

Q1

Given

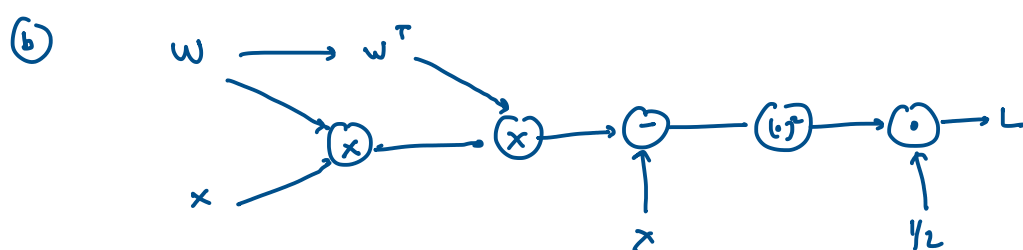
$$x \in \mathbb{R}^n$$

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

$Wx$  is of lower dimensionality than  $x$

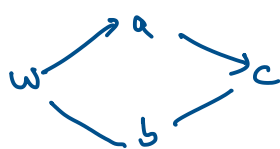
$$L = \frac{1}{2} \|W^T W x - x\|^2$$

- (a) In the loss function, we are minimizing the difference between  $W^T W x$  and  $x$ . Thus the representation  $Wx$  will be a lower dimensional representation of  $x$  and will preserve the information about  $x$

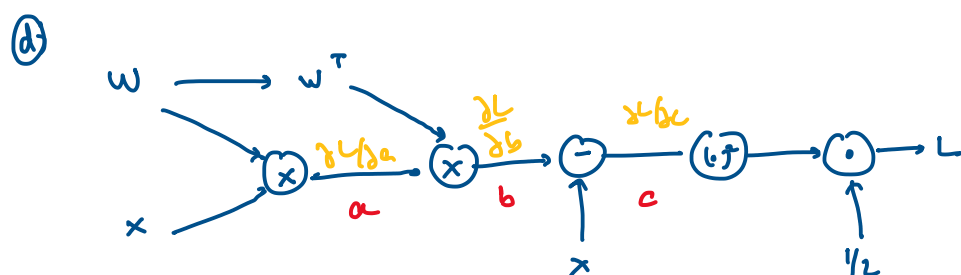


- (c) We will be using law of total derivatives & adding the derivatives from both the paths - as in the example:

eg:



$$\frac{\partial c}{\partial W} = \frac{\partial a}{\partial W} \frac{\partial c}{\partial a} + \frac{\partial b}{\partial W} \frac{\partial c}{\partial b}$$

To find  $\nabla_W L$ 

$$a = Wx$$

$$b = W^T W x$$

$$c = W^T W x - x$$

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

$$\Rightarrow \frac{\partial L}{\partial c} = c = W^T W x - x$$

$$c = b - x$$

$$\Rightarrow \frac{\partial L}{\partial b} = \frac{\partial c}{\partial b} \frac{\partial L}{\partial c} = c = W^T W x - x$$

$$b = W^T a$$

$$\frac{\partial L}{\partial a} = \frac{\partial b}{\partial a} \frac{\partial L}{\partial b} = W^T c$$

$$\frac{\partial L}{\partial W^T} = \frac{\partial b}{\partial W^T} \frac{\partial L}{\partial b} = c a^T$$

$$\begin{aligned}
 &= W(W^T W x - x) & &= (W^T W x - x)(W x)^T \\
 a &= W x & & \\
 \frac{\partial L}{\partial W} &= \frac{\partial a}{\partial W} \frac{\partial L}{\partial a} & & \frac{\partial L}{\partial W} = \left( \frac{\partial L}{\partial W^T} \right)^T \\
 &= \underline{W(W^T W x - x) x^T} & & = \underline{W x (W^T W x - x)^T}
 \end{aligned}$$

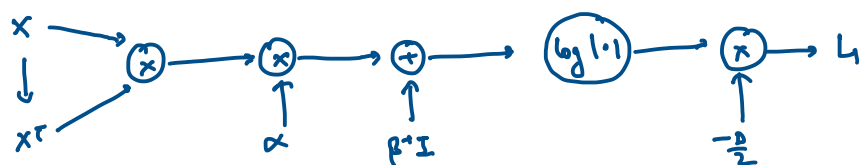
By law of total derivatives,

$$\begin{aligned}
 \Rightarrow \nabla_W L &= \text{backprop to } W + \text{backprop to } W^T \\
 &= W(W^T W x - x) x^T + W x (W^T W x - x)^T
 \end{aligned}$$

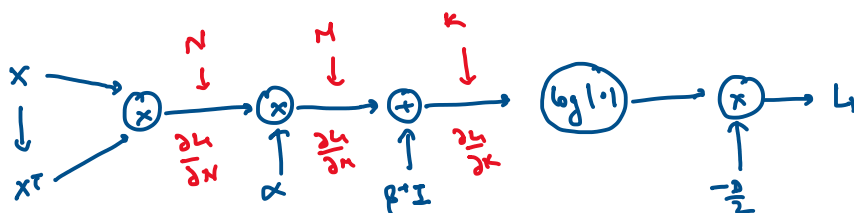
Q2

$$L_1 = -\frac{D}{2} \log |K K^T + \beta^2 I|$$

(a) computational graph for  $L_1$



(b) To compute  $\frac{\partial L_1}{\partial x}$



From computational graph

$$L_1 = -\frac{D}{2} \log |K|$$

Using the matrix cookbook

$$① \quad \frac{\partial L_1}{\partial K} = -\frac{D}{2} (K^{-1})^T = -\frac{D}{2} (K^T)^{-1}$$

$$\begin{aligned} \textcircled{2} \quad K &= M + \beta^+ I \\ \frac{\partial L_1}{\partial M} &= \frac{\partial K}{\partial M} \frac{\partial L_1}{\partial K} \\ \frac{\partial K}{\partial M} &= I \\ \Rightarrow \frac{\partial L_1}{\partial M} &= -\frac{D(K^T)^{-1}}{2} \end{aligned}$$

$\Rightarrow$  Addition operator distributes the gradient.

$$\textcircled{3} \quad \frac{\partial L_1}{\partial N} = \frac{\partial M}{\partial N} \frac{\partial L_1}{\partial M}$$

$$M = \alpha N$$

$\Rightarrow$  Multiplication operator routes the gradient

$$\Rightarrow \frac{\partial L_1}{\partial N} = -\alpha \frac{D(K^T)^{-1}}{2}$$

$$\textcircled{4} \quad N = X X^T$$

using hint given in the problem,

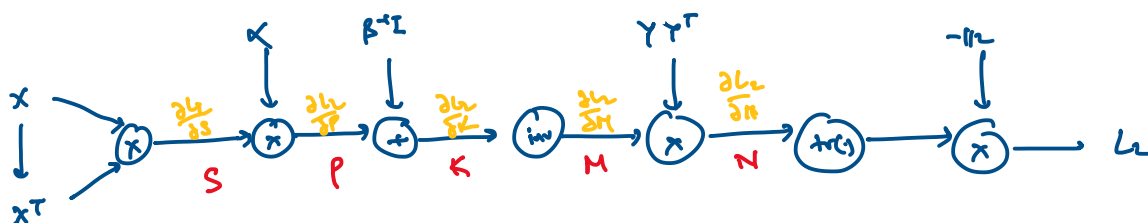
$$\frac{\partial L}{\partial X} = \frac{\partial N}{\partial X} \frac{\partial L}{\partial N}$$

$$= -\alpha D(K^T)^{-1} X$$

$$\text{Hence } \frac{\partial L_1}{\partial X} = -\alpha D(K^T)^{-1} X$$

$$L_2 = -\frac{1}{2} \text{tr}((\alpha X X^T + \beta^+ I)^{-1} Y Y^T)$$

(c) Computational graph for  $L_2$



(d) To compute  $\frac{\partial L_2}{\partial X}$

$$\textcircled{1} \quad L_2 = -\frac{1}{2} \text{tr}(N)$$

using matrix cookbook,

$$\frac{\partial L_2}{\partial N} = -\frac{1}{2} I$$

$$\textcircled{2} \quad N = M Y Y^T$$

$$\begin{aligned} \frac{\partial L_2}{\partial M} &= \frac{\partial N}{\partial M} \frac{\partial L_2}{\partial N} \\ &= -\frac{1}{2} Y Y^T \end{aligned}$$

$$(3) \quad M = K^{-1}$$

$$\begin{aligned} \frac{\partial L_2}{\partial K} &= \frac{\partial M}{\partial K} \frac{\partial L_2}{\partial M} \\ &= \frac{1}{2} K^T Y Y^T K^{-1} \end{aligned}$$

$$(4) \quad K = P + P^T I$$

$$\begin{aligned} \frac{\partial L_2}{\partial P} &= \frac{\partial K}{\partial P} \frac{\partial L_2}{\partial K} \\ &= \frac{1}{2} K^T Y Y^T K^{-1} \end{aligned}$$

$$(5) \quad P = \alpha S$$

$$\begin{aligned} \frac{\partial L_2}{\partial S} &= \frac{\partial P}{\partial S} \frac{\partial L_2}{\partial P} \\ &= \frac{1}{2} \alpha K^T Y Y^T K^{-1} \end{aligned}$$

$$(6) \quad P = X X^T$$

$$\begin{aligned} \frac{\partial L_2}{\partial X} &= \frac{\partial P}{\partial X} \frac{\partial L_2}{\partial P} \\ &= \alpha K^T Y Y^T K^{-1} X \end{aligned}$$

$$\boxed{\text{Hence } \frac{\partial L_2}{\partial X} = \alpha K^T Y Y^T K^{-1} X}$$

