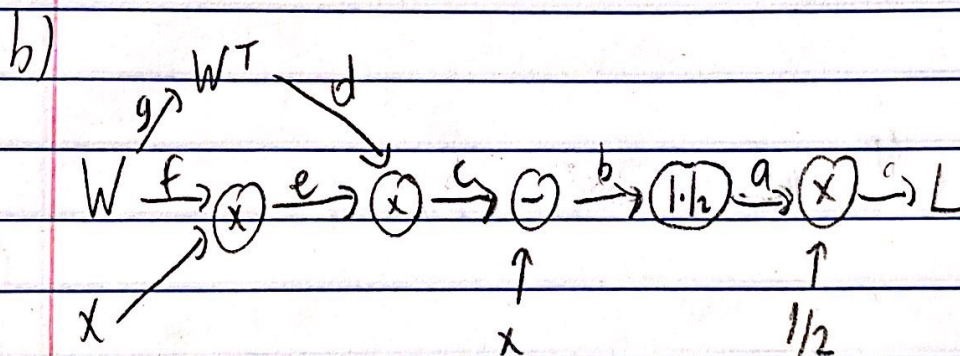


# ECE C247 HW3

1) a) We can think of the  $Wx$  operation as projecting the  $x$  vector into another dimension as in PCA. Then,  $W^T Wx$  is simply the reconstruction. To minimize the error, the reconstruction must be as close as possible to the original  $x$ . In order to have this, the projection  $Wx$  must retain as much information about  $x$  as possible. Therefore, minimizing the loss should find a  $W$  that ought to preserve information about  $x$ .



c) From the law of total derivatives, we know that when the variable we are taking the derivative for is affected by other variables, we should take the derivative with respect to these other variables and sum the results. In other words:

$$\frac{dL}{dx} = \sum_{i=1}^2 \frac{dL}{dq_i} \cdot \frac{dq_i}{dx}$$

$$= \frac{dL}{dq_1} \cdot \frac{dq_1}{dx} + \frac{dL}{dq_2} \cdot \frac{dq_2}{dx}$$

So, we need to take the derivative with respect to the two paths and sum them up.

d) By looking at the graph:

$$\frac{dL}{da} = 1/2 \quad \frac{db}{dc} = \frac{d(c-x)}{dc} = 1$$

$$\frac{da}{db} = 2b \quad \frac{dc}{dd} = \frac{d(W^T W x)}{dW^T} = x^T W^T$$

$$\frac{dc}{de} = \frac{d(W^T W x)}{d(W x)} = W \quad \frac{de}{dW} = \frac{d(W x)}{dW} = x^T$$

$$\frac{dL}{dd} = \frac{dL}{da} \cdot \frac{da}{db} \cdot \frac{db}{dc} \cdot \frac{dc}{dd}$$

$$= 1/2 \cdot 2b \cdot 1 \cdot x^T W = b x^T W^T$$

$$\frac{dL}{de} = \frac{dL}{da} \cdot \frac{da}{db} \cdot \frac{db}{dc} \cdot \frac{dc}{de}$$

$$= 1/2 \cdot 2b \cdot 1 \cdot W$$

$$\frac{dL}{dW} = \frac{dL}{dd} \cdot \frac{dd}{dW} + \frac{dL}{de} \cdot \frac{de}{dW}$$

$$= (b x^T W^T)^T + W b x^T$$

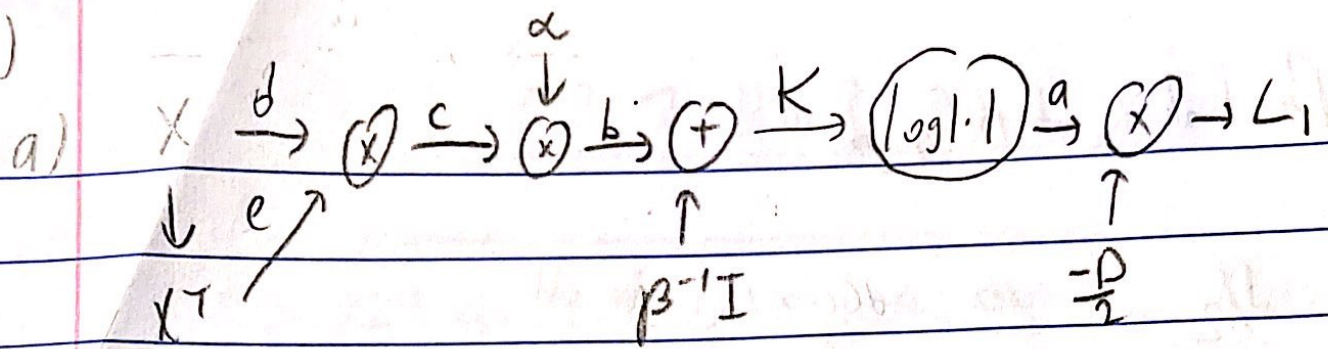
$$= W x b^T + W b x^T$$

$$\Downarrow b = (W^T W x - x)$$

$$\frac{dL}{dW} = W x (W^T W x - x)^T + W (W^T W x - x) x^T$$



2)



b)

$$\frac{dL_1}{da} = \frac{-D}{2} \quad \frac{da}{dK} = \log |K| = K^{-T} \quad \frac{dK}{db} = 1$$

$$\frac{db}{dc} = \alpha \quad \frac{dc}{dd} = x \quad \frac{dc}{de} = x^T$$

$$\frac{dL_1}{dc} = \frac{dL_1}{da} \cdot \frac{da}{dK} \cdot \frac{dK}{db} \cdot \frac{db}{dc}$$

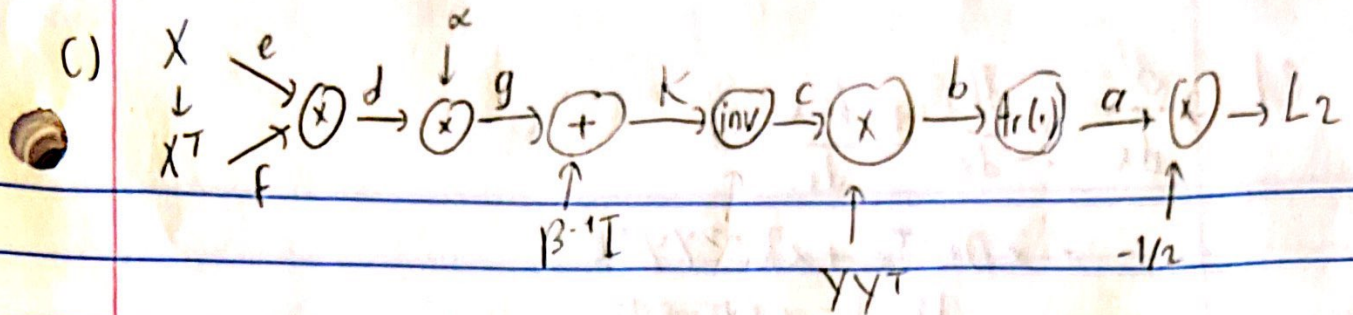
$$= \frac{-D}{2} \cdot K^{-T} \cdot 1 \cdot \alpha$$

$$\frac{dL_1}{dx} = \frac{dL_1}{dc} \cdot \frac{dc}{dx} + \frac{dL_1}{de} \cdot \frac{de}{dx}$$

$$= \frac{-D}{2} K^{-T} \cdot 1 \cdot \alpha (x+x)$$

$$= -\alpha DK^{-T} x$$





$$d) \frac{dL_2}{dK} = -K^{-T} \frac{dL_2}{da} K^{-T}$$

$$\frac{dL_2}{dK^{-1}} = \frac{dL_2}{da} \cdot \frac{da}{dK^{-1}}$$

$$\frac{dL_2}{da} = -1/2 \quad \frac{da}{dK^{-1}} = \frac{d(\text{tr}(K^{-1}YY^T))}{dK^{-1}} = YY^T$$

$$\frac{dL_2}{dK^{-1}} = -\frac{1}{2} YY^T$$

$$\text{So, } \frac{dL_2}{dK} = -K^{-T} - \frac{1}{2} YY^T K^{-T}$$

$$\frac{dK}{dg} = 1, \quad \frac{dg}{da} = \alpha, \quad \frac{da}{de} = X, \quad \frac{de}{df} = X^T$$

$$\frac{dL_2}{dX} = \frac{dL_2}{dK} \cdot \frac{dK}{dg} \cdot \frac{dg}{da} \left( \frac{da}{de} \cdot \frac{de}{dX} + \frac{da}{df} \cdot \frac{df}{dX} \right)$$

$$= -K^{-T} - \frac{1}{2} YY^T K^{-T} \cdot \alpha \cdot 2X$$

$$= \alpha K^{-T} YY^T K^{-T} X$$

$$e) \quad \frac{dL}{dx} = \frac{dL_1}{dx} + \frac{dL_2}{dx}$$

$$= -\alpha D K^{-T} x + \alpha K^{-T} Y Y^T K^{-T} x$$

$$\Downarrow \quad K = \alpha X X^T + \beta^{-1} I$$

$$= -\alpha D (\alpha X X^T + \beta^{-1} I)^{-T} x + \alpha (\alpha X X^T + \beta^{-1} I)^{-T} Y Y^T (\alpha X X^T + \beta^{-1} I)^{-T} x$$

