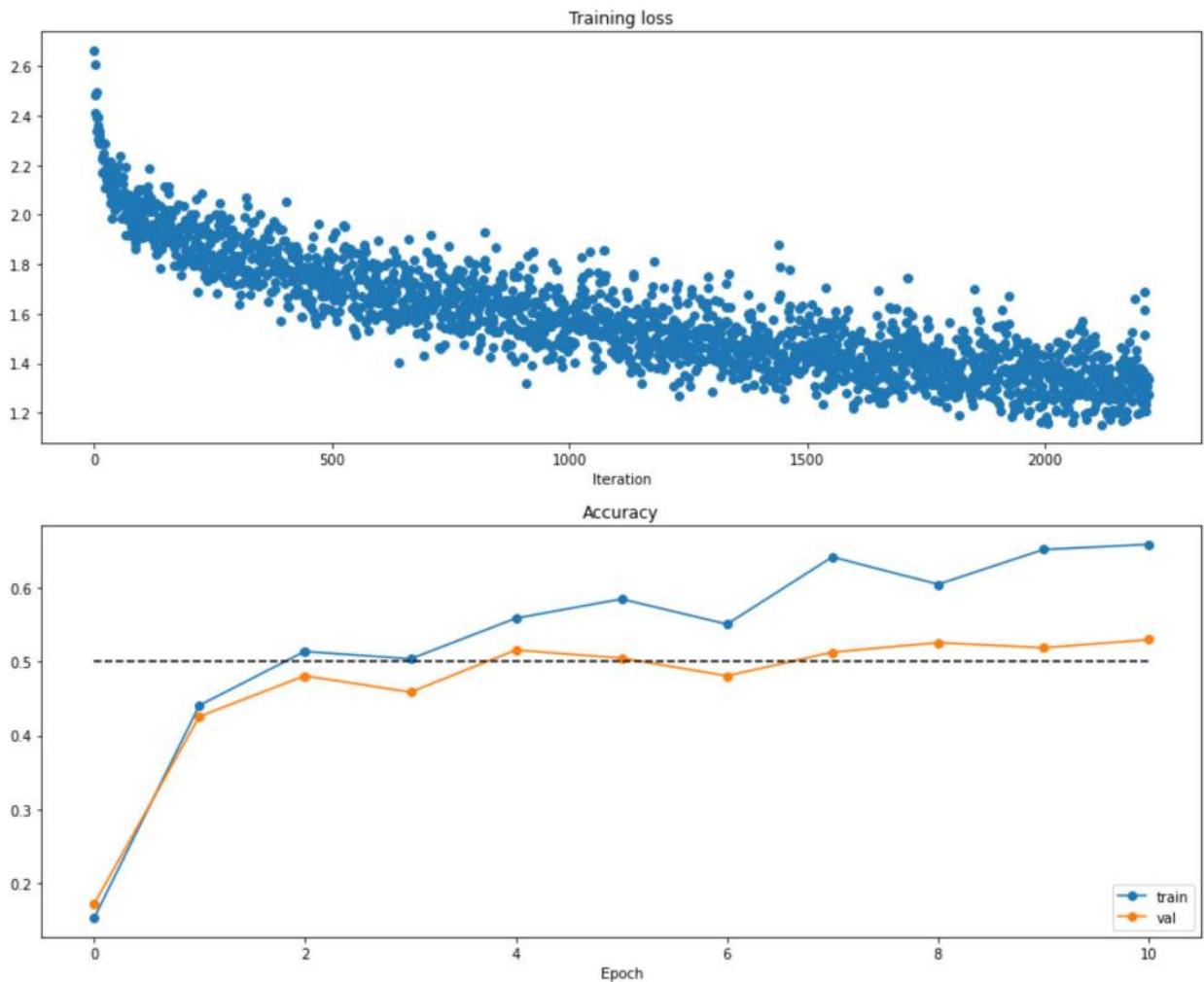
```
(Epoch 9 / 10) train acc: 0.652000; val_acc: 0.519000
(Iteration 2001 / 2220) loss: 1.172493
(Iteration 2101 / 2220) loss: 1.475718
(Iteration 2201 / 2220) loss: 1.285669
(Epoch 10 / 10) train acc: 0.659000; val_acc: 0.530000
```

In [34]:

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

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

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

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

```
Running check with reg = 0
Initial loss: 2.296592349183469
W1 relative error: 7.392250087287479e-07
W2 relative error: 6.902284317100103e-07
W3 relative error: 6.188395661084466e-08
b1 relative error: 4.0166272330831426e-08
b2 relative error: 2.9521518990922087e-09
b3 relative error: 1.6248826225499729e-10
Running check with reg = 3.14
Initial loss: 7.038390430408256
W1 relative error: 1.051749311442295e-08
W2 relative error: 6.891098770538369e-09
W3 relative error: 1.0231239706316546e-08
b1 relative error: 7.952042923096436e-09
b2 relative error: 1.3744697639568225e-09
b3 relative error: 1.641431444993658e-10
```

```
In [52]: # Use the three layer neural network to overfit a small dataset.

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

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

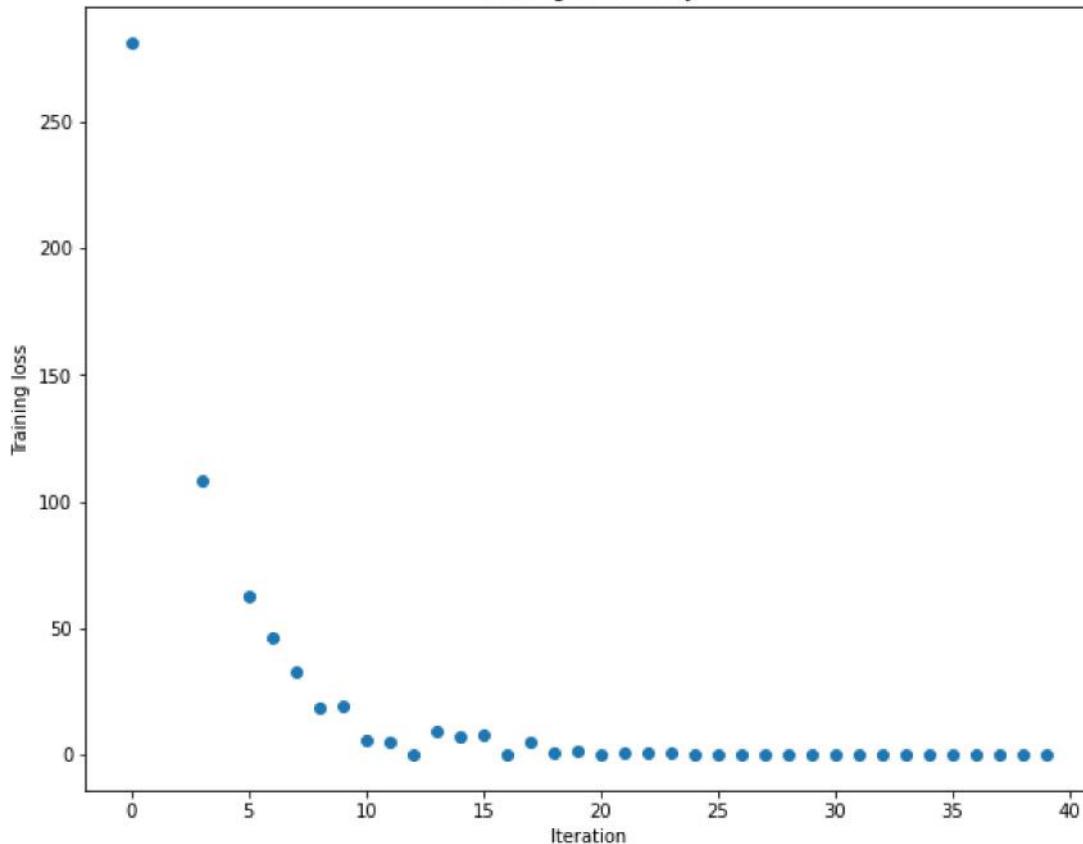
model = FullyConnectedNet([100, 100],
                          weight_scale=weight_scale, dtype=np.float64)
solver = Solver(model, small_data,
                print_every=10, num_epochs=20, batch_size=25,
                update_rule='sgd',
                optim_config={
                    'learning_rate': 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: 281.562713  
(Epoch 0 / 20) train acc: 0.160000; val\_acc: 0.112000

