# Fully connected networks

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

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

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

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

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

# Modular layers

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

```python
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:

```python
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 [30]:
```python
## Import and setups

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

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

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

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

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

In [31]:
```python
# 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.

```python
In [32]:  # Test the affine_forward function

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

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

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

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

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

```
Testing affine_forward function:
difference: 9.769849468192957e-10
```

## Affine layer backward pass

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

In [33]:
```python
# 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: 2.0102720432483927e-10
dw error: 1.786902592840743e-10
db error: 5.168732824993997e-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 [34]:
```python
# 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 [35]:
```python
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: 3.27562144014014e-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 [36]: from nndl.layer_utils import affine_relu_forward, affine_relu_backward

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

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

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

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

```
Testing affine_relu_forward and affine_relu_backward:
dx error: 1.647876847124382e-10
dw error: 7.367935047019492e-10
db error: 6.506831444739067e-11
```

## Softmax and SVM losses

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

```
In [37]:  num_classes, num_inputs = 10, 50
          x = 0.001 * np.random.randn(num_inputs, num_classes)
          y = np.random.randint(num_classes, size=num_inputs)

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

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

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

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

```
Testing svm_loss:
loss: 9.000898228652648
dx error: 1.4021566006651672e-09

Testing softmax_loss:
loss: 2.302675382264634
dx error: 1.0603610683853635e-08
```

# Implementation of a two-layer NN

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

```python
In [38]: 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,  16.09215096],
            [12.05769098,  12.74614105,  13.43459113,  14.1230412,   14.81149128, 15.49994135,  16.18839143],
            [12.58373087,  13.20054771,  13.81736455,  14.43418138,  15.05099822, 15.66781506,  16.2846319 ]])
         scores_diff = np.abs(scores - correct_scores).sum()
         assert scores_diff < 1e-6, 'Problem with test-time forward pass'

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

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

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

           for name in sorted(grads):
             f = lambda _: model.loss(X, y)[0]
             grad_num = eval_numerical_gradient(f, model.params[name], verbose=False)
             print('{} relative error: {}'.format(name, rel_error(grad_num, 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: 1.8336562786695002e-08
W2 relative error: 3.201560569143183e-10
b1 relative error: 9.828315204644842e-09
b2 relative error: 4.329134954569865e-10
Running numeric gradient check with reg = 0.7
W1 relative error: 2.5279152310200606e-07
W2 relative error: 2.8508510893102143e-08
b1 relative error: 1.564679947504764e-08
b2 relative error: 9.089617896905665e-10
```

# Solver

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

In [39]:
```python
model = TwoLayerNet()
solver = None

# ================================================================ #
# YOUR CODE HERE:
#    Declare an instance of a TwoLayerNet and then train
#    it with the Solver. Choose hyperparameters so that your validation
#    accuracy is at least 50%.  We won't have you optimize this further
#    since you did it in the previous notebook.
#
# ================================================================ #

model = TwoLayerNet(input_dim = 3*32*32, hidden_dims = 200, num_classes = 10,
weight_scale = 1e-3)
solver = Solver(model, data, update_rule = 'sgd', optim_config = {'learning_ra
te': 0.0018889},
                lr_decay = 0.9125, num_epochs = 30, batch_size = 100, print_eve
ry = 100000)
solver.train()