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

```
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  $0 \leq y[i] < C$ . This parameter is optional; if it is not passed then we only return scores, and if it is passed then we instead return the loss and gradients.
- `reg`: Regularization strength.

Returns:

If `y` is `None`, return a matrix scores of shape  $(N, C)$  where `scores[i, c]` is

the score for class `c` on input `X[i]`.

If `y` is not `None`, instead return a tuple of:

- `loss`: Loss (data loss and regularization loss) for this batch of training samples.
- `grads`: Dictionary mapping parameter names to gradients of those parameters with respect to the loss function; has the same keys as `self.params`.

```
"""
# Unpack variables from the params dictionary
W1, b1 = self.params['W1'], self.params['b1']
W2, b2 = self.params['W2'], self.params['b2']
N, D = X.shape

# Compute the forward pass
scores = None

# ===== #
# YOUR CODE HERE:
# Calculate the output scores of the neural network. The result
# should be (N, C). As stated in the description for this class,
# there should not be a ReLU layer after the second FC layer.
# The output of the second FC layer is the output scores. Do not
# use a for loop in your implementation.
# ===== #
#b1 = b1.reshape([len(b1),1])
#b2 = b2.reshape([len(b2),1])
layer1Out = W1.dot(X.T) + b1[:,None]
ReluOut = np.maximum(0,layer1Out)
layer2Out = W2.dot(ReluOut) + b2[:,None]
scores = layer2Out.T

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

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

# Compute the loss
loss = None

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

# scores is num_examples by num_classes
soft = np.exp(scores)
sums = np.sum(soft,axis=1)
probs = soft / sums[:,None]
predsForClass = probs[np.arange(y.shape[0]),y]
SoftmaxLoss = np.mean(-np.log(predsForClass))
```



```

l2regularization = np.sum(W1**2) + np.sum(W2**2)
l2regularization = 0.5*reg*l2regularization
loss=SoftmaxLoss+l2regularization
# ===== #
# END YOUR CODE HERE
# ===== #

grads = {}

# ===== #
# YOUR CODE HERE:
# Implement the backward pass. Compute the derivatives of the
# weights and the biases. Store the results in the grads
# dictionary. e.g., grads['W1'] should store the gradient for
# W1, and be of the same size as W1.
# ===== #
#A1 = W1X+B1
#A2 = RELU(A1)
#A3 = W2A2+B2
#A4 = SOFTMAX(A3)
#DL/DA3 = PREDICTIONS-LABELS ONEHOT
#DA3/DW2 = A2
#DA3/DA2 = W2
#DA2/DA1 = 0 OR 1
#DA1/DW1 = A1
grad = probs.copy()
grad[np.arange(y.shape[0]),y] -= 1
dLA3=grad/X.shape[0]
dA3W2=ReluOut.copy()
b2grad = np.sum(dLA3,axis=0)
w2grad = np.dot(dLA3.T, dA3W2.T) + reg*W2
dA3A2 = W2
dA2dA1 = layer10out.copy()
dA2dA1[dA2dA1<0]=0
dA2dA1[dA2dA1>0]=1
dA1dW1 = X.copy()

kronecker = ((dLA3.dot(dA3A2)).T*dA2dA1)
b1grad = np.sum(kronecker.T,axis=0)
w1grad = (kronecker.dot(dA1dW1)) + reg*W1

grads["W1"]=w1grad
grads["b1"]=b1grad
grads["W2"]=w2grad
grads["b2"]=b2grad

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

return loss, grads

def train(self, X, y, X_val, y_val,
          learning_rate=1e-3, learning_rate_decay=0.95,
          reg=1e-5, num_iters=100,
          batch_size=200, verbose=False):
    """

```



Train this neural network using stochastic gradient descent.

Inputs:

- X: A numpy array of shape (N, D) giving training data.
- y: A numpy array of shape (N,) giving training labels;  $y[i] = c$  means that  $X[i]$  has label  $c$ , where  $0 \leq c < C$ .
- 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.
    # ===== #
    indices = np.random.choice(np.arange(num_train), batch_size)
    X_batch = X[indices]
    y_batch = y[indices]
    # ===== #
    # 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).
    # ===== #

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

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

```

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

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

    # Decay learning rate
    learning_rate *= learning_rate_decay

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

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

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

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

    # ===== #
    # YOUR CODE HERE:
    # Predict the class given the input data.
    # ===== #
    layer1out = self.params['W1'].dot(X.T) + self.params['b1'][:,None]
    ReluOut = np.maximum(0, layer1out)
    layer2out = self.params['W2'].dot(ReluOut) + self.params['b2'][:,None]
    y_pred = np.argmax(layer2out, axis = 0)

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

    return y_pred

```

# 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 [5]: 1 import random
2 import numpy as np
3 from utils.data_utils import load_CIFAR10
4 import matplotlib.pyplot as plt
5
6 %matplotlib inline
7 %load_ext autoreload
8 %autoreload 2
9
10 def rel_error(x, y):
11     """ returns relative error """
12     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

```
In [2]: 1 from nndl.neural_net import TwoLayerNet
```

```
In [6]: 1 # Create a small net and some toy data to check your implementations.
2 # Note that we set the random seed for repeatable experiments.
3
4 input_size = 4
5 hidden_size = 10
6 num_classes = 3
7 num_inputs = 5
8
9 def init_toy_model():
10     np.random.seed(0)
11     return TwoLayerNet(input_size, hidden_size, num_classes, std=1e-1)
12
13 def init_toy_data():
14     np.random.seed(1)
15     X = 10 * np.random.randn(num_inputs, input_size)
16     y = np.array([0, 1, 2, 2, 1])
17     return X, y
18
19 net = init_toy_model()
20 X, y = init_toy_data()
```



## Compute forward pass scores

In [7]:

```

1  ## Implement the forward pass of the neural network.
2
3  # Note, there is a statement if y is None: return scores, which is why
4  # the following call will calculate the scores.
5  scores = net.loss(X)
6  print('Your scores:')
7  print(scores)
8  print()
9  print('correct scores:')
10 correct_scores = np.asarray([
11     [-1.07260209,  0.05083871, -0.87253915],
12     [-2.02778743, -0.10832494, -1.52641362],
13     [-0.74225908,  0.15259725, -0.39578548],
14     [-0.38172726,  0.10835902, -0.17328274],
15     [-0.64417314, -0.18886813, -0.41106892]])
16 print(correct_scores)
17 print()
18
19 # The difference should be very small. We get < 1e-7
20 print('Difference between your scores and correct scores:')
21 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:

3.381231248461569e-08

## Forward pass loss

In [8]:

```

1  loss, _ = net.loss(X, y, reg=0.05)
2  correct_loss = 1.071696123862817
3
4  # should be very small, we get < 1e-12
5  print("Loss:", loss)
6  print('Difference between your loss and correct loss:')
7  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 [11]: 1 from utils.gradient_check import eval_numerical_gradient
2
3 # Use numeric gradient checking to check your implementation of the backward
4 # If your implementation is correct, the difference between the numeric and
5 # analytic gradients should be less than 1e-8 for each of W1, W2, b1, and b2
6
7 loss, grads = net.loss(X, y, reg=0.05)
8
9 # these should all be less than 1e-8 or so
10 for param_name in grads:
11     f = lambda W: net.loss(X, y, reg=0.05)[0]
12     param_grad_num = eval_numerical_gradient(f, net.params[param_name], verb
13     print('{} max relative error: {}'.format(param_name, rel_error(param_gra
```