```
(Epoch 1 / 20) train acc: 0.200000; val_acc: 0.113000
(Epoch 2 / 20) train acc: 0.440000; val_acc: 0.136000
(Epoch 3 / 20) train acc: 0.560000; val_acc: 0.164000
(Epoch 4 / 20) train acc: 0.760000; val_acc: 0.176000
(Epoch 5 / 20) train acc: 0.780000; val_acc: 0.177000
(Iteration 11 / 40) loss: 5.939655
(Epoch 6 / 20) train acc: 0.900000; val_acc: 0.175000
(Epoch 7 / 20) train acc: 0.900000; val_acc: 0.175000
(Epoch 8 / 20) train acc: 0.900000; val_acc: 0.173000
(Epoch 9 / 20) train acc: 0.940000; val_acc: 0.164000
(Epoch 10 / 20) train acc: 0.980000; val_acc: 0.178000
(Iteration 21 / 40) loss: 0.000061
(Epoch 11 / 20) train acc: 0.960000; val_acc: 0.179000
(Epoch 12 / 20) train acc: 0.980000; val_acc: 0.179000
(Epoch 13 / 20) train acc: 0.980000; val_acc: 0.179000
(Epoch 14 / 20) train acc: 0.980000; val_acc: 0.179000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.176000
(Iteration 31 / 40) loss: 0.000000
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.176000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.176000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.176000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.176000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.176000
```

Training loss history



In [ ]:

In [ ]:

```

1 import numpy as np
2
3 from .layers import *
4 from .layer_utils import *
5
6
7 class TwoLayerNet(object):
8     """
9         A two-layer fully-connected neural network with ReLU nonlinearity and
10        softmax loss that uses a modular layer design. We assume an input dimension
11        of D, a hidden dimension of H, and perform classification over C classes.
12
13        The architecture should be affine - relu - affine - softmax.
14
15        Note that this class does not implement gradient descent; instead, it
16        will interact with a separate Solver object that is responsible for running
17        optimization.
18
19        The learnable parameters of the model are stored in the dictionary
20        self.params that maps parameter names to numpy arrays.
21        """
22
23    def __init__(self, input_dim=3*32*32, hidden_dims=100, num_classes=10,
24                 dropout=0, weight_scale=1e-3, reg=0.0):
25        """
26            Initialize a new network.
27
28            Inputs:
29            - input_dim: An integer giving the size of the input
30            - hidden_dims: An integer giving the size of the hidden layer
31            - num_classes: An integer giving the number of classes to classify
32            - dropout: Scalar between 0 and 1 giving dropout strength.
33            - weight_scale: Scalar giving the standard deviation for random
34                initialization of the weights.
35            - reg: Scalar giving L2 regularization strength.
36        """
37        self.params = {}
38        self.reg = reg
39
40        # ===== #
41        # YOUR CODE HERE:
42        #   Initialize W1, W2, b1, and b2. Store these as self.params['W1'],
43        #   self.params['W2'], self.params['b1'] and self.params['b2']. The
44        #   biases are initialized to zero and the weights are initialized
45        #   so that each parameter has mean 0 and standard deviation weight_scale.
46        #   The dimensions of W1 should be (input_dim, hidden_dim) and the
47        #   dimensions of W2 should be (hidden_dims, num_classes)
48        # ===== #
49        self.input_dim = input_dim
50        self.hidden_dims = hidden_dims
51        self.num_classes = num_classes
52        self.weight_scale = weight_scale
53        self.params['W1'] = np.random.normal(loc=0.0, scale=weight_scale, size =
54                                              (self.input_dim, self.hidden_dims))
55        self.params['b1'] = np.zeros((self.hidden_dims,))
56        self.params['W2'] = np.random.normal(loc=0.0, scale=weight_scale, size =
57                                              (self.hidden_dims, self.num_classes))
58        self.params['b2'] = np.zeros((self.num_classes,))
59
60        # ===== #
61        # END YOUR CODE HERE
62        # ===== #
63
64    def loss(self, X, y=None):
65        """
66            Compute loss and gradient for a minibatch of data.
67
68            Inputs:
69            - X: Array of input data of shape (N, d_1, ..., d_k)
70            - y: Array of labels, of shape (N,). y[i] gives the label for X[i].
71
72            Returns:

```