# ================================================================ #
# END YOUR CODE HERE
# ================================================================ #
```
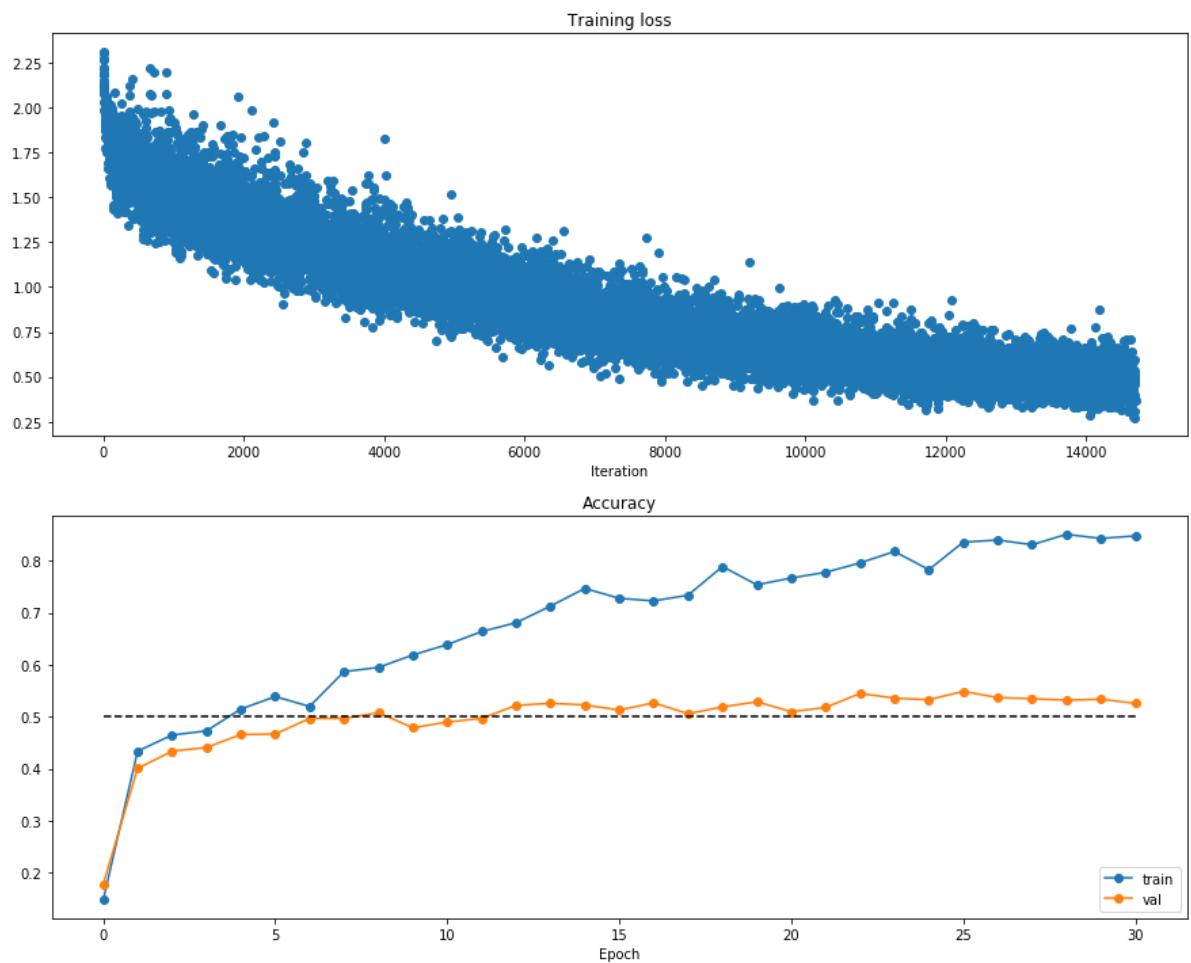
```
(Iteration 1 / 14700) loss: 2.309769
(Epoch 0 / 30) train acc: 0.148000; val_acc: 0.178000
(Epoch 1 / 30) train acc: 0.434000; val_acc: 0.401000
(Epoch 2 / 30) train acc: 0.465000; val_acc: 0.434000
(Epoch 3 / 30) train acc: 0.473000; val_acc: 0.441000
(Epoch 4 / 30) train acc: 0.515000; val_acc: 0.466000
(Epoch 5 / 30) train acc: 0.539000; val_acc: 0.467000
(Epoch 6 / 30) train acc: 0.520000; val_acc: 0.497000
(Epoch 7 / 30) train acc: 0.587000; val_acc: 0.497000
(Epoch 8 / 30) train acc: 0.595000; val_acc: 0.508000
(Epoch 9 / 30) train acc: 0.619000; val_acc: 0.479000
(Epoch 10 / 30) train acc: 0.639000; val_acc: 0.490000
(Epoch 11 / 30) train acc: 0.664000; val_acc: 0.497000
(Epoch 12 / 30) train acc: 0.681000; val_acc: 0.522000
(Epoch 13 / 30) train acc: 0.713000; val_acc: 0.526000
(Epoch 14 / 30) train acc: 0.747000; val_acc: 0.523000
(Epoch 15 / 30) train acc: 0.728000; val_acc: 0.513000
(Epoch 16 / 30) train acc: 0.723000; val_acc: 0.527000
(Epoch 17 / 30) train acc: 0.734000; val_acc: 0.506000
(Epoch 18 / 30) train acc: 0.789000; val_acc: 0.519000
(Epoch 19 / 30) train acc: 0.754000; val_acc: 0.529000
(Epoch 20 / 30) train acc: 0.767000; val_acc: 0.510000
(Epoch 21 / 30) train acc: 0.778000; val_acc: 0.518000
(Epoch 22 / 30) train acc: 0.796000; val_acc: 0.545000
(Epoch 23 / 30) train acc: 0.818000; val_acc: 0.536000
(Epoch 24 / 30) train acc: 0.783000; val_acc: 0.533000
(Epoch 25 / 30) train acc: 0.836000; val_acc: 0.549000
(Epoch 26 / 30) train acc: 0.840000; val_acc: 0.537000
(Epoch 27 / 30) train acc: 0.831000; val_acc: 0.535000
(Epoch 28 / 30) train acc: 0.851000; val_acc: 0.532000
(Epoch 29 / 30) train acc: 0.843000; val_acc: 0.534000
(Epoch 30 / 30) train acc: 0.848000; val_acc: 0.526000
```

In [40]:
```python
# Run this cell to visualize training loss and train / val accuracy

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

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

# Multilayer Neural Network

Now, we implement a multi-layer neural network.

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

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

```python
In [41]: 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: 3.3663538412813647
W1 relative error: 5.2888338248209637e-08
W2 relative error: 1.6791539477084933e-07
b1 relative error: 1.0889278181160834e-09
b2 relative error: 6.782527079089013e-10
Running check with reg = 3.14
Initial loss: 4.656960094658376
W1 relative error: 7.551965269664439e-08
W2 relative error: 4.981635484779869e-08
b1 relative error: 2.541723727607985e-08
b2 relative error: 1.2943627684305898e-09
```

In [50]:
```python
# 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 = 0.00004
learning_rate = 0.0010

model = FullyConnectedNet([100, 100],
                weight_scale=weight_scale, dtype=np.float64)
solver = Solver(model, small_data,
                print_every=10, num_epochs=200, batch_size=25,
                update_rule='sgd',
                optim_config={
                    'learning_rate': learning_rate,
                }
        )
solver.train()

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