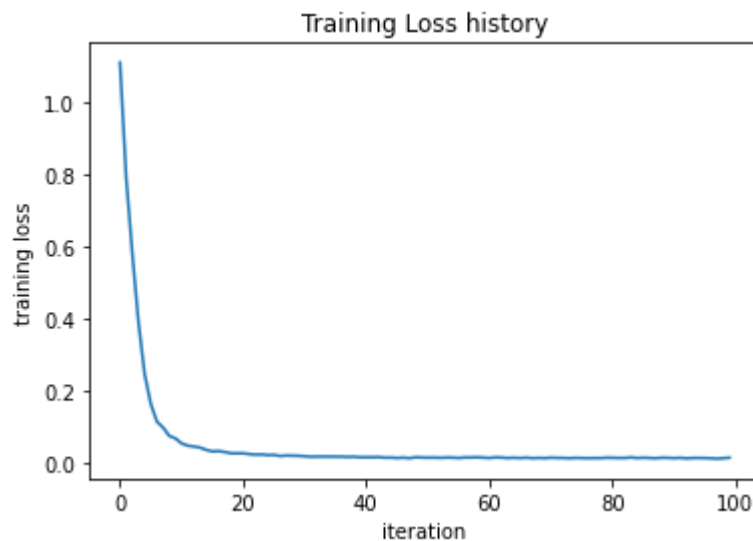
```
W1 max relative error: 1.2832874456864775e-09
b1 max relative error: 3.1726806716844575e-09
W2 max relative error: 2.9632227682005116e-10
b2 max relative error: 1.2482660547101085e-09
```

## Training the network

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

```
In [12]: 1 net = init_toy_model()
2 stats = net.train(X, y, X, y,
3                 learning_rate=1e-1, reg=5e-6,
4                 num_iters=100, verbose=False)
5
6 print('Final training loss: ', stats['loss_history'][-1])
7
8 # plot the loss history
9 plt.plot(stats['loss_history'])
10 plt.xlabel('iteration')
11 plt.ylabel('training loss')
12 plt.title('Training Loss history')
13 plt.show()
```

Final training loss: 0.014497864587765957



## Classify CIFAR-10

Do classification on the CIFAR-10 dataset.



```

In [6]: 1 from utils.data_utils import load_CIFAR10
2
3 def get_CIFAR10_data(num_training=49000, num_validation=1000, num_test=1000)
4     """
5     Load the CIFAR-10 dataset from disk and perform preprocessing to prepare
6     it for the two-layer neural net classifier.
7     """
8     # Load the raw CIFAR-10 data
9     cifar10_dir = 'dataset\cifar-10-batches-py'
10    X_train, y_train, X_test, y_test = load_CIFAR10(cifar10_dir)
11
12    # Subsample the data
13    mask = list(range(num_training, num_training + num_validation))
14    X_val = X_train[mask]
15    y_val = y_train[mask]
16    mask = list(range(num_training))
17    X_train = X_train[mask]
18    y_train = y_train[mask]
19    mask = list(range(num_test))
20    X_test = X_test[mask]
21    y_test = y_test[mask]
22
23    # Normalize the data: subtract the mean image
24    mean_image = np.mean(X_train, axis=0)
25    X_train -= mean_image
26    X_val -= mean_image
27    X_test -= mean_image
28
29    # Reshape data to rows
30    X_train = X_train.reshape(num_training, -1)
31    X_val = X_val.reshape(num_validation, -1)
32    X_test = X_test.reshape(num_test, -1)
33
34    return X_train, y_train, X_val, y_val, X_test, y_test
35
36
37 # Invoke the above function to get our data.
38 X_train, y_train, X_val, y_val, X_test, y_test = get_CIFAR10_data()
39 print('Train data shape: ', X_train.shape)
40 print('Train labels shape: ', y_train.shape)
41 print('Validation data shape: ', X_val.shape)
42 print('Validation labels shape: ', y_val.shape)
43 print('Test data shape: ', X_test.shape)
44 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 [15]: 1 input_size = 32 * 32 * 3
2 hidden_size = 50
3 num_classes = 10
4 net = TwoLayerNet(input_size, hidden_size, num_classes)
5
6 # Train the network
7 stats = net.train(X_train, y_train, X_val, y_val,
8                 num_iters=1000, batch_size=200,
9                 learning_rate=1e-4, learning_rate_decay=0.95,
10                reg=0.25, verbose=True)
11
12 # Predict on the validation set
13 val_acc = (net.predict(X_val) == y_val).mean()
14 print('Validation accuracy: ', val_acc)
15
16 # Save this net as the variable subopt_net for later comparison.
17 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 2.251825904316413
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 [21]: 1 stats['train_acc_history']
```

```
Out[21]: [0.095, 0.15, 0.25, 0.25, 0.315]
```

In [8]:

```
1  #EXPERIMENT TO SEE IF ACCURACY CAN BE IMPROVED
2  input_size = 32 * 32 * 3
3  hidden_size = 50
4  num_classes = 10
5  net = TwoLayerNet(input_size, hidden_size, num_classes)
6
7  # Train the network
8  stats = net.train(X_train, y_train, X_val, y_val,
9                  num_iters=4000, batch_size=300,
10                 learning_rate=1e-3, learning_rate_decay=0.95,
11                 reg=0.25, verbose=True)
12
13 # Predict on the validation set
14 val_acc = (net.predict(X_val) == y_val).mean()
15 print('Validation accuracy: ', val_acc)
```

```
iteration 0 / 4000: loss 2.3028007351054924
iteration 100 / 4000: loss 1.8671721184198595
iteration 200 / 4000: loss 1.735052073342008
iteration 300 / 4000: loss 1.587012373398663
iteration 400 / 4000: loss 1.5551458368146491
iteration 500 / 4000: loss 1.572744830491382
iteration 600 / 4000: loss 1.5569265665437162
iteration 700 / 4000: loss 1.516355265953603
iteration 800 / 4000: loss 1.444393039034542
iteration 900 / 4000: loss 1.4491723548916906
iteration 1000 / 4000: loss 1.5424145271154444
iteration 1100 / 4000: loss 1.3685976268235942
iteration 1200 / 4000: loss 1.4015107602863033
iteration 1300 / 4000: loss 1.4242167509987038
iteration 1400 / 4000: loss 1.4396977624306249
iteration 1500 / 4000: loss 1.4379210168499643
iteration 1600 / 4000: loss 1.4099537449443924
iteration 1700 / 4000: loss 1.3280501659227353
iteration 1800 / 4000: loss 1.4715080321395446
iteration 1900 / 4000: loss 1.3520156785903574
iteration 2000 / 4000: loss 1.464624787212279
iteration 2100 / 4000: loss 1.4232469178947993
iteration 2200 / 4000: loss 1.4221963080109268
iteration 2300 / 4000: loss 1.4088913206433655
iteration 2400 / 4000: loss 1.2836177125029784
iteration 2500 / 4000: loss 1.3359247207809524
iteration 2600 / 4000: loss 1.3921218742224277
iteration 2700 / 4000: loss 1.349099685098755
iteration 2800 / 4000: loss 1.3981975882536557
iteration 2900 / 4000: loss 1.3972649069961904
iteration 3000 / 4000: loss 1.2974790933323546
iteration 3100 / 4000: loss 1.3364343953961237
iteration 3200 / 4000: loss 1.3796549823141266
iteration 3300 / 4000: loss 1.348516916662763
iteration 3400 / 4000: loss 1.2696567716283156
iteration 3500 / 4000: loss 1.2787863728490556
iteration 3600 / 4000: loss 1.3769728934640058
iteration 3700 / 4000: loss 1.3975589168530587
iteration 3800 / 4000: loss 1.2330940850486645
```



```
iteration 3900 / 4000: loss 1.2913398672506897  
Validation accuracy: 0.514
```