(e) To compute  $\frac{\partial L}{\partial X}$

$$L = -C - \frac{D}{2} \log |K| - \frac{1}{2} \text{tr}(K^T Y Y^T)$$

$$= -C + L_1 + L_2$$

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

$$= \boxed{-\alpha D (K^T)^{-1} X + \alpha K^T Y Y^T K^{-1} X}$$

# 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 [267...

```
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 [268...

```
from nndl.neural_net import TwoLayerNet
```

In [269...

```
# 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 [270... ## 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:
3.3812311957259755e-08
```

## Forward pass loss

```
In [271... 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 [272... from utils.gradient_check import eval_numerical_gradient

# Use numeric gradient checking to check your implementation of the backward
# 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, grads[param_name])))

```

```

b2 max relative error: 1.2482633693659668e-09
W2 max relative error: 2.9632233460136427e-10
b1 max relative error: 3.172680285697327e-09
W1 max relative error: 1.28328951808708e-09

```

## Training the network

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

In [273]...

```

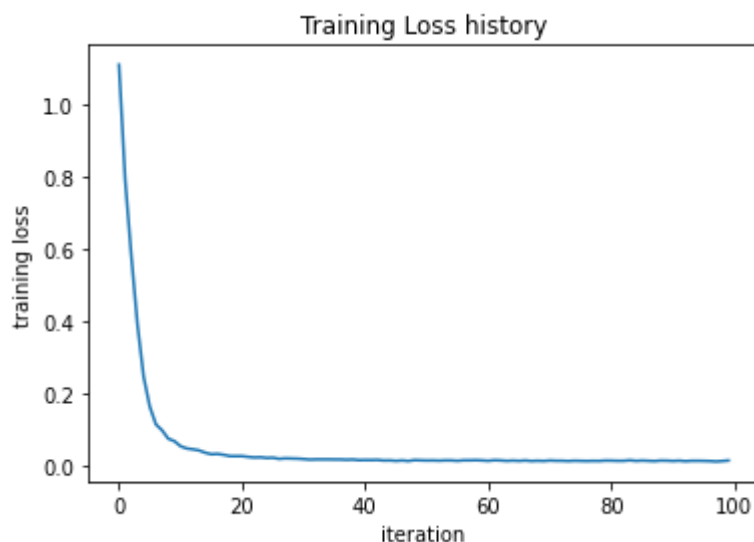
net = init_toy_model()
stats = net.train(X, y, X, y,
                  learning_rate=1e-1, reg=5e-6,
                  num_iters=100, verbose=False)

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

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

```

Final training loss: 0.014497864587765906



## Classify CIFAR-10

Do classification on the CIFAR-10 dataset.

In [274]...

```

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 = '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 [275...

```
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 2.302757518613176
iteration 100 / 1000: loss 2.302120159207236
iteration 200 / 1000: loss 2.2956136007408703
iteration 300 / 1000: loss 2.2518259043164135
iteration 400 / 1000: loss 2.188995235046776
iteration 500 / 1000: loss 2.1162527791897747
iteration 600 / 1000: loss 2.064670827698217
iteration 700 / 1000: loss 1.9901688623083942
iteration 800 / 1000: loss 2.002827640124685
iteration 900 / 1000: loss 1.9465176817856495
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 [254... stats['train_acc_history']
```

```
Out[254... [0.06, 0.145, 0.175, 0.295, 0.26]
```

```
In [277... # ===== #
# 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

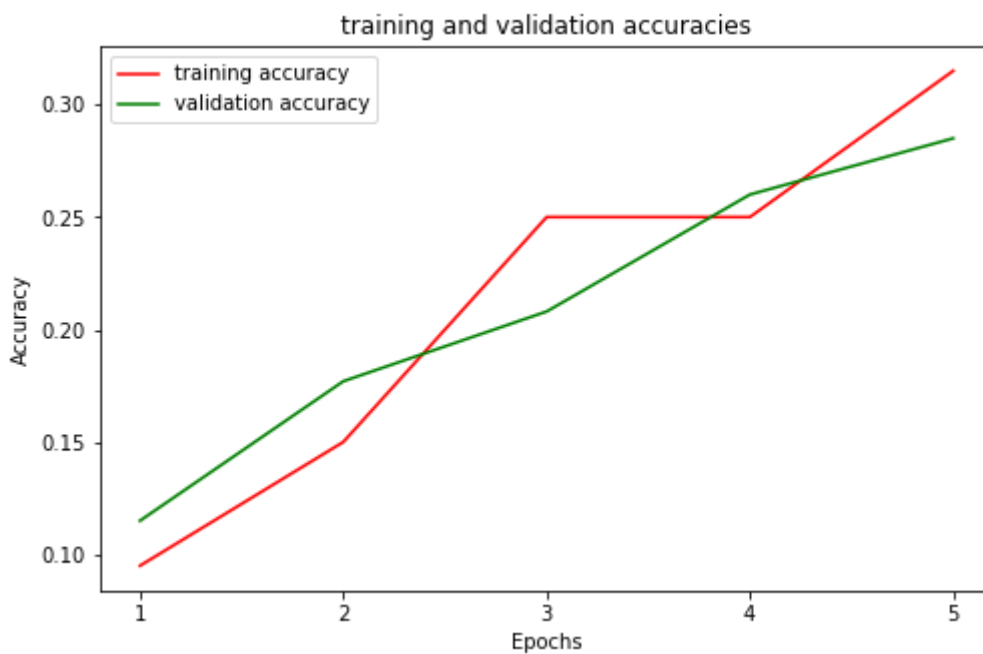
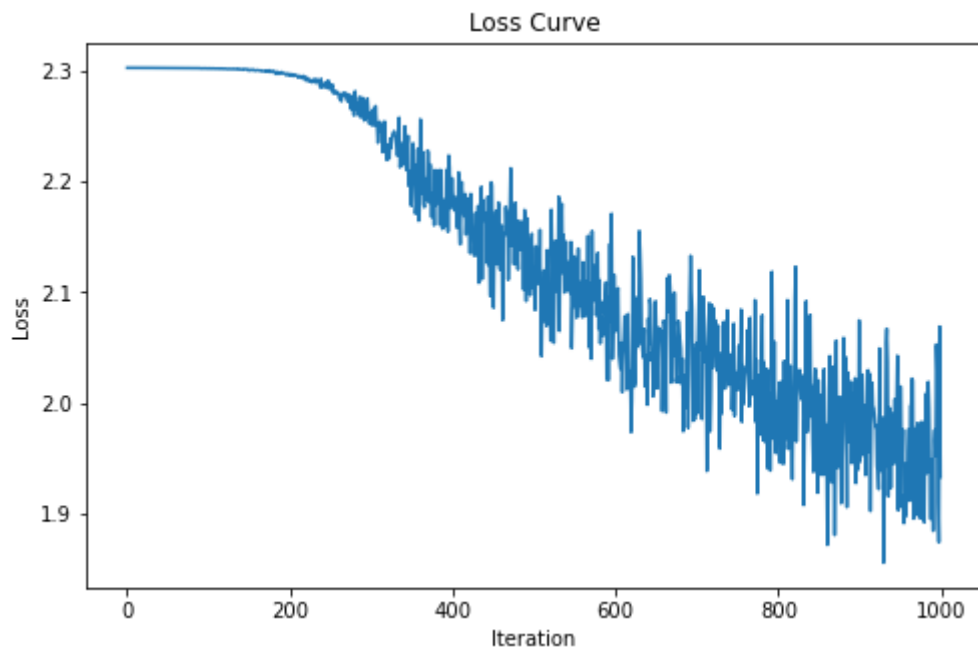
plt.figure(figsize=(8,5))
plt.title("Loss Curve")
plt.plot(stats['loss_history'])
plt.xlabel('Iteration')
plt.ylabel('Loss')

epochs = np.arange(1, len(stats['train_acc_history']) + 1)
plt.figure(figsize=(8,5))
plt.title("training and validation accuracies")
plt.xticks(epochs)
plt.plot(epochs, stats['train_acc_history'], label="training accuracy", color='red')
plt.plot(epochs, stats['val_acc_history'], label="validation accuracy", color='blue')

plt.xlabel("Epochs")
plt.ylabel("Accuracy")
plt.legend()
plt.show()
# ===== #
```

# END YOUR CODE HERE

# ===== #



## Answers:

(1) We can increase iterations or learning rate as the network hasn't reached its best performance. We can see that both training and validation accuracies are still increasing. We should ideally observe a graph where the validation accuracy saturates.