```

71     If y is None, then run a test-time forward pass of the model and return:
72     - scores: Array of shape (N, C) giving classification scores, where
73         scores[i, c] is the classification score for X[i] and class c.
74
75     If y is not None, then run a training-time forward and backward pass and
76     return a tuple of:
77     - loss: Scalar value giving the loss
78     - grads: Dictionary with the same keys as self.params, mapping parameter
79         names to gradients of the loss with respect to those parameters.
80     """
81     scores = None
82
83     # ===== #
84     # YOUR CODE HERE:
85     # Implement the forward pass of the two-layer neural network. Store
86     # the class scores as the variable 'scores'. Be sure to use the layers
87     # you prior implemented.
88     # ===== #
89
90     layer_1, cache_1 = affine_relu_forward(X, self.params['W1'], self.params['b1'])
91     scores, cache_2 = affine_forward(layer_1, self.params['W2'], self.params['b2'])
92
93     # ===== #
94     # END YOUR CODE HERE
95     # ===== #
96
97     # If y is None then we are in test mode so just return scores
98     if y is None:
99         return scores
100
101    loss, grads = 0, {}
102    # ===== #
103    # YOUR CODE HERE:
104    # Implement the backward pass of the two-layer neural net. Store
105    # the loss as the variable 'loss' and store the gradients in the
106    # 'grads' dictionary. For the grads dictionary, grads['W1'] holds
107    # the gradient for W1, grads['b1'] holds the gradient for b1, etc.
108    # i.e., grads[k] holds the gradient for self.params[k].
109    #
110    # Add L2 regularization, where there is an added cost 0.5*self.reg*W^2
111    # for each W. Be sure to include the 0.5 multiplying factor to
112    # match our implementation.
113    #
114    # And be sure to use the layers you prior implemented.
115    # ===== #
116
117    loss, d2 = softmax_loss(scores, y)
118    loss+= 0.5*self.reg*(np.sum(np.square(self.params['W1'])) + 
119    np.sum(np.square(self.params['W2'])))
120    d1,dw2,db2 = affine_backward(d2,cache_2)
121    dx,dw1,db1 = affine_relu_backward(d1,cache_1)
122    grads['W1']= dw1 + self.reg*self.params['W1']
123    grads['b1']= db1
124    grads['W2']= dw2 + self.reg*self.params['W2']
125    grads['b2']= db2
126
127    # ===== #
128    # END YOUR CODE HERE
129    # ===== #
130
131    return loss, grads
132
133 class FullyConnectedNet(object):
134     """
135     A fully-connected neural network with an arbitrary number of hidden layers,
136     ReLU nonlinearities, and a softmax loss function. This will also implement
137     dropout and batch normalization as options. For a network with L layers,
138     the architecture will be
139
140     {affine - [batch norm] - relu - [dropout]} x (L - 1) - affine - softmax
141

```

```

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

```

```

212 # With batch normalization we need to keep track of running means and
213 # variances, so we need to pass a special bn_param object to each batch
214 # normalization layer. You should pass self.bn_params[0] to the forward pass
215 # of the first batch normalization layer, self.bn_params[1] to the forward
216 # pass of the second batch normalization layer, etc.
217 self.bn_params = []
218 if self.use_batchnorm:
219     self.bn_params = [{ 'mode': 'train' } for i in np.arange(self.num_layers - 1)]
220
221 # Cast all parameters to the correct datatype
222 for k, v in self.params.items():
223     self.params[k] = v.astype(dtype)
224
225
226 def loss(self, X, y=None):
227     """
228     Compute loss and gradient for the fully-connected net.
229
230     Input / output: Same as TwoLayerNet above.
231     """
232     X = X.astype(self.dtype)
233     mode = 'test' if y is None else 'train'
234
235     # Set train/test mode for batchnorm params and dropout param since they
236     # behave differently during training and testing.
237     if self.dropout_param is not None:
238         self.dropout_param['mode'] = mode
239     if self.use_batchnorm:
240         for bn_param in self.bn_params:
241             bn_param[mode] = mode
242
243     scores = None
244
245     # ===== #
246     # YOUR CODE HERE:
247     #   Implement the forward pass of the FC net and store the output
248     #   scores as the variable "scores".
249     # ===== #
250     act = []
251     cache = []
252     for i in range(self.num_layers):
253         w = 'W'+str(i+1)
254         b = 'b'+str(i+1)
255         if i==0:
256             a,c = affine_relu_forward(X,self.params[w],self.params[b])
257             act.append(a)
258         elif i==self.num_layers-1:
259             a,c = affine_forward(act[i-1],self.params[w],self.params[b])
260             scores=a
261         else:
262             a,c = affine_relu_forward(act[i-1],self.params[w],self.params[b])
263             act.append(a)
264             cache.append(c)
265
266
267     # ===== #
268     # END YOUR CODE HERE
269     # ===== #
270
271     # If test mode return early
272     if mode == 'test':
273         return scores
274
275
276     loss, grads = 0.0, {}
277     # ===== #
278     # YOUR CODE HERE:
279     #   Implement the backwards pass of the FC net and store the gradients
280     #   in the grads dict, so that grads[k] is the gradient of self.params[k]
281     #   Be sure your L2 regularization includes a 0.5 factor.
282     # ===== #
283     loss=0

```