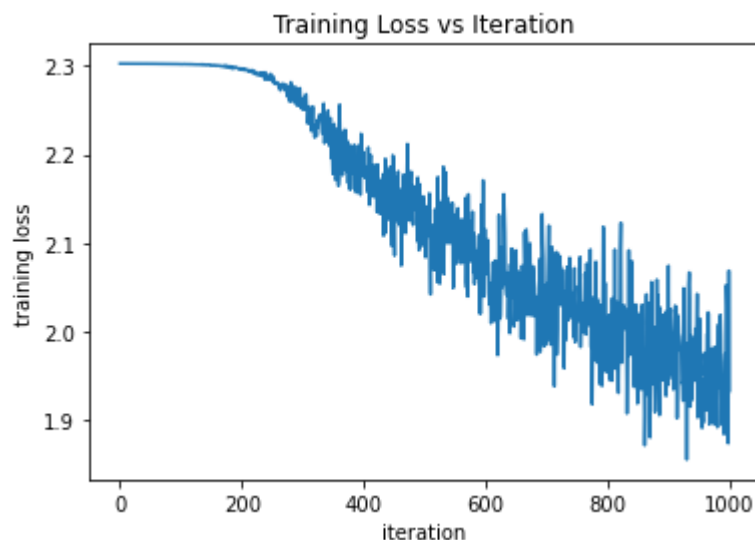


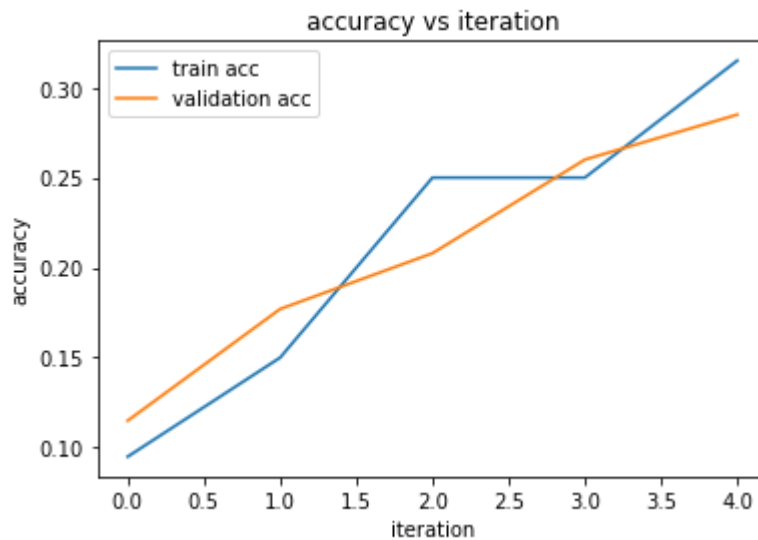
```

In [18]: 1 # ===== #
2 # YOUR CODE HERE:
3 #   Do some debugging to gain some insight into why the optimization
4 #   isn't great.
5 # ===== #
6
7 # Plot the loss function and train / validation accuracies
8
9 plt.figure()
10 plt.plot(stats['loss_history'])
11 plt.xlabel('iteration')
12 plt.ylabel('training loss')
13 plt.title('Training Loss vs Iteration')
14
15 plt.figure()
16 plt.plot(stats['train_acc_history'],label='train acc')
17 plt.plot(stats['val_acc_history'],label='validation acc')
18 plt.xlabel('iteration')
19 plt.ylabel('accuracy')
20 plt.title('accuracy vs iteration')
21 plt.legend()
22 # ===== #
23 # END YOUR CODE HERE
24 # ===== #

```

Out[18]: <matplotlib.legend.Legend at 0x268a3959730>





## Answers:

(1) By examining the progression of the training and validation accuracies, we can determine the most obvious reason for the low accuracy as the small number of iterations the model is trained. Since both training and validation accuracies are in an increasing trend, allowing the model to run for more iterations will improve performance. After some time, we should observe that the training accuracy is increasing even if the validation accuracy does not improve. This implies overfitting and training must be stopped before this point. Furthermore, examination of training loss reveals that there is little improvement in the first 200 iterations. This implies a low learning rate. Increasing the learning rate should allow the model to improve faster and thus increase accuracy. Finally, towards the end we see a large fluctuation in loss. This can imply that because our chosen batch size is too small, we get noisy estimates of the gradient which causes the oscillations in the learning curve. Increasing the batch size should lead to a smoother loss curve, which can lead to a faster decrease in loss and increase in accuracy.

(2) I would increase the number of iterations, increase the learning rate and increase the batch size. In fact, on 2 cells above it has been shown that doing so increases the validation accuracy to 51.4%. Further optimization such as changing the regularization coefficient can also be done to improve 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 [31]: 1 best_net = None # store the best model into this
2
3 # ===== #
4 # YOUR CODE HERE:
5 # Optimize over your hyperparameters to arrive at the best neural
6 # network. You should be able to get over 50% validation accuracy.
7 # For this part of the notebook, we will give credit based on the
8 # accuracy you get. Your score on this question will be multiplied by:
9 # min(floor((X - 28%)) / %22, 1)
10 # where if you get 50% or higher validation accuracy, you get full
11 # points.
12 #
13 # Note, you need to use the same network structure (keep hidden_size = 50)
14 # ===== #
15 input_size = 32 * 32 * 3
16 hidden_size = 50
17 num_classes = 10
18 learningRates = [1e-4, 1e-3, 1e-2, 1e-1]
19 regularizations = [0.3, 0.2, 0.1, 0.01, 0.001]
20 iterations = [10000, 20000, 30000]
21 #iterations=[10]
22 best_net = []
23 best_val_acc=0
24 for lr in learningRates:
25     for reg in regularizations:
26         for it in iterations:
27             net = TwoLayerNet(input_size, hidden_size, num_classes)
28             stats = net.train(X_train, y_train, X_val, y_val,
29                             num_iters=it, batch_size=200,
30                             learning_rate=lr, learning_rate_decay= 0.95,
31                             reg=reg, verbose=False)
32             print("Training complete")
33             val_acc = (net.predict(X_val) == y_val).mean()
34             print("When lr=%f, reg=%f, noIt = %d, validation accuracy=%f"%(lr, reg, it, val_acc))
35             if val_acc > best_val_acc:
36                 best_net=net
37                 best_val_acc=val_acc
38 print("SEARCH COMPLETE")
39 # ===== #
40 # END YOUR CODE HERE
41 # ===== #
42 val_acc = (best_net.predict(X_val) == y_val).mean()
43 print('Validation accuracy: ', val_acc)

```

```

Training complete
When lr=0.100000, reg=0.100000, noIt = 10000, validation accuracy=0.087000
Training complete
When lr=0.100000, reg=0.100000, noIt = 20000, validation accuracy=0.087000
Training complete
When lr=0.100000, reg=0.100000, noIt = 30000, validation accuracy=0.087000
Training complete
When lr=0.100000, reg=0.010000, noIt = 10000, validation accuracy=0.087000
Training complete
When lr=0.100000, reg=0.010000, noIt = 20000, validation accuracy=0.087000
Training complete
When lr=0.100000, reg=0.010000, noIt = 30000, validation accuracy=0.087000
Training complete
When lr=0.100000, reg=0.001000, noIt = 10000, validation accuracy=0.087000

```

```

when lr=0.100000, reg=0.001000, noIt = 10000, validation accuracy=0.087000
Training complete
When lr=0.100000, reg=0.001000, noIt = 20000, validation accuracy=0.087000
Training complete
When lr=0.100000, reg=0.001000, noIt = 30000, validation accuracy=0.087000
SEARCH COMPLETE
Validation accuracy: 0.528

```

```

In [32]: 1 val_acc = (best_net.predict(X_val) == y_val).mean()
          2 print('Validation accuracy: ', val_acc)