(2) We need to optimize the hyperparameters like iterations, learning rate, batch size, learning rate decay to improve model performance.

## 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 [278...

```

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)!
# ===== #
input_size = 32 * 32 * 3
num_classes = 10
hidden_size = 50

num_iterss = [1000, 3000, 5000]
batch_sizes = [100, 200, 300]
learning_rates = [5e-4, 1e-4, 1e-5]
learning_rate_decays = [0.93, 0.95, 0.97, 0.99]

hyperparams = ((num_iters, batch_size, learning_rate, learning_rate_decay)
                for num_iters in num_iterss
                for batch_size in batch_sizes
                for learning_rate in learning_rates
                for learning_rate_decay in learning_rate_decays)

for (num_iters, batch_size, learning_rate, learning_rate_decay) in hyperparams:
    print ((num_iters, batch_size, learning_rate, learning_rate_decay))
    net = TwoLayerNet(input_size, hidden_size, num_classes)

    # Train the network
    stats = net.train(X_train, y_train, X_val, y_val,
                      num_iters=num_iters, batch_size=batch_size,
                      learning_rate=learning_rate, learning_rate_decay=learning_rate_decay,
                      reg=0.25, verbose=False)

    # 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.
    if (best_net == None):
        best_net = net
    elif ((best_net.predict(X_val) == y_val).mean() < val_acc):
        best_net = net

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

```

```

(1000, 100, 0.0005, 0.93)
Validation accuracy: 0.436
(1000, 100, 0.0005, 0.95)
Validation accuracy: 0.444
(1000, 100, 0.0005, 0.97)
Validation accuracy: 0.423

```

```
(1000, 100, 0.0005, 0.99)
Validation accuracy: 0.444
(1000, 100, 0.0001, 0.93)
Validation accuracy: 0.28
(1000, 100, 0.0001, 0.95)
Validation accuracy: 0.279
(1000, 100, 0.0001, 0.97)
Validation accuracy: 0.287
(1000, 100, 0.0001, 0.99)
Validation accuracy: 0.297
(1000, 100, 1e-05, 0.93)
Validation accuracy: 0.206
(1000, 100, 1e-05, 0.95)
Validation accuracy: 0.184
(1000, 100, 1e-05, 0.97)
Validation accuracy: 0.24
(1000, 100, 1e-05, 0.99)
Validation accuracy: 0.231
(1000, 200, 0.0005, 0.93)
Validation accuracy: 0.438
(1000, 200, 0.0005, 0.95)
Validation accuracy: 0.462
(1000, 200, 0.0005, 0.97)
Validation accuracy: 0.45
(1000, 200, 0.0005, 0.99)
Validation accuracy: 0.474
(1000, 200, 0.0001, 0.93)
Validation accuracy: 0.27
(1000, 200, 0.0001, 0.95)
Validation accuracy: 0.281
(1000, 200, 0.0001, 0.97)
Validation accuracy: 0.28
(1000, 200, 0.0001, 0.99)
Validation accuracy: 0.299
(1000, 200, 1e-05, 0.93)
Validation accuracy: 0.247
(1000, 200, 1e-05, 0.95)
Validation accuracy: 0.223
(1000, 200, 1e-05, 0.97)
Validation accuracy: 0.253
(1000, 200, 1e-05, 0.99)
Validation accuracy: 0.206
(1000, 300, 0.0005, 0.93)
Validation accuracy: 0.465
(1000, 300, 0.0005, 0.95)
Validation accuracy: 0.443
(1000, 300, 0.0005, 0.97)
Validation accuracy: 0.458
(1000, 300, 0.0005, 0.99)
Validation accuracy: 0.471
(1000, 300, 0.0001, 0.93)
Validation accuracy: 0.291
(1000, 300, 0.0001, 0.95)
Validation accuracy: 0.298
(1000, 300, 0.0001, 0.97)
Validation accuracy: 0.293
(1000, 300, 0.0001, 0.99)
Validation accuracy: 0.295
(1000, 300, 1e-05, 0.93)
Validation accuracy: 0.187
(1000, 300, 1e-05, 0.95)
Validation accuracy: 0.236
(1000, 300, 1e-05, 0.97)
Validation accuracy: 0.199
(1000, 300, 1e-05, 0.99)
Validation accuracy: 0.192
(3000, 100, 0.0005, 0.93)
Validation accuracy: 0.495
(3000, 100, 0.0005, 0.95)
```

Validation accuracy: 0.488  
(3000, 100, 0.0005, 0.97)  
Validation accuracy: 0.492  
(3000, 100, 0.0005, 0.99)  
Validation accuracy: 0.479  
(3000, 100, 0.0001, 0.93)  
Validation accuracy: 0.404  
(3000, 100, 0.0001, 0.95)  
Validation accuracy: 0.404  
(3000, 100, 0.0001, 0.97)  
Validation accuracy: 0.396  
(3000, 100, 0.0001, 0.99)  
Validation accuracy: 0.429  
(3000, 100, 1e-05, 0.93)  
Validation accuracy: 0.187  
(3000, 100, 1e-05, 0.95)  
Validation accuracy: 0.206  
(3000, 100, 1e-05, 0.97)  
Validation accuracy: 0.196  
(3000, 100, 1e-05, 0.99)  
Validation accuracy: 0.2  
(3000, 200, 0.0005, 0.93)  
Validation accuracy: 0.511  
(3000, 200, 0.0005, 0.95)  
Validation accuracy: 0.488  
(3000, 200, 0.0005, 0.97)  
Validation accuracy: 0.489  
(3000, 200, 0.0005, 0.99)  
Validation accuracy: 0.505  
(3000, 200, 0.0001, 0.93)  
Validation accuracy: 0.374  
(3000, 200, 0.0001, 0.95)  
Validation accuracy: 0.387  
(3000, 200, 0.0001, 0.97)  
Validation accuracy: 0.403  
(3000, 200, 0.0001, 0.99)  
Validation accuracy: 0.421  
(3000, 200, 1e-05, 0.93)  
Validation accuracy: 0.192  
(3000, 200, 1e-05, 0.95)  
Validation accuracy: 0.163  
(3000, 200, 1e-05, 0.97)  
Validation accuracy: 0.185  
(3000, 200, 1e-05, 0.99)  
Validation accuracy: 0.197  
(3000, 300, 0.0005, 0.93)  
Validation accuracy: 0.494  
(3000, 300, 0.0005, 0.95)  
Validation accuracy: 0.501  
(3000, 300, 0.0005, 0.97)  
Validation accuracy: 0.507  
(3000, 300, 0.0005, 0.99)  
Validation accuracy: 0.49  
(3000, 300, 0.0001, 0.93)  
Validation accuracy: 0.414  
(3000, 300, 0.0001, 0.95)  
Validation accuracy: 0.417  
(3000, 300, 0.0001, 0.97)  
Validation accuracy: 0.42  
(3000, 300, 0.0001, 0.99)  
Validation accuracy: 0.426  
(3000, 300, 1e-05, 0.93)  
Validation accuracy: 0.195  
(3000, 300, 1e-05, 0.95)  
Validation accuracy: 0.179  
(3000, 300, 1e-05, 0.97)  
Validation accuracy: 0.194  
(3000, 300, 1e-05, 0.99)  
Validation accuracy: 0.199

```
(5000, 100, 0.0005, 0.93)
Validation accuracy: 0.517
(5000, 100, 0.0005, 0.95)
Validation accuracy: 0.496
(5000, 100, 0.0005, 0.97)
Validation accuracy: 0.511
(5000, 100, 0.0005, 0.99)
Validation accuracy: 0.504
(5000, 100, 0.0001, 0.93)
Validation accuracy: 0.446
(5000, 100, 0.0001, 0.95)
Validation accuracy: 0.44
(5000, 100, 0.0001, 0.97)
Validation accuracy: 0.453
(5000, 100, 0.0001, 0.99)
Validation accuracy: 0.452
(5000, 100, 1e-05, 0.93)
Validation accuracy: 0.198
(5000, 100, 1e-05, 0.95)
Validation accuracy: 0.206
(5000, 100, 1e-05, 0.97)
Validation accuracy: 0.205
(5000, 100, 1e-05, 0.99)
Validation accuracy: 0.233
(5000, 200, 0.0005, 0.93)
Validation accuracy: 0.503
(5000, 200, 0.0005, 0.95)
Validation accuracy: 0.51
(5000, 200, 0.0005, 0.97)
Validation accuracy: 0.507
(5000, 200, 0.0005, 0.99)
Validation accuracy: 0.521
(5000, 200, 0.0001, 0.93)
Validation accuracy: 0.406
(5000, 200, 0.0001, 0.95)
Validation accuracy: 0.43
(5000, 200, 0.0001, 0.97)
Validation accuracy: 0.459
(5000, 200, 0.0001, 0.99)
Validation accuracy: 0.452
(5000, 200, 1e-05, 0.93)
Validation accuracy: 0.198
(5000, 200, 1e-05, 0.95)
Validation accuracy: 0.17
(5000, 200, 1e-05, 0.97)
Validation accuracy: 0.186
(5000, 200, 1e-05, 0.99)
Validation accuracy: 0.22
(5000, 300, 0.0005, 0.93)
Validation accuracy: 0.521
(5000, 300, 0.0005, 0.95)
Validation accuracy: 0.524
(5000, 300, 0.0005, 0.97)
Validation accuracy: 0.511
(5000, 300, 0.0005, 0.99)
Validation accuracy: 0.524
(5000, 300, 0.0001, 0.93)
Validation accuracy: 0.451
(5000, 300, 0.0001, 0.95)
Validation accuracy: 0.457
(5000, 300, 0.0001, 0.97)
Validation accuracy: 0.454
(5000, 300, 0.0001, 0.99)
Validation accuracy: 0.458
(5000, 300, 1e-05, 0.93)
Validation accuracy: 0.212
(5000, 300, 1e-05, 0.95)
Validation accuracy: 0.22
(5000, 300, 1e-05, 0.97)
```