```
(Iteration 1 / 400) loss: 4.605172
(Epoch 0 / 200) train acc: 0.220000; val_acc: 0.107000
(Epoch 1 / 200) train acc: 0.220000; val_acc: 0.092000
(Epoch 2 / 200) train acc: 0.280000; val_acc: 0.154000
(Epoch 3 / 200) train acc: 0.340000; val_acc: 0.140000
(Epoch 4 / 200) train acc: 0.220000; val_acc: 0.105000
(Epoch 5 / 200) train acc: 0.220000; val_acc: 0.100000
(Iteration 11 / 400) loss: 4.599741
(Epoch 6 / 200) train acc: 0.260000; val_acc: 0.122000
(Epoch 7 / 200) train acc: 0.200000; val_acc: 0.135000
(Epoch 8 / 200) train acc: 0.220000; val_acc: 0.141000
(Epoch 9 / 200) train acc: 0.220000; val_acc: 0.085000
(Epoch 10 / 200) train acc: 0.200000; val_acc: 0.116000
(Iteration 21 / 400) loss: 3.889044
(Epoch 11 / 200) train acc: 0.240000; val_acc: 0.088000
(Epoch 12 / 200) train acc: 0.280000; val_acc: 0.149000
(Epoch 13 / 200) train acc: 0.260000; val_acc: 0.128000
(Epoch 14 / 200) train acc: 0.220000; val_acc: 0.111000
(Epoch 15 / 200) train acc: 0.280000; val_acc: 0.154000
(Iteration 31 / 400) loss: 2.468571
(Epoch 16 / 200) train acc: 0.320000; val_acc: 0.141000
(Epoch 17 / 200) train acc: 0.340000; val_acc: 0.159000
(Epoch 18 / 200) train acc: 0.320000; val_acc: 0.158000
(Epoch 19 / 200) train acc: 0.380000; val_acc: 0.169000
(Epoch 20 / 200) train acc: 0.500000; val_acc: 0.160000
(Iteration 41 / 400) loss: 1.344916
(Epoch 21 / 200) train acc: 0.560000; val_acc: 0.166000
(Epoch 22 / 200) train acc: 0.660000; val_acc: 0.169000
(Epoch 23 / 200) train acc: 0.700000; val_acc: 0.184000
(Epoch 24 / 200) train acc: 0.700000; val_acc: 0.173000
(Epoch 25 / 200) train acc: 0.780000; val_acc: 0.181000
(Iteration 51 / 400) loss: 0.706996
(Epoch 26 / 200) train acc: 0.800000; val_acc: 0.194000
(Epoch 27 / 200) train acc: 0.840000; val_acc: 0.201000
(Epoch 28 / 200) train acc: 0.840000; val_acc: 0.175000
(Epoch 29 / 200) train acc: 0.840000; val_acc: 0.192000
(Epoch 30 / 200) train acc: 0.820000; val_acc: 0.171000
(Iteration 61 / 400) loss: 0.688226
(Epoch 31 / 200) train acc: 0.900000; val_acc: 0.176000
(Epoch 32 / 200) train acc: 0.900000; val_acc: 0.183000
(Epoch 33 / 200) train acc: 0.840000; val_acc: 0.162000
(Epoch 34 / 200) train acc: 0.900000; val_acc: 0.188000
(Epoch 35 / 200) train acc: 0.860000; val_acc: 0.163000
(Iteration 71 / 400) loss: 0.445107
(Epoch 36 / 200) train acc: 0.920000; val_acc: 0.178000
(Epoch 37 / 200) train acc: 0.960000; val_acc: 0.170000
(Epoch 38 / 200) train acc: 0.960000; val_acc: 0.179000
(Epoch 39 / 200) train acc: 0.940000; val_acc: 0.186000
(Epoch 40 / 200) train acc: 0.940000; val_acc: 0.160000
(Iteration 81 / 400) loss: 0.494688
(Epoch 41 / 200) train acc: 0.960000; val_acc: 0.170000
(Epoch 42 / 200) train acc: 0.940000; val_acc: 0.186000
(Epoch 43 / 200) train acc: 0.960000; val_acc: 0.170000
(Epoch 44 / 200) train acc: 0.940000; val_acc: 0.191000
(Epoch 45 / 200) train acc: 0.980000; val_acc: 0.180000
(Iteration 91 / 400) loss: 0.109107
(Epoch 46 / 200) train acc: 0.960000; val_acc: 0.170000
```
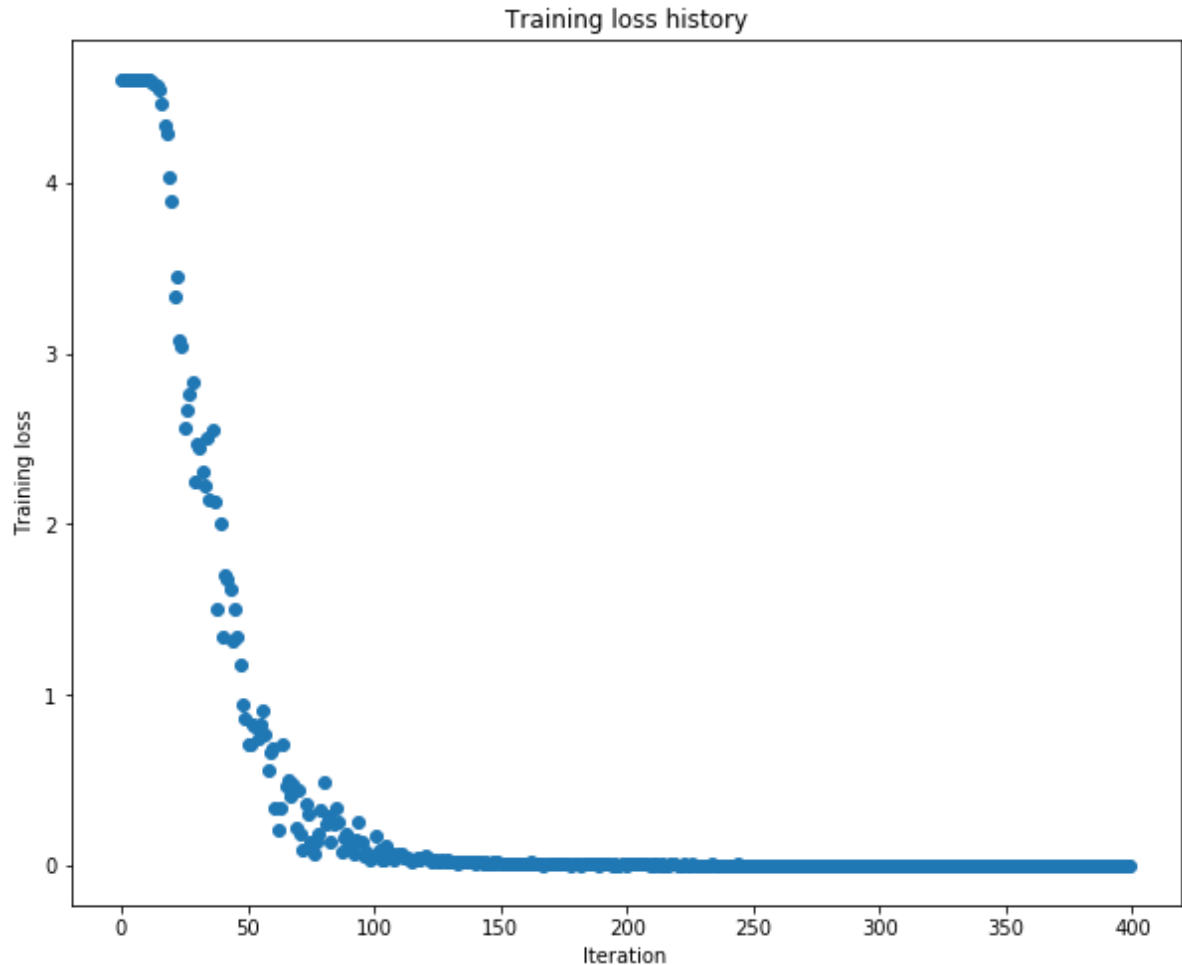
```
(Epoch 47 / 200) train acc: 0.960000; val_acc: 0.181000
(Epoch 48 / 200) train acc: 1.000000; val_acc: 0.192000
(Epoch 49 / 200) train acc: 0.980000; val_acc: 0.189000
(Epoch 50 / 200) train acc: 1.000000; val_acc: 0.175000
(Iteration 101 / 400) loss: 0.059865
(Epoch 51 / 200) train acc: 0.980000; val_acc: 0.194000
(Epoch 52 / 200) train acc: 1.000000; val_acc: 0.182000
(Epoch 53 / 200) train acc: 1.000000; val_acc: 0.184000
(Epoch 54 / 200) train acc: 1.000000; val_acc: 0.188000
(Epoch 55 / 200) train acc: 1.000000; val_acc: 0.189000
(Iteration 111 / 400) loss: 0.054751
(Epoch 56 / 200) train acc: 1.000000; val_acc: 0.185000
(Epoch 57 / 200) train acc: 1.000000; val_acc: 0.186000
(Epoch 58 / 200) train acc: 1.000000; val_acc: 0.183000
(Epoch 59 / 200) train acc: 1.000000; val_acc: 0.186000
(Epoch 60 / 200) train acc: 1.000000; val_acc: 0.188000
(Iteration 121 / 400) loss: 0.063551
(Epoch 61 / 200) train acc: 1.000000; val_acc: 0.187000
(Epoch 62 / 200) train acc: 1.000000; val_acc: 0.184000
(Epoch 63 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 64 / 200) train acc: 1.000000; val_acc: 0.182000
(Epoch 65 / 200) train acc: 1.000000; val_acc: 0.184000
(Iteration 131 / 400) loss: 0.024325
(Epoch 66 / 200) train acc: 1.000000; val_acc: 0.179000
(Epoch 67 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 68 / 200) train acc: 1.000000; val_acc: 0.181000
(Epoch 69 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 70 / 200) train acc: 1.000000; val_acc: 0.182000
(Iteration 141 / 400) loss: 0.016950
(Epoch 71 / 200) train acc: 1.000000; val_acc: 0.184000
(Epoch 72 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 73 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 74 / 200) train acc: 1.000000; val_acc: 0.174000
(Epoch 75 / 200) train acc: 1.000000; val_acc: 0.174000
(Iteration 151 / 400) loss: 0.015496
(Epoch 76 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 77 / 200) train acc: 1.000000; val_acc: 0.171000
(Epoch 78 / 200) train acc: 1.000000; val_acc: 0.174000
(Epoch 79 / 200) train acc: 1.000000; val_acc: 0.171000
(Epoch 80 / 200) train acc: 1.000000; val_acc: 0.172000
(Iteration 161 / 400) loss: 0.012433
(Epoch 81 / 200) train acc: 1.000000; val_acc: 0.172000
(Epoch 82 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 83 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 84 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 85 / 200) train acc: 1.000000; val_acc: 0.170000
(Iteration 171 / 400) loss: 0.015077
(Epoch 86 / 200) train acc: 1.000000; val_acc: 0.170000
(Epoch 87 / 200) train acc: 1.000000; val_acc: 0.173000
(Epoch 88 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 89 / 200) train acc: 1.000000; val_acc: 0.172000
(Epoch 90 / 200) train acc: 1.000000; val_acc: 0.172000
(Iteration 181 / 400) loss: 0.014952
(Epoch 91 / 200) train acc: 1.000000; val_acc: 0.174000
(Epoch 92 / 200) train acc: 1.000000; val_acc: 0.174000
(Epoch 93 / 200) train acc: 1.000000; val_acc: 0.175000
(Epoch 94 / 200) train acc: 1.000000; val_acc: 0.175000
```