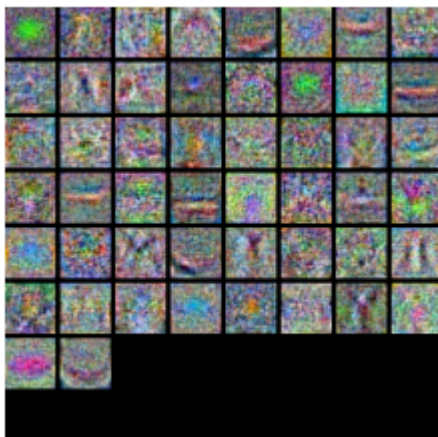
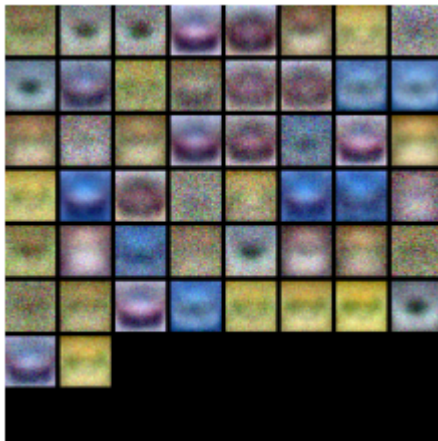
```

Validation accuracy: 0.528

```

In [33]: 1 from utils.vis_utils import visualize_grid
          2
          3 # Visualize the weights of the network
          4
          5 def show_net_weights(net):
          6     W1 = net.params['W1']
          7     W1 = W1.T.reshape(32, 32, 3, -1).transpose(3, 0, 1, 2)
          8     plt.imshow(visualize_grid(W1, padding=3).astype('uint8'))
          9     plt.gca().axis('off')
         10     plt.show()
         11
         12 show_net_weights(subopt_net)
         13 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 of the suboptimal net appear to be noisy and without any distinguishable characteristics. It is not easy to understand the features the weights are looking at and they all look averaged in terms of shape. On the other hand, the weights of the best net offer a clear shape which implies that they learned specific features to look at. The weights are discernible from each other with distinct shapes and it is easy to understand what each weight looks at in an image.

## Evaluate on test set

```
In [34]: 1 test_acc = (best_net.predict(X_test) == y_test).mean()  
         2 print('Test accuracy: ', test_acc)
```

Test accuracy: 0.51

```

import numpy as np
import pdb

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

    The input x has shape (N, d_1, ..., d_k) and contains a minibatch of N
    examples, where each example x[i] has shape (d_1, ..., d_k). We will
    reshape each input into a vector of dimension  $\bar{D} = d_1 * \dots * d_k$ , 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.
    # ===== #
    correctDimX = x.reshape(x.shape[0], -1)
    out = correctDimX.dot(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)

    Returns a tuple of:
    - dx: Gradient with respect to x, of shape (N, d1, ..., d_k)
    - dw: Gradient with respect to w, of shape (D, M)
    - db: Gradient with respect to b, of shape (M,)
    """

```

```

x, w, b = cache
dx, dw, db = None, None, None

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

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

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

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

return dx, dw, db

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

    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.
    # ===== #

    out = np.maximum(x,0)
    # ===== #
    # 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 (ReLU).

    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
correctDimX = x.reshape(x.shape[0], -1)
correctDimX[correctDimX<0]=0
correctDimX[correctDimX>0]=1
dx = dout*correctDimX

# ===== #
# 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  $0 \leq y[i] < C$

Returns a tuple of:

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

"""

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

```

import numpy as np

from .layers import *
from .layer_utils import *

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

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

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

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

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

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

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

        self.params['W1'] = np.random.normal(loc=0.0, scale=weight_scale, size = (input_dim, hidden_dims))
        self.params['b1'] = np.zeros(hidden_dims)
        self.params['W2'] = np.random.normal(loc=0.0, scale=weight_scale, size = (hidden_dims, num_classes))
        self.params['b2'] = np.zeros(num_classes)

        # ===== #

```

```

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

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

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

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

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

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

    a, fc_cache = affine_forward(X, W1, b1)
    out, relu_cache = relu_forward(a)
    cache_hidden = (fc_cache, relu_cache)
    scores, cache_z = affine_forward(out, W2, 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, dz = softmax_loss(scores, y)
loss = loss + 0.5*self.reg*(np.sum(W1**2) + np.sum(W2**2))

dhidden, dw2, db2 = affine_backward(dz, cache_z)
fc_cache, relu_cache = cache_hidden
da = relu_backward(dhidden, relu_cache)
dx, dw1, db1 = affine_backward(da, fc_cache)

grads['W1'] = dw1 + self.reg * W1
grads['b1'] = db1
grads['W2'] = dw2 + self.reg * W2
grads['b2'] = db2
# ===== #
# END YOUR CODE HERE
# ===== #

return loss, grads

```

```
class FullyConnectedNet(object):
```

```

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

```

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

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

```

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

```

```

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

```

```
    """
```

```
    Initialize a new FullyConnectedNet.
```

```
    Inputs:
```

- hidden\_dims: A list of integers giving the size of each hidden layer.
- input\_dim: An integer giving the size of the input.
- num\_classes: An integer giving the number of classes to classify.
- dropout: Scalar between 0 and 1 giving dropout strength. If dropout=0 then the network should not use dropout at all.
- use\_batchnorm: Whether or not the network should use batch normalization.
- reg: Scalar giving L2 regularization strength.



- `weight_scale`: Scalar giving the standard deviation for random initialization of the weights.
- `dtype`: A numpy datatype object; all computations will be performed using this datatype. float32 is faster but less accurate, so you should use float64 for numeric gradient checking.
- `seed`: If not None, then pass this random seed to the dropout layers. This will make the dropout layers deterministic 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`.

# ===== #

```
for i in range(0, self.num_layers):
```

```
    name_W = 'W'+str(i+1)
```

```
    name_b = 'b'+str(i+1)
```

```
    if i == 0: #First
```

```
        self.params[name_W] = np.random.normal(loc=0.0, scale=weight_scale, size = (input
```

```
        self.params[name_b] = np.zeros(hidden_dims[i])
```

```
    elif i == self.num_layers-1: #Last
```

```
        self.params[name_W] = np.random.normal(loc=0.0, scale=weight_scale, size = (hidde
```

```
        self.params[name_b] = np.zeros(num_classes)
```

```
    else: #Between
```

```
        self.params[name_W] = np.random.normal(loc=0.0, scale=weight_scale, size = (hidde
```

```
        self.params[name_b] = np.zeros(hidden_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".
    # ===== #

    H = []
    cache_h = []
    for i in np.arange(0, self.num_layers):
        name_W = 'W'+str(i+1)
        name_b = 'b'+str(i+1)

        if i == 0: #First
            a, fc_cache = affine_forward(X, self.params[name_W], self.params[name_b])
            out, relu_cache = relu_forward(a)
            cH = (fc_cache, relu_cache)
            H.append(out)
            cache_h.append(cH)
        elif i == self.num_layers-1: #Last
            scores = affine_forward(H[i-1], self.params[name_W], self.params[name_b])[0]
            cache_h.append(affine_forward(H[i-1], self.params[name_W], self.params[name_b]))
        else: #Between
            a, fc_cache = affine_forward(H[i-1], self.params[name_W], self.params[name_b])
            out, relu_cache = relu_forward(a)
            cH = (fc_cache, relu_cache)
            H.append(out)
            cache_h.append(cH)

    # ===== #

```

```

# 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, dz = softmax_loss(scores, y)
for i in range(self.num_layers, 0, -1):
    name_W = 'W'+str(i)
    name_b = 'b'+str(i)
    loss = loss + (0.5 * self.reg * np.sum(self.params[name_W]*self.params[name_W]))

    if i == self.num_layers:
        dh1, grads[name_W], grads[name_b] = affine_backward(dz, cache_h[self.num_layers])
    else:
        fc_cache, relu_cache = cache_h[i-1]
        da = relu_backward(dh1, relu_cache)
        dh1, grads[name_W], grads[name_b] = affine_backward(da, fc_cache)

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

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

```

# 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 [1]:

```
1  ## Import and setups
2
3  import time
4  import numpy as np
5  import matplotlib.pyplot as plt
6  from nn1.fc_net import *
7  from utils.data_utils import get_CIFAR10_data
8  from utils.gradient_check import eval_numerical_gradient, eval_numerical_grads
9  from utils.solver import Solver
10
11 %matplotlib inline
12 plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
13 plt.rcParams['image.interpolation'] = 'nearest'
14 plt.rcParams['image.cmap'] = 'gray'
15
16 # for auto-reloading external modules
17 # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
18 %load_ext autoreload
19 %autoreload 2
20
21 def rel_error(x, y):
22     """ returns relative error """
23     return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