Validation accuracy: 0.228  
 (5000, 300, 1e-05, 0.99)  
 Validation accuracy: 0.233  
 Validation accuracy: 0.524

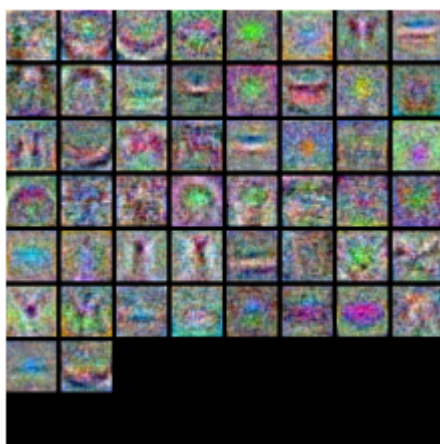
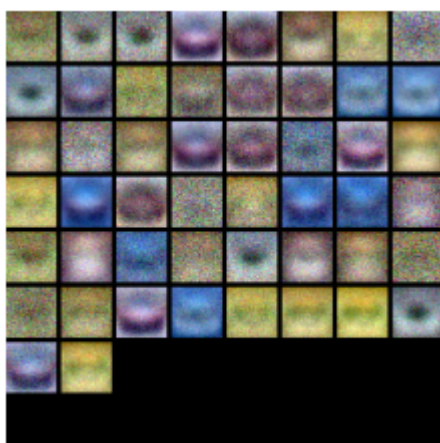
In [279]...

```
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) The suboptimal network has visual features which are very difficult to distinguish and smooth whereas we can even observe onshapes in the best net. The best net has more features to distinguish the images

## Evaluate on test set

In [280...

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

Test accuracy: 0.502

In [ ]:



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         self.params = {}
37         self.params['W1'] = std * np.random.randn(hidden_size, input_size)
38         self.params['b1'] = np.zeros(hidden_size)
39         self.params['W2'] = std * np.random.randn(output_size, hidden_size)
40         self.params['b2'] = np.zeros(output_size)
41
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 parameter
64       with respect to the loss function; has the same keys as self.params.
65     """
66     # Unpack variables from the params dictionary
67     W1, b1 = self.params['W1'], self.params['b1']
68     W2, b2 = self.params['W2'], self.params['b2']
69     N, D = X.shape
70
71     # Compute the forward pass
72     scores = None
73
74     # ===== #
75     # YOUR CODE HERE:
76     # Calculate the output scores of the neural network. The result
77     # should be (N, C). As stated in the description for this class,
78     # there should not be a ReLU layer after the second FC layer.
79     # The output of the second FC layer is the output scores. Do not
80     # use a for loop in your implementation.
81     # ===== #
82
83     relu = lambda x : x * (x>0)
84
85     r = X @ W1.T + b1
86     h1 = relu(r)
87     scores = h1 @ W2.T + b2
88
89     # ===== #
90     # END YOUR CODE HERE
91     # ===== #
92
93
94     # If the targets are not given then jump out, we're done
95     if y is None:
96         return scores
97
98     # Compute the loss
99     loss = None
100
101     # ===== #
102     # YOUR CODE HERE:
103     # Calculate the loss of the neural network. This includes the
104     # softmax loss and the L2 regularization for W1 and W2. Store the
105     # total loss in the variable loss. Multiply the regularization
106     # loss by 0.5 (in addition to the factor reg).
107     # ===== #
108
109     # scores is num_examples by num_classes
110
111     # max_score = np.max(scores, axis=1, keepdims=True)
112     e_x = np.exp(scores)
113     e_x_sum = np.sum(e_x, axis=1, keepdims=True)
114     softmax = e_x / e_x_sum
115     loss = -np.log(np.choose(y, softmax.T))
116     loss = np.mean(loss)
117
118     # adding L2 regularization loss
119     loss += 0.5*reg*((np.sum(W1**2)) + np.sum(W2**2))
120

```

```

121
122 # ===== #
123 # END YOUR CODE HERE
124 # ===== #
125
126 grads = {}
127
128 # ===== #
129 # YOUR CODE HERE:
130 # Implement the backward pass. Compute the derivatives of the
131 # weights and the biases. Store the results in the grads
132 # dictionary. e.g., grads['W1'] should store the gradient for
133 # W1, and be of the same size as W1.
134 # ===== #
135
136 indicator = lambda x : 1 * (x>0)
137
138
139 softmax[np.arange(X.shape[0]),y] -= 1
140 softmax = softmax / X.shape[0]
141
142 grads['b2'] = np.sum(softmax, axis=0)
143 grads['W2'] = softmax.T @ h1 + reg * W2
144
145
146 grad_layer1 = np.multiply(indicator(r), softmax @ W2)
147 grads['b1'] = np.sum(grad_layer1, axis=0)
148 grads['W1'] = grad_layer1.T @ X + reg * W1
149
150
151
152
153 # ===== #
154 # END YOUR CODE HERE
155 # ===== #
156
157 return loss, grads
158
159 def train(self, X, y, X_val, y_val,
160           learning_rate=1e-3, learning_rate_decay=0.95,
161           reg=1e-5, num_iters=100,
162           batch_size=200, verbose=False):
163     """
164     Train this neural network using stochastic gradient descent.
165
166     Inputs:
167     - X: A numpy array of shape (N, D) giving training data.
168     - y: A numpy array of shape (N,) giving training labels; y[i] = c means that
169       X[i] has label c, where 0 <= c < C.
170     - X_val: A numpy array of shape (N_val, D) giving validation data.
171     - y_val: A numpy array of shape (N_val,) giving validation labels.
172     - learning_rate: Scalar giving learning rate for optimization.
173     - learning_rate_decay: Scalar giving factor used to decay the learning rate
174       after each epoch.
175     - reg: Scalar giving regularization strength.
176     - num_iters: Number of steps to take when optimizing.
177     - batch_size: Number of training examples to use per step.
178     - verbose: boolean; if true print progress during optimization.
179     """
180     num_train = X.shape[0]
181     iterations_per_epoch = max(num_train / batch_size, 1)