```
(Epoch 95 / 200) train acc: 1.000000; val_acc: 0.176000
(Iteration 191 / 400) loss: 0.009590
(Epoch 96 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 97 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 98 / 200) train acc: 1.000000; val_acc: 0.181000
(Epoch 99 / 200) train acc: 1.000000; val_acc: 0.181000
(Epoch 100 / 200) train acc: 1.000000; val_acc: 0.180000
(Iteration 201 / 400) loss: 0.006620
(Epoch 101 / 200) train acc: 1.000000; val_acc: 0.179000
(Epoch 102 / 200) train acc: 1.000000; val_acc: 0.179000
(Epoch 103 / 200) train acc: 1.000000; val_acc: 0.182000
(Epoch 104 / 200) train acc: 1.000000; val_acc: 0.182000
(Epoch 105 / 200) train acc: 1.000000; val_acc: 0.183000
(Iteration 211 / 400) loss: 0.009187
(Epoch 106 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 107 / 200) train acc: 1.000000; val_acc: 0.182000
(Epoch 108 / 200) train acc: 1.000000; val_acc: 0.182000
(Epoch 109 / 200) train acc: 1.000000; val_acc: 0.182000
(Epoch 110 / 200) train acc: 1.000000; val_acc: 0.182000
(Iteration 221 / 400) loss: 0.005830
(Epoch 111 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 112 / 200) train acc: 1.000000; val_acc: 0.172000
(Epoch 113 / 200) train acc: 1.000000; val_acc: 0.173000
(Epoch 114 / 200) train acc: 1.000000; val_acc: 0.176000
(Epoch 115 / 200) train acc: 1.000000; val_acc: 0.174000
(Iteration 231 / 400) loss: 0.007205
(Epoch 116 / 200) train acc: 1.000000; val_acc: 0.174000
(Epoch 117 / 200) train acc: 1.000000; val_acc: 0.172000
(Epoch 118 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 119 / 200) train acc: 1.000000; val_acc: 0.176000
(Epoch 120 / 200) train acc: 1.000000; val_acc: 0.174000
(Iteration 241 / 400) loss: 0.005731
(Epoch 121 / 200) train acc: 1.000000; val_acc: 0.174000
(Epoch 122 / 200) train acc: 1.000000; val_acc: 0.174000
(Epoch 123 / 200) train acc: 1.000000; val_acc: 0.175000
(Epoch 124 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 125 / 200) train acc: 1.000000; val_acc: 0.178000
(Iteration 251 / 400) loss: 0.005841
(Epoch 126 / 200) train acc: 1.000000; val_acc: 0.173000
(Epoch 127 / 200) train acc: 1.000000; val_acc: 0.171000
(Epoch 128 / 200) train acc: 1.000000; val_acc: 0.173000
(Epoch 129 / 200) train acc: 1.000000; val_acc: 0.172000
(Epoch 130 / 200) train acc: 1.000000; val_acc: 0.171000
(Iteration 261 / 400) loss: 0.006214
(Epoch 131 / 200) train acc: 1.000000; val_acc: 0.176000
(Epoch 132 / 200) train acc: 1.000000; val_acc: 0.172000
(Epoch 133 / 200) train acc: 1.000000; val_acc: 0.175000
(Epoch 134 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 135 / 200) train acc: 1.000000; val_acc: 0.175000
(Iteration 271 / 400) loss: 0.004203
(Epoch 136 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 137 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 138 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 139 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 140 / 200) train acc: 1.000000; val_acc: 0.179000
(Iteration 281 / 400) loss: 0.004103
(Epoch 141 / 200) train acc: 1.000000; val_acc: 0.179000
```