In [2]:

```
1  # Load the (preprocessed) CIFAR10 data.
2
3  data = get_CIFAR10_data()
4  for k in data.keys():
5      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 `nnd1/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 [4]: 1 # Test the affine_forward function
2
3 num_inputs = 2
4 input_shape = (4, 5, 6)
5 output_dim = 3
6
7 input_size = num_inputs * np.prod(input_shape)
8 weight_size = output_dim * np.prod(input_shape)
9
10 x = np.linspace(-0.1, 0.5, num=input_size).reshape(num_inputs, *input_shape)
11 w = np.linspace(-0.2, 0.3, num=weight_size).reshape(np.prod(input_shape), ou
12 b = np.linspace(-0.3, 0.1, num=output_dim)
13
14 out, _ = affine_forward(x, w, b)
15 correct_out = np.array([[ 1.49834967,  1.70660132,  1.91485297],
16                          [ 3.25553199,  3.5141327,  3.77273342]])
17
18 # Compare your output with ours. The error should be around 1e-9.
19 print('Testing affine_forward function:')
20 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 [5]: 1 # Test the affine_backward function
2
3 x = np.random.randn(10, 2, 3)
4 w = np.random.randn(6, 5)
5 b = np.random.randn(5)
6 dout = np.random.randn(10, 5)
7
8 dx_num = eval_numerical_gradient_array(lambda x: affine_forward(x, w, b)[0],
9 dw_num = eval_numerical_gradient_array(lambda w: affine_forward(x, w, b)[0],
10 db_num = eval_numerical_gradient_array(lambda b: affine_forward(x, w, b)[0],
11
12 _, cache = affine_forward(x, w, b)
13 dx, dw, db = affine_backward(dout, cache)
14
15 # The error should be around 1e-10
16 print('Testing affine_backward function:')
17 print('dx error: {}'.format(rel_error(dx_num, dx)))
18 print('dw error: {}'.format(rel_error(dw_num, dw)))
19 print('db error: {}'.format(rel_error(db_num, db)))
```

```
Testing affine_backward function:
dx error: 3.8090866555111607e-11
dw error: 1.7375717532082819e-10
db error: 4.787050739675008e-10
```

## Activation layers

In this section you'll implement the ReLU activation.

### ReLU forward pass

Implement the `relu_forward` function in `nnd1/layers.py` and then test your code by running the following cell.

```
In [6]: 1 # Test the relu_forward function
2
3 x = np.linspace(-0.5, 0.5, num=12).reshape(3, 4)
4
5 out, _ = relu_forward(x)
6 correct_out = np.array([[ 0.,          0.,          0.,          0.,
7                          [ 0.,          0.,          0.04545455, 0.13636364,
8                          [ 0.22727273, 0.31818182, 0.40909091, 0.5,
9
10 # Compare your output with ours. The error should be around 1e-8
11 print('Testing relu_forward function:')
12 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 `nnd1/layers.py` and then test your code by running the following cell.

```
In [7]: 1 x = np.random.randn(10, 10)
2 dout = np.random.randn(*x.shape)
3
4 dx_num = eval_numerical_gradient_array(lambda x: relu_forward(x)[0], x, dout
5
6 _, cache = relu_forward(x)
7 dx = relu_backward(dout, cache)
8
9 # The error should be around 1e-12
10 print('Testing relu_backward function:')
11 print('dx error: {}'.format(rel_error(dx_num, dx)))
```

```
Testing relu_backward function:
dx error: 3.2756336458611084e-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 `nnd1/layer_utils.py`.

### Affine-ReLU layers

We've implemented `affine_relu_forward()` and `affine_relu_backward` in `nnd1/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 [11]: 1 from nndl.layer_utils import affine_relu_forward, affine_relu_backward
2
3 x = np.random.randn(2, 3, 4)
4 w = np.random.randn(12, 10)
5 b = np.random.randn(10)
6 dout = np.random.randn(2, 10)
7
8 out, cache = affine_relu_forward(x, w, b)
9 dx, dw, db = affine_relu_backward(dout, cache)
10
11 dx_num = eval_numerical_gradient_array(lambda x: affine_relu_forward(x, w, b)
12 dw_num = eval_numerical_gradient_array(lambda w: affine_relu_forward(x, w, b)
13 db_num = eval_numerical_gradient_array(lambda b: affine_relu_forward(x, w, b)
14
15 print('Testing affine_relu_forward and affine_relu_backward:')
16 print('dx error: {}'.format(rel_error(dx_num, dx)))
17 print('dw error: {}'.format(rel_error(dw_num, dw)))
18 print('db error: {}'.format(rel_error(db_num, db)))
```

Testing affine\_relu\_forward and affine\_relu\_backward:

dx error: 6.382280630077004e-11

dw error: 7.413924268923301e-10

db error: 2.1600192477321958e-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 [9]: 1 num_classes, num_inputs = 10, 50
2 x = 0.001 * np.random.randn(num_inputs, num_classes)
3 y = np.random.randint(num_classes, size=num_inputs)
4
5
6
7 dx_num = eval_numerical_gradient(lambda x: softmax_loss(x, y)[0], x, verbose
8 loss, dx = softmax_loss(x, y)
9
10 # Test softmax_loss function. Loss should be 2.3 and dx error should be 1e-8
11 print('\nTesting softmax_loss:')
12 print('loss: {}'.format(loss))
13 print('dx error: {}'.format(rel_error(dx_num, dx)))
```

Testing softmax\_loss:

loss: 2.3026204658474767

dx error: 8.120600213352476e-09

## Implementation of a two-layer NN

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

In [20]:

```

1 N, D, H, C = 3, 5, 50, 7
2 X = np.random.randn(N, D)
3 y = np.random.randint(C, size=N)
4
5 std = 1e-2
6 model = TwoLayerNet(input_dim=D, hidden_dims=H, num_classes=C, weight_scale=
7
8 print('Testing initialization ... ')
9 W1_std = abs(model.params['W1'].std() - std)
10 b1 = model.params['b1']
11 W2_std = abs(model.params['W2'].std() - std)
12 b2 = model.params['b2']
13 assert W1_std < std / 10, 'First layer weights do not seem right'
14 assert np.all(b1 == 0), 'First layer biases do not seem right'
15 assert W2_std < std / 10, 'Second layer weights do not seem right'
16 assert np.all(b2 == 0), 'Second layer biases do not seem right'
17
18 print('Testing test-time forward pass ... ')
19 model.params['W1'] = np.linspace(-0.7, 0.3, num=D*H).reshape(D, H)
20 model.params['b1'] = np.linspace(-0.1, 0.9, num=H)
21 model.params['W2'] = np.linspace(-0.3, 0.4, num=H*C).reshape(H, C)
22 model.params['b2'] = np.linspace(-0.9, 0.1, num=C)
23 X = np.linspace(-5.5, 4.5, num=N*D).reshape(D, N).T
24 scores = model.loss(X)
25 correct_scores = np.asarray(
26     [[11.53165108, 12.2917344, 13.05181771, 13.81190102, 14.57198434, 15.
27      [12.05769098, 12.74614105, 13.43459113, 14.1230412, 14.81149128, 15.
28      [12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.05099822, 15.
29 scores_diff = np.abs(scores - correct_scores).sum()
30 assert scores_diff < 1e-6, 'Problem with test-time forward pass'
31
32 print('Testing training loss (no regularization)')
33 y = np.asarray([0, 5, 1])
34 loss, grads = model.loss(X, y)
35 correct_loss = 3.4702243556
36 assert abs(loss - correct_loss) < 1e-10, 'Problem with training-time loss'
37
38 model.reg = 1.0
39 loss, grads = model.loss(X, y)
40 correct_loss = 26.5948426952
41 assert abs(loss - correct_loss) < 1e-10, 'Problem with regularization loss'
42
43 for reg in [0.0, 0.7]:
44     print('Running numeric gradient check with reg = {}'.format(reg))
45     model.reg = reg
46     loss, grads = model.loss(X, y)
47
48     for name in sorted(grads):
49         f = lambda _: model.loss(X, y)[0]
50         grad_num = eval_numerical_gradient(f, model.params[name], verbose=False)
51         print('{} relative error: {}'.format(name, rel_error(grad_num, grads[nam