```

```

182
183     # Use SGD to optimize the parameters in self.model
184     loss_history = []
185     train_acc_history = []
186     val_acc_history = []
187
188     for it in np.arange(num_iters):
189         X_batch = None
190         y_batch = None
191
192         # ===== #
193         # YOUR CODE HERE:
194         # Create a minibatch by sampling batch_size samples randomly.
195         # ===== #
196         idx = np.random.randint(X.shape[0], size=batch_size)
197         X_batch = X[idx, :]
198         y_batch = y[idx]
199
200         # ===== #
201         # END YOUR CODE HERE
202         # ===== #
203
204         # Compute loss and gradients using the current minibatch
205         loss, grads = self.loss(X_batch, y=y_batch, reg=reg)
206         loss_history.append(loss)
207
208         # ===== #
209         # YOUR CODE HERE:
210         # Perform a gradient descent step using the minibatch to update
211         # all parameters (i.e., W1, W2, b1, and b2).
212         # ===== #
213
214         self.params['W1'] = self.params['W1'] - learning_rate * grads['W1']
215         self.params['b1'] = self.params['b1'] - learning_rate * grads['b1']
216         self.params['W2'] = self.params['W2'] - learning_rate * grads['W2']
217         self.params['b2'] = self.params['b2'] - learning_rate * grads['b2']
218
219         # ===== #
220         # END YOUR CODE HERE
221         # ===== #
222
223         if verbose and it % 100 == 0:
224             print('iteration {} / {}: loss {}'.format(it, num_iters, loss))
225
226         # Every epoch, check train and val accuracy and decay learning rate.
227         if it % iterations_per_epoch == 0:
228             # Check accuracy
229             train_acc = (self.predict(X_batch) == y_batch).mean()
230             val_acc = (self.predict(X_val) == y_val).mean()
231             train_acc_history.append(train_acc)
232             val_acc_history.append(val_acc)
233
234             # Decay learning rate
235             learning_rate *= learning_rate_decay
236
237     return {
238         'loss_history': loss_history,
239         'train_acc_history': train_acc_history,
240         'val_acc_history': val_acc_history,
241     }
242

```

```

243 def predict(self, X):
244     """
245     Use the trained weights of this two-layer network to predict labels for
246     data points. For each data point we predict scores for each of the C
247     classes, and assign each data point to the class with the highest score.
248
249     Inputs:
250     - X: A numpy array of shape (N, D) giving N D-dimensional data points to
251         classify.
252
253     Returns:
254     - y_pred: A numpy array of shape (N,) giving predicted labels for each of
255         the elements of X. For all i, y_pred[i] = c means that X[i] is predicted
256         to have class c, where 0 <= c < C.
257     """
258     y_pred = None
259
260     # ===== #
261     # YOUR CODE HERE:
262     # Predict the class given the input data.
263     # ===== #
264
265     scores = self.loss(X)
266     max_score = np.max(scores, axis=1, keepdims=True)
267     e_x = np.exp(scores - max_score)
268     e_x_sum = np.sum(e_x, axis=1, keepdims=True)
269     softmax = e_x / e_x_sum
270     y_pred = np.argmax(softmax, axis=1)
271
272
273     # ===== #
274     # END YOUR CODE HERE
275     # ===== #
276
277     return y_pred
278

```

In [ ]:

1

# 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 [69]: ## 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_grad
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-ipy
%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 [70]: # 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 [71]: # 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: 9.7698500479884e-10

```

## Affine layer backward pass

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

In [72]:

```

# 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: 2.382650916371664e-10
dw error: 3.101548428922205e-10
db error: 1.4248231174136619e-11

```

## Activation layers

In this section you'll implement the ReLU activation.

### ReLU forward pass

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



```
In [73]: # 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: 4.999999798022158e-08
```

## ReLU backward pass

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

```
In [74]: 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: 3.275602043584319e-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 [75]: 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)
dw_num = eval_numerical_gradient_array(lambda w: affine_relu_forward(x, w, b)
```

```
db_num = eval_numerical_gradient_array(lambda b: affine_relu_forward(x, w, b)

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: 3.1196798386884486e-10
dw error: 1.3444421240797957e-10
db error: 6.449843882269766e-11
```

## Softmax losses

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

In [76]:

```
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=1,
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: 2.3019611530262
dx error: 8.651759089870189e-09
```

## Implementation of a two-layer NN

In `nn/nnl/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 [77]:

```
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, 15.3
      [12.05769098, 12.74614105, 13.43459113, 14.1230412, 14.81149128, 15.4
      [12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.05099822, 15.6
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: 1.2236151215593397e-08
W2 relative error: 3.3429539606923665e-10
b1 relative error: 4.7288944058018464e-09
b2 relative error: 4.3291285233961314e-10
Running numeric gradient check with reg = 0.7
W1 relative error: 2.527915286171985e-07
W2 relative error: 1.3678335722105113e-07
b1 relative error: 1.5646801749611563e-08
b2 relative error: 9.089621155678095e-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 [79]:

```

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.

```

```

#
# ===== #

weight_scale = 1e-3
learning_rate = 5e-4

model = TwoLayerNet(input_dim=3*32*32, hidden_dims=200, num_classes=10, weight_scale=weight_scale)
solver = Solver(model, data, print_every=100, num_epochs=20, batch_size=200,
                 update_rule='sgd',
                 optim_config={
                     'learning_rate': learning_rate,
                 })
solver.train()

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

```

```

(Iteration 1 / 4900) loss: 2.296319
(Epoch 0 / 20) train acc: 0.147000; val_acc: 0.149000
(Iteration 101 / 4900) loss: 1.906182
(Iteration 201 / 4900) loss: 1.825798
(Epoch 1 / 20) train acc: 0.386000; val_acc: 0.408000
(Iteration 301 / 4900) loss: 1.731690
(Iteration 401 / 4900) loss: 1.607905
(Epoch 2 / 20) train acc: 0.428000; val_acc: 0.449000
(Iteration 501 / 4900) loss: 1.547835
(Iteration 601 / 4900) loss: 1.521764
(Iteration 701 / 4900) loss: 1.493333
(Epoch 3 / 20) train acc: 0.469000; val_acc: 0.459000
(Iteration 801 / 4900) loss: 1.476093
(Iteration 901 / 4900) loss: 1.439115
(Epoch 4 / 20) train acc: 0.494000; val_acc: 0.489000
(Iteration 1001 / 4900) loss: 1.413100
(Iteration 1101 / 4900) loss: 1.386595
(Iteration 1201 / 4900) loss: 1.400718
(Epoch 5 / 20) train acc: 0.532000; val_acc: 0.483000
(Iteration 1301 / 4900) loss: 1.291695
(Iteration 1401 / 4900) loss: 1.354334
(Epoch 6 / 20) train acc: 0.561000; val_acc: 0.507000
(Iteration 1501 / 4900) loss: 1.405150
(Iteration 1601 / 4900) loss: 1.211215
(Iteration 1701 / 4900) loss: 1.305215
(Epoch 7 / 20) train acc: 0.545000; val_acc: 0.499000
(Iteration 1801 / 4900) loss: 1.289038
(Iteration 1901 / 4900) loss: 1.259552
(Epoch 8 / 20) train acc: 0.509000; val_acc: 0.533000
(Iteration 2001 / 4900) loss: 1.235271
(Iteration 2101 / 4900) loss: 1.293656
(Iteration 2201 / 4900) loss: 1.213470
(Epoch 9 / 20) train acc: 0.584000; val_acc: 0.520000
(Iteration 2301 / 4900) loss: 1.294880
(Iteration 2401 / 4900) loss: 1.101703
(Epoch 10 / 20) train acc: 0.580000; val_acc: 0.520000
(Iteration 2501 / 4900) loss: 1.250590
(Iteration 2601 / 4900) loss: 1.198588
(Epoch 11 / 20) train acc: 0.582000; val_acc: 0.530000
(Iteration 2701 / 4900) loss: 1.181495
(Iteration 2801 / 4900) loss: 1.299144
(Iteration 2901 / 4900) loss: 1.033070
(Epoch 12 / 20) train acc: 0.577000; val_acc: 0.541000
(Iteration 3001 / 4900) loss: 1.138828
(Iteration 3101 / 4900) loss: 1.046784
(Epoch 13 / 20) train acc: 0.607000; val_acc: 0.529000
(Iteration 3201 / 4900) loss: 1.211506
(Iteration 3301 / 4900) loss: 1.250867

```

```

(Iteration 3401 / 4900) loss: 1.133314
(Epoch 14 / 20) train acc: 0.616000; val_acc: 0.536000
(Iteration 3501 / 4900) loss: 1.121289
(Iteration 3601 / 4900) loss: 1.217634
(Epoch 15 / 20) train acc: 0.611000; val_acc: 0.519000
(Iteration 3701 / 4900) loss: 0.985237
(Iteration 3801 / 4900) loss: 1.038370
(Iteration 3901 / 4900) loss: 0.999046
(Epoch 16 / 20) train acc: 0.642000; val_acc: 0.511000
(Iteration 4001 / 4900) loss: 0.996137
(Iteration 4101 / 4900) loss: 1.139573
(Epoch 17 / 20) train acc: 0.629000; val_acc: 0.524000
(Iteration 4201 / 4900) loss: 0.969272
(Iteration 4301 / 4900) loss: 1.077427
(Iteration 4401 / 4900) loss: 0.977316
(Epoch 18 / 20) train acc: 0.670000; val_acc: 0.537000
(Iteration 4501 / 4900) loss: 1.120773
(Iteration 4601 / 4900) loss: 0.966483
(Epoch 19 / 20) train acc: 0.654000; val_acc: 0.531000
(Iteration 4701 / 4900) loss: 1.156264
(Iteration 4801 / 4900) loss: 1.114885
(Epoch 20 / 20) train acc: 0.699000; val_acc: 0.533000