```
(Epoch 142 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 143 / 200) train acc: 1.000000; val_acc: 0.175000
(Epoch 144 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 145 / 200) train acc: 1.000000; val_acc: 0.181000
(Iteration 291 / 400) loss: 0.003961
(Epoch 146 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 147 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 148 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 149 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 150 / 200) train acc: 1.000000; val_acc: 0.176000
(Iteration 301 / 400) loss: 0.005293
(Epoch 151 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 152 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 153 / 200) train acc: 1.000000; val_acc: 0.182000
(Epoch 154 / 200) train acc: 1.000000; val_acc: 0.179000
(Epoch 155 / 200) train acc: 1.000000; val_acc: 0.179000
(Iteration 311 / 400) loss: 0.003632
(Epoch 156 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 157 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 158 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 159 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 160 / 200) train acc: 1.000000; val_acc: 0.177000
(Iteration 321 / 400) loss: 0.003431
(Epoch 161 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 162 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 163 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 164 / 200) train acc: 1.000000; val_acc: 0.181000
(Epoch 165 / 200) train acc: 1.000000; val_acc: 0.180000
(Iteration 331 / 400) loss: 0.003142
(Epoch 166 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 167 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 168 / 200) train acc: 1.000000; val_acc: 0.181000
(Epoch 169 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 170 / 200) train acc: 1.000000; val_acc: 0.179000
(Iteration 341 / 400) loss: 0.004564
(Epoch 171 / 200) train acc: 1.000000; val_acc: 0.176000
(Epoch 172 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 173 / 200) train acc: 1.000000; val_acc: 0.179000
(Epoch 174 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 175 / 200) train acc: 1.000000; val_acc: 0.178000
(Iteration 351 / 400) loss: 0.003547
(Epoch 176 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 177 / 200) train acc: 1.000000; val_acc: 0.179000
(Epoch 178 / 200) train acc: 1.000000; val_acc: 0.179000
(Epoch 179 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 180 / 200) train acc: 1.000000; val_acc: 0.178000
(Iteration 361 / 400) loss: 0.003727
(Epoch 181 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 182 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 183 / 200) train acc: 1.000000; val_acc: 0.178000
(Epoch 184 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 185 / 200) train acc: 1.000000; val_acc: 0.177000
(Iteration 371 / 400) loss: 0.002589
(Epoch 186 / 200) train acc: 1.000000; val_acc: 0.179000
(Epoch 187 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 188 / 200) train acc: 1.000000; val_acc: 0.179000
(Epoch 189 / 200) train acc: 1.000000; val_acc: 0.179000
```

```
(Epoch 190 / 200) train acc: 1.000000; val_acc: 0.180000
(Iteration 381 / 400) loss: 0.003914
(Epoch 191 / 200) train acc: 1.000000; val_acc: 0.180000
(Epoch 192 / 200) train acc: 1.000000; val_acc: 0.179000
(Epoch 193 / 200) train acc: 1.000000; val_acc: 0.181000
(Epoch 194 / 200) train acc: 1.000000; val_acc: 0.179000
(Epoch 195 / 200) train acc: 1.000000; val_acc: 0.178000
(Iteration 391 / 400) loss: 0.003073
(Epoch 196 / 200) train acc: 1.000000; val_acc: 0.177000
(Epoch 197 / 200) train acc: 1.000000; val_acc: 0.179000
(Epoch 198 / 200) train acc: 1.000000; val_acc: 0.179000
(Epoch 199 / 200) train acc: 1.000000; val_acc: 0.179000
(Epoch 200 / 200) train acc: 1.000000; val_acc: 0.179000
```



Training loss history

In [ ]:

In [ ]:

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

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

Please print out the workbook entirely when completed.

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

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

```python
In [2]:  import random
         import numpy as np
         from cs231n.data_utils import load_CIFAR10
         import matplotlib.pyplot as plt

         %matplotlib inline
         %load_ext autoreload
         %autoreload 2

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

# Toy example

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

```python
In [3]:  from nndl.neural_net import TwoLayerNet
```

In [4]:
```python
# 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 [5]:
```python
## 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.
#My debug code
print('Xshape: ', X.shape)
#End my debug code

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)))
```

```
Xshape:  (5, 4)
Your scores:
[[-1.07260209  0.05083871 -0.87253915]
 [-2.02778743 -0.10832494 -1.52641362]
 [-0.74225908  0.15259725 -0.39578548]
 [-0.38172726  0.10835902 -0.17328274]
 [-0.64417314 -0.18886813 -0.41106892]]

correct scores:
[[-1.07260209  0.05083871 -0.87253915]
 [-2.02778743 -0.10832494 -1.52641362]
 [-0.74225908  0.15259725 -0.39578548]
 [-0.38172726  0.10835902 -0.17328274]
 [-0.64417314 -0.18886813 -0.41106892]]

Difference between your scores and correct scores:
3.381231233889892e-08
```

## Forward pass loss

In [6]: 
```
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 [7]: 
```
from cs231n.gradient_check import eval_numerical_gradient

# Use numeric gradient checking to check your implementation of the backward p
ass.
# 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)

#print(loss, grads)

# 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], verbos
e=False)
    print('{} max relative error: {}'.format(param_name, rel_error(param_grad_
num, grads[param_name])))
```

```
W2 max relative error: 2.9632227682005116e-10
b2 max relative error: 1.2482714253983918e-09
W1 max relative error: 1.2832823337649917e-09
b1 max relative error: 3.172680092703762e-09
```