```

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 `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 [21]:

```

1 model = TwoLayerNet()
2 solver = None
3
4 # ===== #
5 # YOUR CODE HERE:
6 #   Declare an instance of a TwoLayerNet and then train
7 #   it with the Solver. Choose hyperparameters so that your validation
8 #   accuracy is at least 50%. We won't have you optimize this further
9 #   since you did it in the previous notebook.
10 #
11 # ===== #
12
13 model = TwoLayerNet(hidden_dims=200, reg = 0.1)
14 solver = Solver(model, data, update_rule='sgd',
15                 optim_config={
16                     'learning_rate': 1e-3,
17                 }, lr_decay=0.95, num_epochs=10, batch_size=200, print_eve
18 solver.train()
19
20 # ===== #
21 # END YOUR CODE HERE
22 # ===== #

```

```

(Iteration 1 / 2450) loss: 2.333385
(Epoch 0 / 10) train acc: 0.126000; val_acc: 0.150000
(Iteration 101 / 2450) loss: 1.642614
(Iteration 201 / 2450) loss: 1.757910
(Epoch 1 / 10) train acc: 0.426000; val_acc: 0.420000
(Iteration 301 / 2450) loss: 1.565339
(Iteration 401 / 2450) loss: 1.498052
(Epoch 2 / 10) train acc: 0.515000; val_acc: 0.484000
(Iteration 501 / 2450) loss: 1.501815
(Iteration 601 / 2450) loss: 1.357544
(Iteration 701 / 2450) loss: 1.434842
(Epoch 3 / 10) train acc: 0.515000; val_acc: 0.496000
(Iteration 801 / 2450) loss: 1.401399
(Iteration 901 / 2450) loss: 1.583114
(Epoch 4 / 10) train acc: 0.554000; val_acc: 0.496000
(Iteration 1001 / 2450) loss: 1.467264
(Iteration 1101 / 2450) loss: 1.431473
(Iteration 1201 / 2450) loss: 1.303135
(Epoch 5 / 10) train acc: 0.554000; val_acc: 0.517000
(Iteration 1301 / 2450) loss: 1.271931
(Iteration 1401 / 2450) loss: 1.295709
(Epoch 6 / 10) train acc: 0.545000; val_acc: 0.501000
(Iteration 1501 / 2450) loss: 1.234075
(Iteration 1601 / 2450) loss: 1.277352
(Iteration 1701 / 2450) loss: 1.159497
(Epoch 7 / 10) train acc: 0.589000; val_acc: 0.498000
(Iteration 1801 / 2450) loss: 1.199592
(Iteration 1901 / 2450) loss: 1.333886
(Epoch 8 / 10) train acc: 0.565000; val_acc: 0.488000
(Iteration 2001 / 2450) loss: 1.206221
(Iteration 2101 / 2450) loss: 1.153008
(Iteration 2201 / 2450) loss: 1.252871
(Epoch 9 / 10) train acc: 0.578000; val_acc: 0.499000
(Iteration 2301 / 2450) loss: 1.359032

```

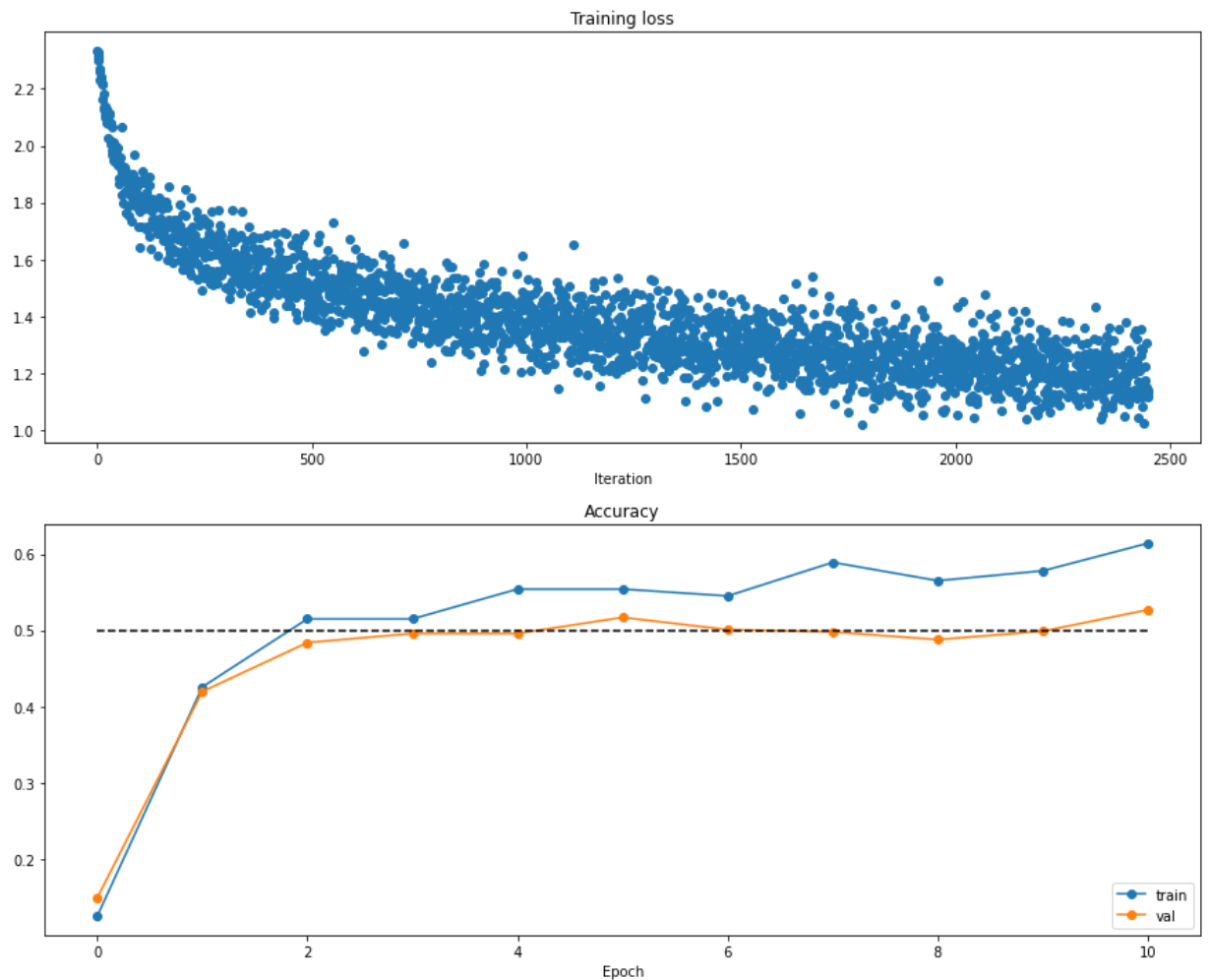
(Iteration 2401 / 2450) loss: 1.185255

(Epoch 10 / 10) train acc: 0.614000; val\_acc: 0.527000

```

In [22]: 1 # Run this cell to visualize training loss and train / val accuracy
2
3 plt.subplot(2, 1, 1)
4 plt.title('Training loss')
5 plt.plot(solver.loss_history, 'o')
6 plt.xlabel('Iteration')
7
8 plt.subplot(2, 1, 2)
9 plt.title('Accuracy')
10 plt.plot(solver.train_acc_history, '-o', label='train')
11 plt.plot(solver.val_acc_history, '-o', label='val')
12 plt.plot([0.5] * len(solver.val_acc_history), 'k--')
13 plt.xlabel('Epoch')
14 plt.legend(loc='lower right')
15 plt.gcf().set_size_inches(15, 12)
16 plt.show()

```



## Multilayer Neural Network

Now, we implement a multi-layer neural network.