```

In [80]:

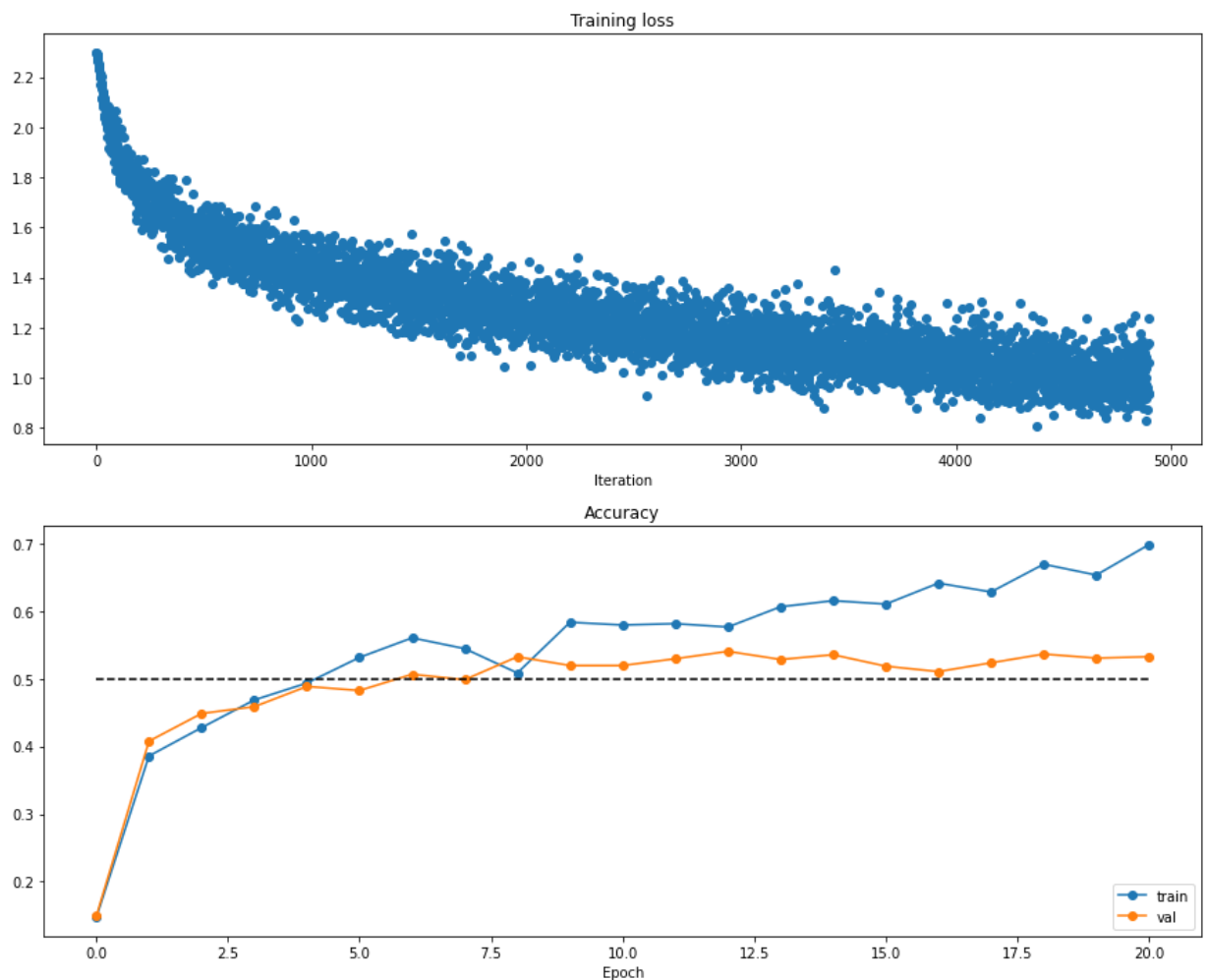
```
# Run this cell to visualize training loss and train / val accuracy
```

```

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 [81]: 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,
        print('{} relative error: {}'.format(name, rel_error(grad_num, grads[name]

Running check with reg = 0
Initial loss: 2.300773749401901
W1 relative error: 2.1340236665142505e-06
```

```

W2 relative error: 2.2743309431166997e-07
W3 relative error: 6.197934169904083e-07
b1 relative error: 2.3588216651711635e-08
b2 relative error: 3.9595855196785e-09
b3 relative error: 1.1959783484217903e-10
Running check with reg = 3.14
Initial loss: 7.2065399362318505
W1 relative error: 9.086197491210612e-08
W2 relative error: 3.921022766550705e-08
W3 relative error: 1.6464860687585295e-08
b1 relative error: 5.370112431068358e-09
b2 relative error: 5.1965296229205134e-08
b3 relative error: 1.335207956377867e-10

```

In [85]:

```

# 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()

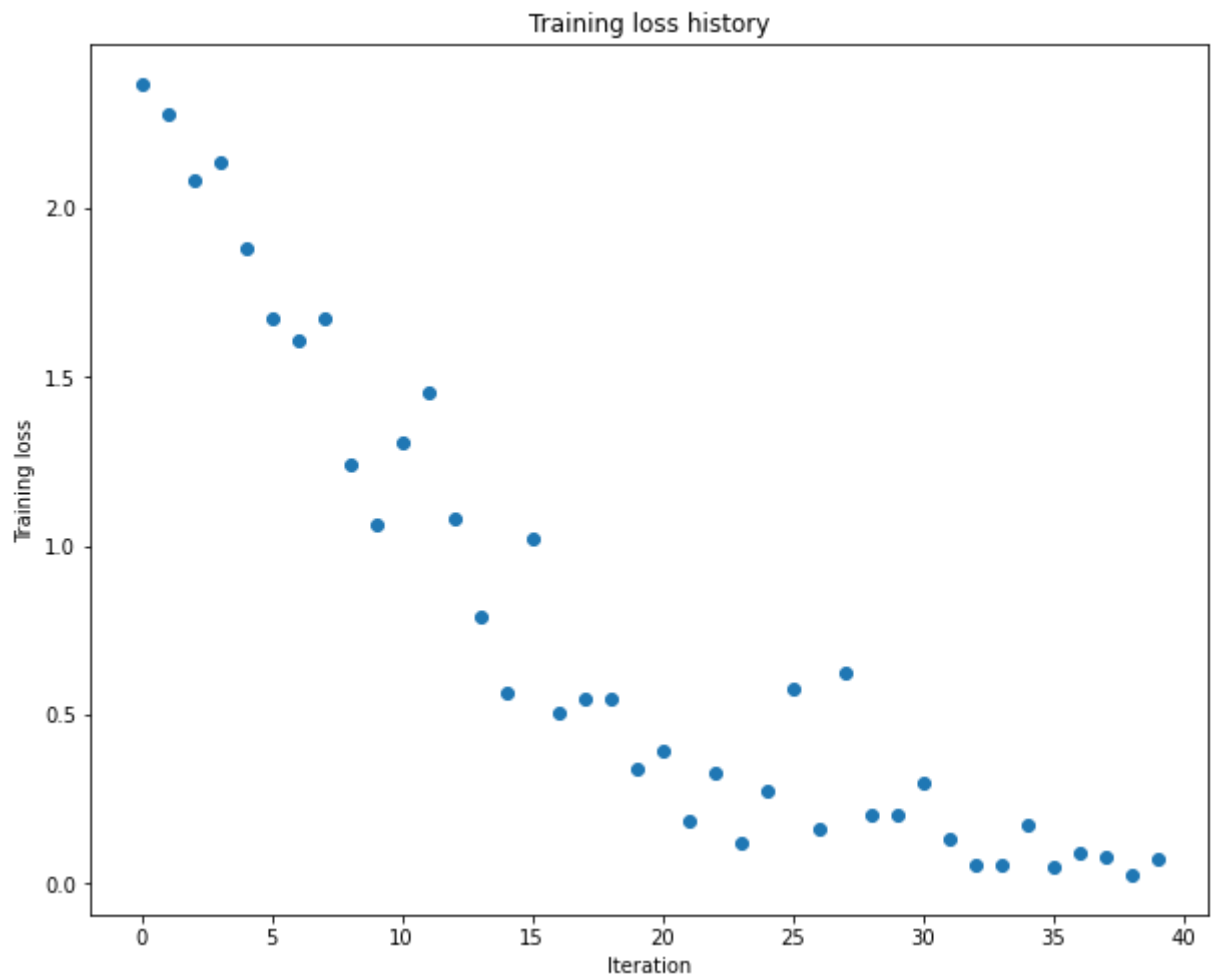
```

```

(Iteration 1 / 40) loss: 2.369362
(Epoch 0 / 20) train acc: 0.260000; val_acc: 0.111000
(Epoch 1 / 20) train acc: 0.340000; val_acc: 0.140000
(Epoch 2 / 20) train acc: 0.280000; val_acc: 0.082000
(Epoch 3 / 20) train acc: 0.420000; val_acc: 0.142000
(Epoch 4 / 20) train acc: 0.680000; val_acc: 0.154000
(Epoch 5 / 20) train acc: 0.680000; val_acc: 0.178000
(Iteration 11 / 40) loss: 1.304191
(Epoch 6 / 20) train acc: 0.660000; val_acc: 0.170000
(Epoch 7 / 20) train acc: 0.820000; val_acc: 0.193000
(Epoch 8 / 20) train acc: 0.900000; val_acc: 0.172000
(Epoch 9 / 20) train acc: 0.920000; val_acc: 0.200000
(Epoch 10 / 20) train acc: 0.920000; val_acc: 0.197000
(Iteration 21 / 40) loss: 0.395965
(Epoch 11 / 20) train acc: 0.940000; val_acc: 0.206000
(Epoch 12 / 20) train acc: 0.960000; val_acc: 0.205000
(Epoch 13 / 20) train acc: 0.880000; val_acc: 0.173000
(Epoch 14 / 20) train acc: 0.900000; val_acc: 0.199000