## Training the network

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

```
In [8]:  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:  0.014498406590265567



# Classify CIFAR-10

Do classification on the CIFAR-10 dataset.

In [9]:
```python
from cs231n.data_utils import load_CIFAR10

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

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

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

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

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


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

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

## Running SGD

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

```python
In [10]: 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 2.302122329647926
iteration 200 / 1000: loss 2.2956767854707882
iteration 300 / 1000: loss 2.2523144504019696
iteration 400 / 1000: loss 2.1896338140489533
iteration 500 / 1000: loss 2.117053945819248
iteration 600 / 1000: loss 2.0653486572337925
iteration 700 / 1000: loss 1.9915273825850979
iteration 800 / 1000: loss 2.0040533587870257
iteration 900 / 1000: loss 1.9480758500797803
Validation accuracy:  0.282
```

# 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)?

```python
In [11]: stats['train_acc_history']
```

```
Out[11]: [0.095, 0.15, 0.255, 0.25, 0.32]
```

In [12]:
```python
# ================================================================ #
# YOUR CODE HERE:
#   Do some debugging to gain some insight into why the optimization
#   isn't great.
# ================================================================ #

# Plot the loss function and train / validation accuracies

plt.plot(stats['loss_history'])
plt.xlabel('Each Iteration')
plt.ylabel('Training Loss')
plt.title('Training Loss History')
plt.show()

plt.plot(stats['train_acc_history'])
plt.xlabel('Each Iteration')
plt.ylabel('Training Accuracy')
plt.title('Training Accuracy History')
plt.show()

plt.plot(stats['val_acc_history'])
plt.xlabel('Each Iteration')
plt.ylabel('Validation Accuracy')
plt.title('Validation Accuracy History')
plt.show()

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

## Training Loss History



## Training Accuracy History



## Validation Accuracy History

# Answers:

(1) As can be observed, the training loss starts zig zagging about 220 iterations in. This indicates that at about 220 iterations, our learning rate is too high. This causes the weights to overcorrect each step it takes, which resulted in the zig zag behavior seen in the graph. Additionally, the training and validation accuracies show linear behavior after 3.0. This suggests that we could train our model for more iterations until the slopes of these accuracies start to plateau.

(2) One method that many of my colleagues have used at a previous internship were adaptive learning rates. There are different ways, but Adagrad, Adam Optimization, RMSprop, momentum, etc. Additionally, the learning rate decay could be increased in order to reduce the zigzagging over each iteration. We could also increase the number of iterations to increase the validation accuracy.

# 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 [13]:
```python
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)!
# ================================================================ #
import time
t = time.time()
print('starting')
best_val_acc = 0.5

learning_rates = np.linspace(1e-4, 5e-3, 10)
num_iterations = [1500]
m_batch = 200
learning_rate_decays = np.linspace(0.95,0,9,5)
regs = np.linspace(0.15, 0.25, 3)
best_hyperparameters = list(range(4))

break_allloops = False

for learning_rate in learning_rates:
    if break_allloops:
        print('learning_rate')
        break
    for iters in num_iterations:
        print('num_iter')
        if break_allloops:
            break
        for decay in learning_rate_decays:
            print('decay')
            if break_allloops:
                break
            for reg in regs:
                print('reg')
                if break_allloops:
                    break
                mNet = TwoLayerNet(input_size, hidden_size, num_classes)

                stats = mNet.train(X_train, y_train, X_val, y_val, num_iters = iters,
                                    batch_size=m_batch, learning_rate = learning_rate, learning_rate_decay=decay,
                                    reg=reg, verbose = False)
                val_acc = np.amax(stats['val_acc_history'])
                epoch = np.argmax(stats['val_acc_history'])
                m_iteration = 1500

                if val_acc > best_val_acc:
```

```
                        best_net = mNet

                        best_val_acc = val_acc
                        best_hyperparameters = [learning_rate, iters, decay, reg]
                        break_allloops = True
                        print(best_hyperparameters)
                        print(best_val_acc)


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

#Hyperparameters calculated are: [0.0022777777777777774, 1500, 0.8312499999999
999, 0.15]
#With validation accuracy 0.503
```

```
starting
num_iter
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
num_iter
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
```

```
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
num_iter
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
num_iter
decay
```

```
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
num_iter
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
```

```
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
decay
reg
reg
reg
num_iter
decay
reg
reg
reg
decay
reg
reg
[0.002822222222222222, 1500, 0.8312499999999999, 0.2]
0.504
reg
decay
learning_rate
Validation accuracy:  0.492
```

In [14]:
```python
from cs231n.vis_utils import visualize_grid

# Visualize the weights of the network

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

show_net_weights(subopt_net)
show_net_weights(best_net)
```





# Question:

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

## Answer:

(1) The weights in the suboptimal network look very similar. There aren't many significant color changes within the suboptimal network, and the shapes of the weights seem similar as well. The best network that my debugging code above found had more varied weights in regards to shape and color.

## Evaluate on test set

```
In [15]: test_acc = (best_net.predict(X_test) == y_test).mean()
         print('Test accuracy: ', test_acc)
```

```
Test accuracy:  0.49
```

In [ ]:

```python
import numpy as np
import pdb

"""
This code was originally written for CS 231n at Stanford University
(cs231n.stanford.edu).  It has been modified in various areas for use in the
ECE 239AS class at UCLA.  This includes the descriptions of what code to
implement as well as some slight potential changes in variable names to be
consistent with class nomenclature.  We thank Justin Johnson & Serena Yeung for
permission to use this code.  To see the original version, please visit
cs231n.stanford.edu.
"""


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

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

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

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

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

  x_r = x.reshape(x.shape[0], -1)
  out = np.dot(x_r, w) + b

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

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


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

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

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


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

    x_r = x.reshape(x.shape[0], -1)
    db = np.sum(dout, axis = 0)
    dw = np.dot(x_r.T, dout)
    dx = np.dot(dout, w.T).reshape(x.shape)

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


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

    return dx, dw, db

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

    Input:
    - x: Inputs, of any shape

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

    relu = lambda x: x * (x > 0)
    out = relu(x)
    # ================================================================ #
    # END YOUR CODE HERE
    # ================================================================ #

    cache = x
    return out, cache


def relu_backward(dout, cache):
    """
    Computes the backward pass for a layer of rectified linear units (ReLUs).