Read through the `FullyConnectedNet` class in the file `nnd1/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 [23]: 1 N, D, H1, H2, C = 2, 15, 20, 30, 10
2 X = np.random.randn(N, D)
3 y = np.random.randint(C, size=(N,))
4
5 for reg in [0, 3.14]:
6     print('Running check with reg = {}'.format(reg))
7     model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
8                               reg=reg, weight_scale=5e-2, dtype=np.float64)
9
10    loss, grads = model.loss(X, y)
11    print('Initial loss: {}'.format(loss))
12
13    for name in sorted(grads):
14        f = lambda _: model.loss(X, y)[0]
15        grad_num = eval_numerical_gradient(f, model.params[name], verbose=False,
16        print('{} relative error: {}'.format(name, rel_error(grad_num, grads[nam
```

```
Running check with reg = 0
Initial loss: 2.29898895452082
W1 relative error: 7.834318374346021e-08
W2 relative error: 1.8468443517067148e-06
W3 relative error: 1.5479893849745557e-06
b1 relative error: 9.868988656484848e-08
b2 relative error: 9.025216577061229e-09
b3 relative error: 1.1983247432319214e-10
Running check with reg = 3.14
Initial loss: 7.050918253913044
W1 relative error: 2.0816416576698898e-08
W2 relative error: 2.4808336653037244e-08
W3 relative error: 6.719881886242843e-09
b1 relative error: 4.6631644980123366e-08
b2 relative error: 2.171782017826178e-09
b3 relative error: 1.9404166927610103e-10
```

```

In [25]: 1 # Use the three layer neural network to overfit a small dataset.
2
3 num_train = 50
4 small_data = {
5     'X_train': data['X_train'][:num_train],
6     'y_train': data['y_train'][:num_train],
7     'X_val': data['X_val'],
8     'y_val': data['y_val'],
9 }
10
11
12 ##### !!!!!!!
13 # Play around with the weight_scale and learning_rate so that you can overfi
14 # Your training accuracy should be 1.0 to receive full credit on this part.
15 weight_scale = 1e-2
16 learning_rate = 1e-2
17
18 model = FullyConnectedNet([100, 100],
19                             weight_scale=weight_scale, dtype=np.float64)
20 solver = Solver(model, small_data,
21                 print_every=10, num_epochs=20, batch_size=25,
22                 update_rule='sgd',
23                 optim_config={
24                     'learning_rate': learning_rate,
25                 })
26
27 solver.train()
28
29 plt.plot(solver.loss_history, 'o')
30 plt.title('Training loss history')
31 plt.xlabel('Iteration')
32 plt.ylabel('Training loss')
33 plt.show()

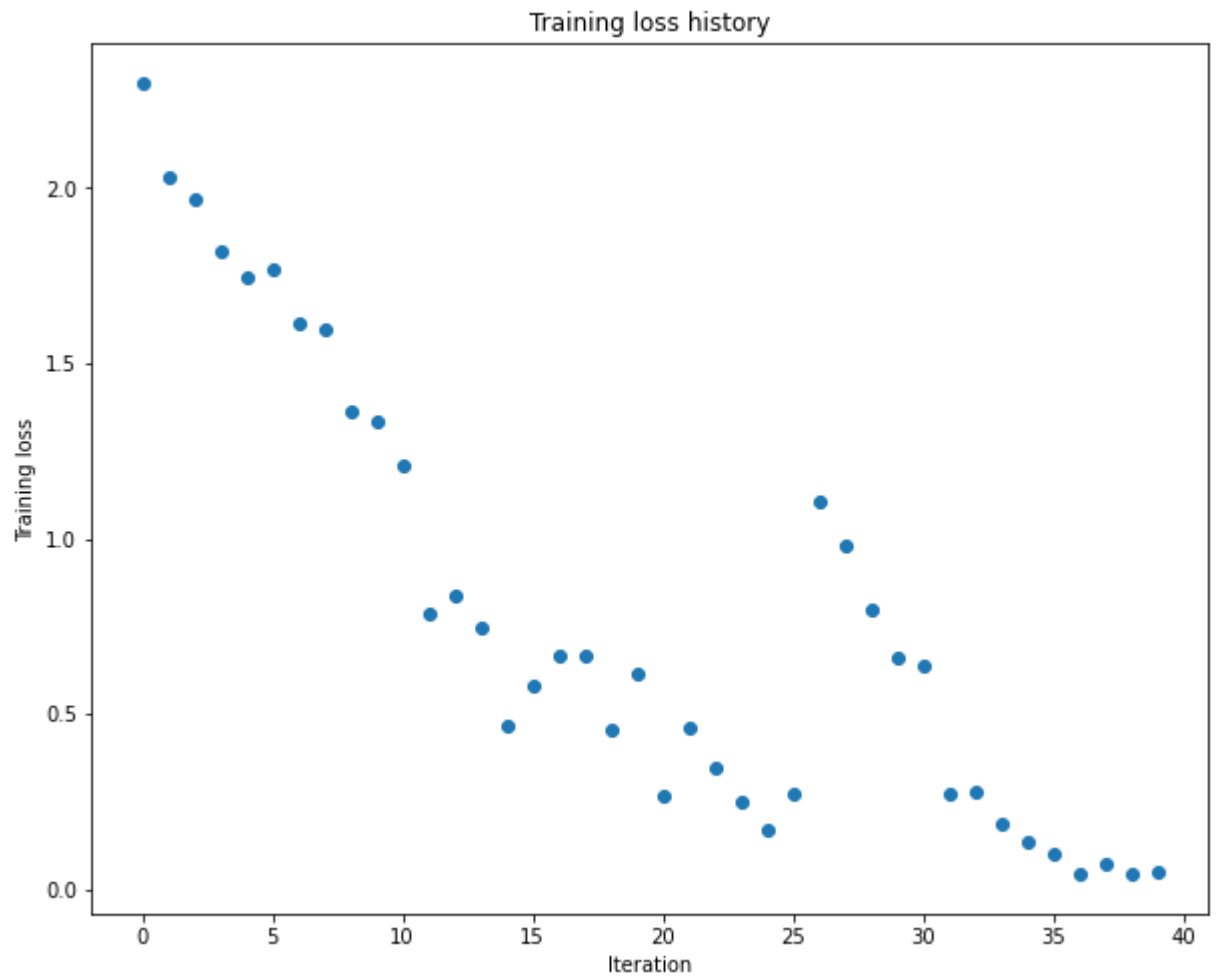
```

```

(Iteration 1 / 40) loss: 2.300841
(Epoch 0 / 20) train acc: 0.300000; val_acc: 0.118000
(Epoch 1 / 20) train acc: 0.340000; val_acc: 0.125000
(Epoch 2 / 20) train acc: 0.420000; val_acc: 0.163000
(Epoch 3 / 20) train acc: 0.480000; val_acc: 0.161000
(Epoch 4 / 20) train acc: 0.580000; val_acc: 0.144000
(Epoch 5 / 20) train acc: 0.780000; val_acc: 0.188000
(Iteration 11 / 40) loss: 1.210870
(Epoch 6 / 20) train acc: 0.820000; val_acc: 0.177000
(Epoch 7 / 20) train acc: 0.880000; val_acc: 0.173000
(Epoch 8 / 20) train acc: 0.880000; val_acc: 0.194000
(Epoch 9 / 20) train acc: 0.880000; val_acc: 0.185000
(Epoch 10 / 20) train acc: 0.960000; val_acc: 0.163000
(Iteration 21 / 40) loss: 0.264498
(Epoch 11 / 20) train acc: 0.920000; val_acc: 0.165000
(Epoch 12 / 20) train acc: 0.960000; val_acc: 0.176000
(Epoch 13 / 20) train acc: 0.900000; val_acc: 0.168000
(Epoch 14 / 20) train acc: 0.760000; val_acc: 0.152000
(Epoch 15 / 20) train acc: 0.760000; val_acc: 0.170000
(Iteration 31 / 40) loss: 0.635391
(Epoch 16 / 20) train acc: 0.940000; val_acc: 0.167000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.160000

```

(Epoch 18 / 20) train acc: 1.000000; val\_acc: 0.174000  
(Epoch 19 / 20) train acc: 1.000000; val\_acc: 0.173000  
(Epoch 20 / 20) train acc: 1.000000; val\_acc: 0.166000



In [ ]: 1

In [ ]:

1