```

```
(Epoch 15 / 20) train acc: 0.960000; val_acc: 0.202000
(Iteration 31 / 40) loss: 0.299047
(Epoch 16 / 20) train acc: 0.980000; val_acc: 0.205000
(Epoch 17 / 20) train acc: 0.960000; val_acc: 0.201000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.200000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.196000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.190000
```

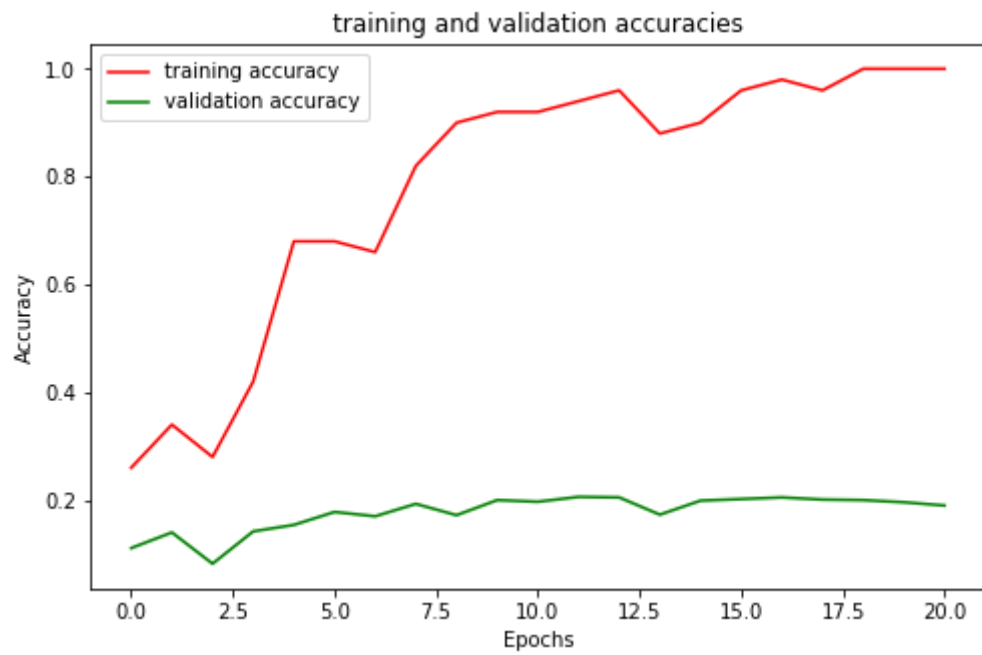


In [86]:

```
plt.figure(figsize=(8,5))
plt.title("training and validation accuracies")
# plt.xticks(epochs)
plt.plot( solver.train_acc_history, label="training accuracy", color='r')
plt.plot( solver.val_acc_history, label="validation accuracy", color='g')

plt.xlabel("Epochs")
plt.ylabel("Accuracy")
plt.legend()
plt.show()
```





In [ ]:

In [ ]:

```

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

```

```

60  # YOUR CODE HERE:
61  #   Calculate the gradients for the backward pass.
62  # ===== #
63
64  # dout is N x M
65  # dx should be N x d1 x ... x dk; it relates to dout through multiplication w
66  # dw should be D x M; it relates to dout through multiplication with x, which
67  # db should be M; it is just the sum over dout examples
68
69  dx = (dout @ w.T).reshape(x.shape)
70  dw = x.reshape(x.shape[0], -1).T @ dout
71  db = np.sum(dout, axis = 0)
72
73  # ===== #
74  # END YOUR CODE HERE
75  # ===== #
76
77  return dx, dw, db
78
79 def relu_forward(x):
80     """
81     Computes the forward pass for a layer of rectified linear units (ReLU).
82
83     Input:
84     - x: Inputs, of any shape
85
86     Returns a tuple of:
87     - out: Output, of the same shape as x
88     - cache: x
89     """
90     # ===== #
91     # YOUR CODE HERE:
92     #   Implement the ReLU forward pass.
93     # ===== #
94
95     out = x * (x>0)
96
97     # ===== #
98     # END YOUR CODE HERE
99     # ===== #
100
101     cache = x
102     return out, cache
103
104
105 def relu_backward(dout, cache):
106     """
107     Computes the backward pass for a layer of rectified linear units (ReLU).
108
109     Input:
110     - dout: Upstream derivatives, of any shape
111     - cache: Input x, of same shape as dout
112
113     Returns:
114     - dx: Gradient with respect to x
115     """
116     x = cache
117
118     # ===== #
119     # YOUR CODE HERE:
120     #   Implement the ReLU backward pass

```

```

121 # ===== #
122
123 # ReLU directs linearly to those > 0
124 dx = np.multiply(dout, 1 * (x>0))
125
126 # ===== #
127 # END YOUR CODE HERE
128 # ===== #
129
130 return dx
131
132 def svm_loss(x, y):
133     """
134     Computes the loss and gradient using for multiclass SVM classification.
135
136     Inputs:
137     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
138       for the ith input.
139     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
140       0 <= y[i] < C
141
142     Returns a tuple of:
143     - loss: Scalar giving the loss
144     - dx: Gradient of the loss with respect to x
145     """
146     N = x.shape[0]
147     correct_class_scores = x[np.arange(N), y]
148     margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
149     margins[np.arange(N), y] = 0
150     loss = np.sum(margins) / N
151     num_pos = np.sum(margins > 0, axis=1)
152     dx = np.zeros_like(x)
153     dx[margins > 0] = 1
154     dx[np.arange(N), y] -= num_pos
155     dx /= N
156     return loss, dx
157
158
159 def softmax_loss(x, y):
160     """
161     Computes the loss and gradient for softmax classification.
162
163     Inputs:
164     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
165       for the ith input.
166     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
167       0 <= y[i] < C
168
169     Returns a tuple of:
170     - loss: Scalar giving the loss
171     - dx: Gradient of the loss with respect to x
172     """
173
174     probs = np.exp(x - np.max(x, axis=1, keepdims=True))
175     probs /= np.sum(probs, axis=1, keepdims=True)
176     N = x.shape[0]
177     loss = -np.sum(np.log(probs[np.arange(N), y])) / N
178     dx = probs.copy()
179     dx[np.arange(N), y] -= 1
180     dx /= N
181     return loss, dx

```



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
50         self.params = {}
51         self.params['W1'] = np.random.normal(loc=0, scale=weight_scale, size=(input_dim, hidden_dims))
52         self.params['b1'] = np.zeros(hidden_dims)
53         self.params['W2'] = np.random.normal(loc=0, scale=weight_scale, size=(hidden_dims, num_classes))
54         self.params['b2'] = np.zeros(num_classes)
55
56         # ===== #
57         # END YOUR CODE HERE
58         # ===== #
59

```