    Input:
```

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

        Returns:
        - dx: Gradient with respect to x
        """
        x = cache

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

        # ReLU directs linearly to those > 0
        dx = dout * (x > 0)

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

        return dx

def svm_loss(x, y):
    """
    Computes the loss and gradient using for multiclass SVM classification.

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

    Returns a tuple of:
    - loss: Scalar giving the loss
    - dx: Gradient of the loss with respect to x
    """
    N = x.shape[0]
    correct_class_scores = x[np.arange(N), y]
    margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
    margins[np.arange(N), y] = 0
    loss = np.sum(margins) / N
    num_pos = np.sum(margins > 0, axis=1)
    dx = np.zeros_like(x)
    dx[margins > 0] = 1
    dx[np.arange(N), y] -= num_pos
    dx /= N
    return loss, dx


def softmax_loss(x, y):
    """
    Computes the loss and gradient for softmax classification.

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

    Returns a tuple of:
    - loss: Scalar giving the loss
    - dx: Gradient of the loss with respect to x
```

```
    """

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

```python
import numpy as np
import matplotlib.pyplot as plt

"""
This code was originally written for CS 231n at Stanford University
(cs231n.stanford.edu).  It has been modified in various areas for use in the
ECE 239AS class at UCLA.  This includes the descriptions of what code to
implement as well as some slight potential changes in variable names to be
consistent with class nomenclature.  We thank Justin Johnson & Serena Yeung for
permission to use this code.  To see the original version, please visit
cs231n.stanford.edu.
"""


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

    In other words, the network has the following architecture:

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

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

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

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

        Inputs:
        - input_size: The dimension D of the input data.
        - hidden_size: The number of neurons H in the hidden layer.
        - output_size: The number of classes C.
        """
        self.params = {}
        self.params['W1'] = std * np.random.randn(hidden_size, input_size)
        self.params['b1'] = np.zeros(hidden_size)
        self.params['W2'] = std * np.random.randn(output_size, hidden_size)
        self.params['b2'] = np.zeros(output_size)


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

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

```
        Returns:
        If y is None, return a matrix scores of shape (N, C) where scores[i, c] is
        the score for class c on input X[i].

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

        # Compute the forward
        scores = None


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

        relu = lambda x: x * (x > 0)
        h1 = relu(np.dot(X, W1.T) + b1)
        out = np.dot(h1, W2.T) + b2
        scores = np.copy(out)

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


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

        # Compute the loss
        loss = None


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

        # scores is num_examples by num_classes
        p = np.exp(scores - np.max(scores, axis = 1, keepdims=True))
        p /= np.sum(p, axis=1, keepdims = True)

        loss = -np.sum(np.log(p[np.arange(N), y])) / N

        ds = p.copy()
        ds[np.arange(N), y] -= 1
        ds /= N

        dreg = reg*0.5*(np.sum(W1**2) + np.sum(W2**2))
```

```
        loss += dreg
        # ============================================================ #
        # END YOUR CODE HERE
        # ============================================================ #

        grads = {}

        # ============================================================ #
        # YOUR CODE HERE:
        #       Implement the backward .  Compute the derivatives of the
        #       weights and the biases.  Store the results in the grads
        #       dictionary.  e.g., grads['W1'] should store the gradient for
        #       W1, and be of the same size as W1.
        # ============================================================ #

        grads['W2'] = np.dot(ds.T ,h1) + reg * W2
        grads['b2'] = np.dot(np.ones(N), ds)

        dh1 = np.dot(ds, W2)
        dh1[h1==0] = 0

        grads['W1'] = np.dot(dh1.T, X) + reg*W1
        grads['b1'] = np.dot(np.ones(N), dh1)

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

        return loss, grads

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

        Inputs:
        - X: A numpy array of shape (N, D) giving training data.
        - y: A numpy array f shape (N,) giving training labels; y[i] = c means that
                X[i] has label c, where 0 <= c < C.
        - X_val: A numpy array of shape (N_val, D) giving validation data.
        - y_val: A numpy array of shape (N_val,) giving validation labels.
        - learning_rate: Scalar giving learning rate for optimization.
        - learning_rate_decay: Scalar giving factor used to decay the learning rate
                after each epoch.
        - reg: Scalar giving regularization strength.
        - num_iters: Number of steps to take when optimizing.
        - batch_size: Number of training examples to use per step.
        - verbose: boolean; if true print progress during optimization.
        """
        num_train = X.shape[0]
        iterations_per_epoch = max(num_train / batch_size, 1)

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

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

                # ============================================================ #
                # YOUR CODE HERE:
```

```
                    #        Create a minibatch by sampling batch_size samples randomly.
                    # ============================================================== #
                    batch_indexes = np.random.choice(list(range(len(X))), size=batch_size, replace=True)

                    #print(batch_indexes)

                    X_batch = [X[i] for i in batch_indexes]
                    y_batch = [y[i] for i in batch_indexes]

                    X_batch = np.vstack(X_batch)

                    #print(X_batch)

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

                     # Compute loss and gradients using the current minibatch
                    loss, grads = self.loss(X_batch, y=y_batch, reg=reg)
                    loss_history.append(loss)

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

                    reg_fact = 1 - learning_rate * reg
                    #for key in self.params:
                    #        self.params[key] = reg_fact * (self.params[key]) - learning_rate * grads[key]

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

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

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

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

                            # Decay learning rate
                            learning_rate *= learning_rate_decay

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

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

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

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

    # ================================================================ #
    # YOUR CODE HERE:
    #       Predict the class given the input data.
    # ================================================================ #

    y_pred = np.argmax(self.loss(X), axis=1)



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

    return y_pred
```

```python
import numpy as np

from .layers import *
from .layer_utils import *

"""
This code was originally written for CS 231n at Stanford University
(cs231n.stanford.edu).  It has been modified in various areas for use in the
ECE 239AS class at UCLA.  This includes the descriptions of what code to
implement as well as some slight potential changes in variable names to be
consistent with class nomenclature.  We thank Justin Johnson & Serena Yeung for
permission to use this code.  To see the original version, please visit
cs231n.stanford.edu.
"""

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

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

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

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

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

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

    # ============================================================== #
    # YOUR CODE HERE:
    #   Initialize W1, W2, b1, and b2.  Store these as self.params['W1'],
    #   self.params['W2'], self.params['b1'] and self.params['b2']. The
    #   biases are initialized to zero and the weights are initialized
    #   so that each parameter has mean 0 and standard deviation weight_scale.
    #   The dimensions of W1 should be (input_dim, hidden_dim) and the
    #   dimensions of W2 should be (hidden_dims, num_classes)
    # ============================================================== #

    self.params['W2'] = weight_scale * np.random.randn(hidden_dims, num_classes)
    self.params['b2'] = np.zeros(num_classes)
    self.params['W1'] = weight_scale * np.random.randn(input_dim, hidden_dims)
    self.params['b1'] = np.zeros(hidden_dims)

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

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

    Inputs:
    - X: Array of input data of shape (N, d_1, ..., d_k)
    - y: Array of labels, of shape (N,). y[i] gives the label for X[i].

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

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