```
284 loss, ds = softmax_loss(scores,y)
285 for i in np.arange(self.num_layers-1,-1,-1):
286     w = 'W'+str(i+1)
287     b = 'b'+str(i+1)
288     loss += self.reg*0.5*np.sum(np.square(self.params[w]))
289     if i==self.num_layers-1:
290         da,grads[w], grads[b] = affine_backward(ds,cache[i])
291     else:
292         da,grads[w], grads[b] = affine_relu_backward(da,cache[i])
293     grads[w] += self.reg*self.params[w]
294
295
296
297 # ===== #
298 # END YOUR CODE HERE
299 # ===== #
300
301 return loss, grads
```

```

1 import numpy as np
2 import pdb
3
4
5
6
7 def affine_forward(x, w, b):
8     """
9         Computes the forward pass for an affine (fully-connected) layer.
10
11        The input x has shape (N, d_1, ..., d_k) and contains a minibatch of N
12        examples, where each example x[i] has shape (d_1, ..., d_k). We will
13        reshape each input into a vector of dimension D = d_1 * ... * d_k, and
14        then transform it to an output vector of dimension M.
15
16    Inputs:
17    - x: A numpy array containing input data, of shape (N, d_1, ..., d_k)
18    - w: A numpy array of weights, of shape (D, M)
19    - b: A numpy array of biases, of shape (M,)
20
21    Returns a tuple of:
22    - out: output, of shape (N, M)
23    - cache: (x, w, b)
24    """
25
26    # ===== #
27    # YOUR CODE HERE:
28    #   Calculate the output of the forward pass. Notice the dimensions
29    #   of w are D x M, which is the transpose of what we did in earlier
30    #   assignments.
31    # ===== #
32
33    out = x.reshape(x.shape[0], -1).dot(w)+b
34
35    # ===== #
36    # END YOUR CODE HERE
37    # ===== #
38
39    cache = (x, w, b)
40    return out, cache
41
42
43 def affine_backward(dout, cache):
44     """
45         Computes the backward pass for an affine layer.
46
47     Inputs:
48     - dout: Upstream derivative, of shape (N, M)
49     - cache: Tuple of:
50         - x: Input data, of shape (N, d_1, ..., d_k)
51         - w: Weights, of shape (D, M)
52
53     Returns a tuple of:
54     - dx: Gradient with respect to x, of shape (N, d1, ..., d_k)
55     - dw: Gradient with respect to w, of shape (D, M)
56     - db: Gradient with respect to b, of shape (M,)
57     """
58
59     x, w, b = cache
60     dx, dw, db = None, None, None
61
62    # ===== #
63    # YOUR CODE HERE:
64    #   Calculate the gradients for the backward pass.
65    # ===== #
66
67    # dout is N x M
68    # dx should be N x d1 x ... x dk; it relates to dout through multiplication with
69    # w, which is D x M
70    # dw should be D x M; it relates to dout through multiplication with x, which is N
71    # x D after reshaping
72    # db should be M; it is just the sum over dout examples

```

```

71 #print(dout.shape,w.shape,x.shape)
72 dx = dout.dot(w.T).reshape(x.shape)
73 dw = x.reshape(x.shape[0],-1).T.dot(dout).reshape(w.shape)
74 db = np.sum(dout, axis=0)
75
76
77 # ===== #
78 # END YOUR CODE HERE
79 # ===== #
80
81 return dx, dw, db
82
83 def relu_forward(x):
84 """
85     Computes the forward pass for a layer of rectified linear units (ReLUs).
86
87     Input:
88     - x: Inputs, of any shape
89
90     Returns a tuple of:
91     - out: Output, of the same shape as x
92     - cache: x
93 """
94 # ===== #
95 # YOUR CODE HERE:
96 #     Implement the ReLU forward pass.
97 # ===== #
98
99 out = np.maximum(x,0)
100 # ===== #
101 # END YOUR CODE HERE
102 # ===== #
103
104 cache = x
105 return out, cache
106
107
108 def relu_backward(dout, cache):
109 """
110     Computes the backward pass for a layer of rectified linear units (ReLUs).
111
112     Input:
113     - dout: Upstream derivatives, of any shape
114     - cache: Input x, of same shape as dout
115
116     Returns:
117     - dx: Gradient with respect to x
118 """
119 x = cache
120
121 # ===== #
122 # YOUR CODE HERE:
123 #     Implement the ReLU backward pass
124 # ===== #
125
126 # ReLU directs linearly to those > 0
127 dx = dout
128 dx[x<0]=0
129
130 # ===== #
131 # END YOUR CODE HERE
132 # ===== #
133
134 return dx
135
136 def svm_loss(x, y):
137 """
138     Computes the loss and gradient using for multiclass SVM classification.
139
140     Inputs:
141     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
142         for the ith input.

```

```

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

```