```

60 def loss(self, X, y=None):
61     """
62     Compute loss and gradient for a minibatch of data.
63
64     Inputs:
65     - X: Array of input data of shape (N, d_1, ..., d_k)
66     - y: Array of labels, of shape (N,). y[i] gives the label for X[i].
67
68     Returns:
69     If y is None, then run a test-time forward pass of the model and return:
70     - scores: Array of shape (N, C) giving classification scores, where
71       scores[i, c] is the classification score for X[i] and class c.
72
73     If y is not None, then run a training-time forward and backward pass and
74     return a tuple of:
75     - loss: Scalar value giving the loss
76     - grads: Dictionary with the same keys as self.params, mapping parameter
77       names to gradients of the loss with respect to those parameters.
78     """
79     scores = None
80
81     # ===== #
82     # YOUR CODE HERE:
83     # Implement the forward pass of the two-layer neural network. Store
84     # the class scores as the variable 'scores'. Be sure to use the layers
85     # you prior implemented.
86     # ===== #
87
88     W1, b1 = self.params['W1'], self.params['b1']
89     W2, b2 = self.params['W2'], self.params['b2']
90
91     h, cache_h = affine_relu_forward(X, W1, b1);
92     z, cache_z = affine_forward(h, W2, b2)
93     scores = z
94
95     # ===== #
96     # END YOUR CODE HERE
97     # ===== #
98
99     # If y is None then we are in test mode so just return scores
100    if y is None:
101        return scores
102
103    loss, grads = 0, {}
104    # ===== #
105    # YOUR CODE HERE:
106    # Implement the backward pass of the two-layer neural net. Store
107    # the loss as the variable 'loss' and store the gradients in the
108    # 'grads' dictionary. For the grads dictionary, grads['W1'] holds
109    # the gradient for W1, grads['b1'] holds the gradient for b1, etc.
110    # i.e., grads[k] holds the gradient for self.params[k].
111    #
112    # Add L2 regularization, where there is an added cost 0.5*self.reg*W^2
113    # for each W. Be sure to include the 0.5 multiplying factor to
114    # match our implementation.
115    #
116    # And be sure to use the layers you prior implemented.
117    # ===== #
118
119    loss, dz = softmax_loss(scores, y)
120    loss += 0.5*self.reg*(np.sum(W1**2) + np.sum(W2**2))

```

```

121
122
123     dh, dw2, db2 = affine_backward(dz, cache_z)
124     dx, dw1, db1 = affine_relu_backward(dh, cache_h)
125
126     grads['b1'] = db1
127     grads['W1'] = dw1 + self.reg * W1
128     grads['b2'] = db2
129     grads['W2'] = dw2 + self.reg * W2
130
131
132
133     # ===== #
134     # END YOUR CODE HERE
135     # ===== #
136
137     return loss, grads
138
139
140 class FullyConnectedNet(object):
141     """
142     A fully-connected neural network with an arbitrary number of hidden layers,
143     ReLU nonlinearities, and a softmax loss function. This will also implement
144     dropout and batch normalization as options. For a network with L layers,
145     the architecture will be
146
147     {affine - [batch norm] - relu - [dropout]} x (L - 1) - affine - softmax
148
149     where batch normalization and dropout are optional, and the {...} block is
150     repeated L - 1 times.
151
152     Similar to the TwoLayerNet above, learnable parameters are stored in the
153     self.params dictionary and will be learned using the Solver class.
154     """
155
156     def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,
157                 dropout=0, use_batchnorm=False, reg=0.0,
158                 weight_scale=1e-2, dtype=np.float32, seed=None):
159         """
160         Initialize a new FullyConnectedNet.
161
162         Inputs:
163         - hidden_dims: A list of integers giving the size of each hidden layer.
164         - input_dim: An integer giving the size of the input.
165         - num_classes: An integer giving the number of classes to classify.
166         - dropout: Scalar between 0 and 1 giving dropout strength. If dropout=0 the
167           the network should not use dropout at all.
168         - use_batchnorm: Whether or not the network should use batch normalization.
169         - reg: Scalar giving L2 regularization strength.
170         - weight_scale: Scalar giving the standard deviation for random
171           initialization of the weights.
172         - dtype: A numpy datatype object; all computations will be performed using
173           this datatype. float32 is faster but less accurate, so you should use
174           float64 for numeric gradient checking.
175         - seed: If not None, then pass this random seed to the dropout layers. This
176           will make the dropout layers deterministic so we can gradient check the
177           model.
178         """
179         self.use_batchnorm = use_batchnorm
180         self.use_dropout = dropout > 0
181         self.reg = reg

```



```

182 self.num_layers = 1 + len(hidden_dims)
183 self.dtype = dtype
184 self.params = {}
185
186 # ===== #
187 # YOUR CODE HERE:
188 # Initialize all parameters of the network in the self.params dictionary.
189 # The weights and biases of layer 1 are W1 and b1; and in general the
190 # weights and biases of layer i are Wi and bi. The
191 # biases are initialized to zero and the weights are initialized
192 # so that each parameter has mean 0 and standard deviation weight_scale.
193 # ===== #
194
195 for i in range(self.num_layers):
196     W_name = 'W'+str(i+1)
197     b_name = 'b'+str(i+1)
198     if (i == 0):
199         self.params[W_name] = np.random.normal(loc=0, scale=weight_scale, size=
200         self.params[b_name] = np.zeros(hidden_dims[i])
201     elif (i == self.num_layers-1):
202         self.params[W_name] = np.random.normal(loc=0, scale=weight_scale, size=
203         self.params[b_name] = np.zeros(num_classes)
204     else:
205         self.params[W_name] = np.random.normal(loc=0, scale=weight_scale, size=
206         self.params[b_name] = np.zeros(hidden_dims[i])
207
208 # ===== #
209 # END YOUR CODE HERE
210 # ===== #
211
212 # When using dropout we need to pass a dropout_param dictionary to each
213 # dropout layer so that the layer knows the dropout probability and the mod
214 # (train / test). You can pass the same dropout_param to each dropout layer
215 self.dropout_param = {}
216 if self.use_dropout:
217     self.dropout_param = {'mode': 'train', 'p': dropout}
218     if seed is not None:
219         self.dropout_param['seed'] = seed
220
221 # With batch normalization we need to keep track of running means and
222 # variances, so we need to pass a special bn_param object to each batch
223 # normalization layer. You should pass self.bn_params[0] to the forward pas
224 # of the first batch normalization layer, self.bn_params[1] to the forward
225 # pass of the second batch normalization layer, etc.
226 self.bn_params = []
227 if self.use_batchnorm:
228     self.bn_params = [{'mode': 'train'} for i in np.arange(self.num_layers -
229
230 # Cast all parameters to the correct datatype
231 for k, v in self.params.items():
232     self.params[k] = v.astype(dtype)
233
234
235 def loss(self, X, y=None):
236     """
237     Compute loss and gradient for the fully-connected net.
238
239     Input / output: Same as TwoLayerNet above.
240     """
241     X = X.astype(self.dtype)
242     mode = 'test' if y is None else 'train'

```

```

243
244     # Set train/test mode for batchnorm params and dropout param since they
245     # behave differently during training and testing.
246     if self.dropout_param is not None:
247         self.dropout_param['mode'] = mode
248     if self.use_batchnorm:
249         for bn_param in self.bn_params:
250             bn_param[mode] = mode
251
252     scores = None
253
254     # ===== #
255     # YOUR CODE HERE:
256     # Implement the forward pass of the FC net and store the output
257     # scores as the variable "scores".
258     # ===== #
259
260     h_list = []
261     cache_list = []
262
263     for i in range(self.num_layers):
264         W_name = 'W'+str(i+1)
265         b_name = 'b'+str(i+1)
266         if (i == 0):
267             h, cache = affine_relu_forward(X, self.params[W_name], self.params[b_name])
268         elif (i == self.num_layers-1):
269             h, cache = affine_forward(h, self.params[W_name], self.params[b_name])
270         else:
271             h, cache = affine_relu_forward(h, self.params[W_name], self.params[b_name])
272         h_list.append(h)
273         cache_list.append(cache)
274
275     scores = h_list[-1]
276
277     # ===== #
278     # END YOUR CODE HERE
279     # ===== #
280
281     # If test mode return early
282     if mode == 'test':
283         return scores
284
285     loss, grads = 0.0, {}
286     # ===== #
287     # YOUR CODE HERE:
288     # Implement the backwards pass of the FC net and store the gradients
289     # in the grads dict, so that grads[k] is the gradient of self.params[k]
290     # Be sure your L2 regularization includes a 0.5 factor.
291     # ===== #
292
293     loss, dz = softmax_loss(scores, y)
294
295     for i in range(self.num_layers-1, -1, -1):
296         W_name = 'W'+str(i+1)
297         b_name = 'b'+str(i+1)
298
299         if (i == self.num_layers-1):
300             dh, dw, db = affine_backward(dz, cache_list[i])
301         else:
302             dh, dw, db = affine_relu_backward(dh, cache_list[i])
303

```

```
304     grads[W_name] = dw + self.reg * self.params[W_name]
305     grads[b_name] = db
306
307     loss += 0.5*self.reg*(np.sum(self.params[W_name]**2))
308
309     # ===== #
310     # END YOUR CODE HERE
311     # ===== #
312     return loss, grads
313
```

---