```
    """
    scores = None

    # =========================================================== #
    # YOUR CODE HERE:
    #   Implement the forward pass of the two-layer neural network. Store
    #   the class scores as the variable 'scores'.  Be sure to use the layers
    #   you prior implemented.
    # =========================================================== #

    h1, cache1 = affine_relu_forward(X, self.params['W1'], self.params['b1'])
    scores, cache2 = affine_forward(h1, self.params['W2'], self.params['b2'])
    # =========================================================== #
    # END YOUR CODE HERE
    # =========================================================== #

    # If y is None then we are in test mode so just return scores
    if y is None:
      return scores

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

    loss, ds = softmax_loss(scores, y)
    dreg = self.reg * 0.5*(np.sum(self.params['W1'] ** 2) + np.sum(self.params['W2'] ** 2))
    loss += dreg

    d_h1, grads['W2'], grads['b2'] = affine_backward(ds, cache2)
    grads['W2'] += self.reg * self.params['W2']

    dx, grads['W1'], grads['b1'] = affine_relu_backward(d_h1, cache1)
    grads['W1'] += self.reg * self.params['W1']

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

    return loss, grads


class FullyConnectedNet(object):
  """
  A fully-connected neural network with an arbitrary number of hidden layers,
  ReLU nonlinearities, and a softmax loss function. This will also implement
  dropout and batch normalization as options. For a network with L layers,
  the architecture will be

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

  where batch normalization and dropout are optional, and the {...} block is
  repeated L - 1 times.

  Similar to the TwoLayerNet above, learnable parameters are stored in the
  self.params dictionary and will be learned using the Solver class.
  """

  def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,
               dropout=0, use_batchnorm=False, reg=0.0,
               weight_scale=1e-2, dtype=np.float32, seed=None):
    """
    Initialize a new FullyConnectedNet.

    Inputs:
    - hidden_dims: A list of integers giving the size of each hidden layer.
    - input_dim: An integer giving the size of the input.
    - num_classes: An integer giving the number of classes to classify.
    - dropout: Scalar between 0 and 1 giving dropout strength. If dropout=0 then
      the network should not use dropout at all.
    - use_batchnorm: Whether or not the network should use batch normalization.
    - reg: Scalar giving L2 regularization strength.
    - weight_scale: Scalar giving the standard deviation for random
      initialization of the weights.
    - dtype: A numpy datatype object; all computations will be performed using
```

```python
    this datatype. float32 is faster but less accurate, so you should use
    float64 for numeric gradient checking.
    - seed: If not None, then pass this random seed to the dropout layers. This
      will make the dropout layers deterimnstic so we can gradient check the
      model.
    """
    self.use_batchnorm = use_batchnorm
    self.use_dropout = dropout > 0
    self.reg = reg
    self.num_layers = 1 + len(hidden_dims)
    self.dtype = dtype
    self.params = {}

    # ============================================================ #
    # YOUR CODE HERE:
    #   Initialize all parameters of the network in the self.params dictionary.
    #   The weights and biases of layer 1 are W1 and b1; and in general the
    #   weights and biases of layer i are Wi and bi. The
    #   biases are initialized to zero and the weights are initialized
    #   so that each parameter has mean 0 and standard deviation weight_scale.
    # ============================================================ #

    layer_dims = np.hstack((input_dim, hidden_dims, num_classes))

    for i in list(range(1, self.num_layers)):
      Wi = 'W' + str(i)
      bi = 'b' + str(i)
      self.params[Wi] = weight_scale * np.random.randn(layer_dims[i-1], layer_dims[i])
      self.params[bi] = np.zeros(layer_dims[i])

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

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

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

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


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

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

    # Set train/test mode for batchnorm params and dropout param since they
    # behave differently during training and testing.
    if self.dropout_param is not None:
      self.dropout_param['mode'] = mode
    if self.use_batchnorm:
      for bn_param in self.bn_params:
        bn_param[mode] = mode

    scores = None

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

    cache = {}
    hidden, cache1 = affine_relu_forward(X, self.params['W1'], self.params['b1'])
    cache['c1'] = cache1
    for i in range(1, self.num_layers - 1):
```

```python
        Wi = 'W' + str(i+1)
        bi = 'b' + str(i + 1)
        ci = 'c' + str(i+1)
        if i == self.num_layers:
          hidden, cachei = affine_forward(hidden, self.params[Wi], self.params[bi])
        else:
          hidden, cachei = affine_relu_forward(hidden, self.params[Wi], self.params[bi])
        cache[ci] = cachei
    scores = hidden
    # ========================================================== #
    # END YOUR CODE HERE
    # ========================================================== #

    # If test mode return early
    if mode == 'test':
      return scores

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

    loss, ds = softmax_loss(scores, y)
    dh, grads['W' + str(self.num_layers - 1)], grads['b' + str(self.num_layers - 1)] = affine_relu_backward(ds, cache['c'+str(self.num_layers-1)])
    #dreg = np.sum(self.params['W' + str(self.num_layers - 1)]**2)

    for i in range(self.num_layers - 2, 0, -1):
      Wi = 'W' + str(i)
      bi = 'b' + str(i)
      ci = 'c' + str(i)
      #dreg += np.sum(self.params[Wi]**2)
      if i == self.num_layers:
        dh, grads[Wi], grads[bi] = affine_backward(dh, cache[ci])
      else:
        dh, grads[Wi], grads[bi] = affine_relu_backward(dh, cache[ci])
      loss += 0.5 * self.reg * np.sum(self.params[Wi]**2)


      grads[Wi] += self.reg * self.params[Wi]

    dx = dh


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