

EEE-443 Neural Networks Project-2

Name: Ayberk Yarkin

Surname: Yıldız

Id: 21803386

Date: November 29, 2021

Question 1

In this question, it is asked to implement a neural network that classifies the images either to cat or car using Stochastic Gradient Descent on mini-batches. Mean squared and mean classification errors are used. Mean squared error (MSE) is defined below:

$$MSE = \frac{1}{N} \sum_{i=1}^N (y^i - \hat{y}^i)^2$$

where there is a total of N samples to be calculated.

Mean classification error (MCE) is determined with the percentage of correctly classified images.

a) In part a, it is asked to design a multi-layer neural network with a single hidden and output layer using backpropagation algorithm. It is assumed a hyperbolic tangent activation function for all neurons. The tanh activation function can be derived below:

$$y = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Although activation function can be derived as the formula above, I used tanh function of numpy in my code. The derivative of the activation function that will be used in backpropagation is shown below:

$$\frac{\partial \tanh(x)}{\partial x} = 1 - \tanh^2(x) = 1 - y^2$$

The update formula of the weights is shown below:

$$\Delta W = \eta * \delta * X^T$$

where the delta $\delta = \frac{\partial \tanh(x)}{\partial x} \odot \frac{\partial E}{\partial Y}$ is calculated individually for hidden and output layers, where the error is calculated as:

$$E = \frac{1}{2} (Y - \hat{Y})$$

and X is the input for the hidden or the output layer.

My optimized values for the network are shown below:

Parameter	Value
Hidden Neuron Number	30
Output Neuron Number	1
Learning Rate	0.4
Batch Size	50
Std of Weights	0.01
Epochs	300

Table 1: Optimized Parameters of Q1-part a

I calculated and plotted the mean squared and mean classification errors as follows:

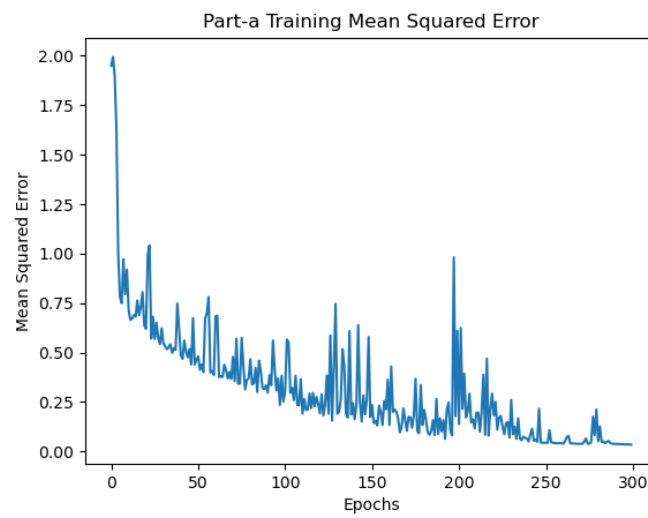


Figure 1: Training Mean Squared Error

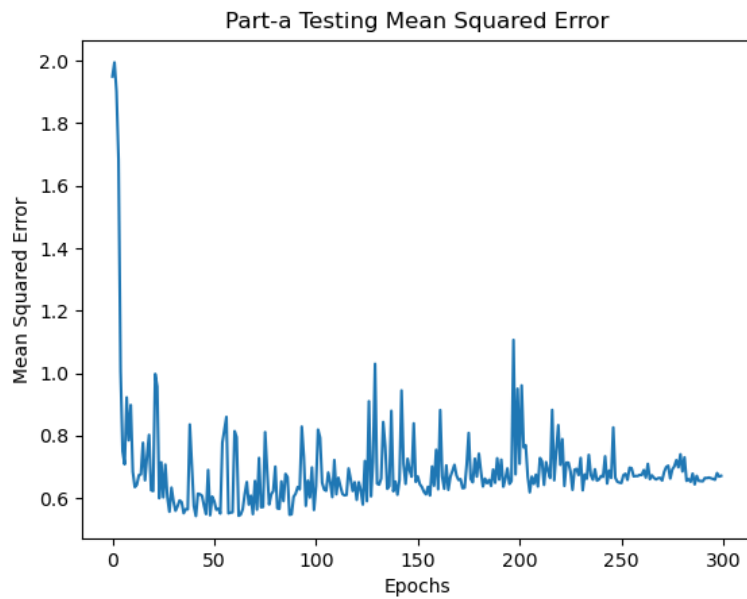


Figure 2: Testing Mean Squared Error

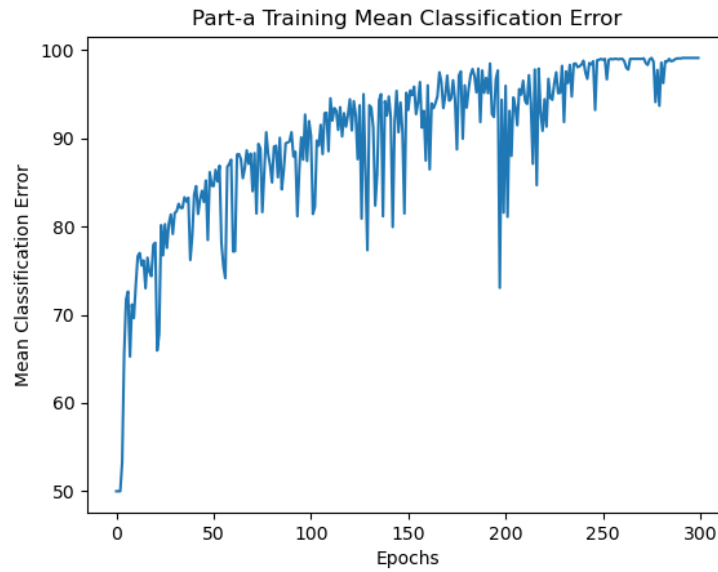


Figure 3: Training Mean Classification Error

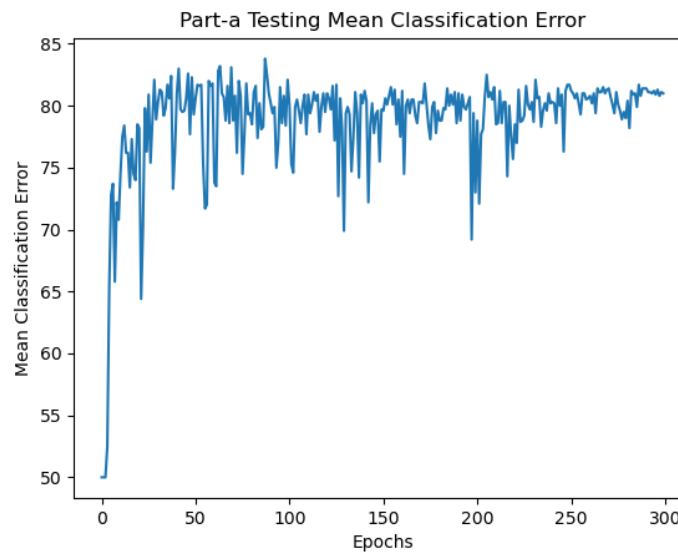


Figure 4: Testing Mean Classification Error

It can be seen that training MSE converges to zero as epoch increases, and testing MSE is converges to nearly between 0.6 and 0.7.

The training MCE converges to nearly 100% and testing MCE converges to 81%. This variance occur due to the variance in the data that the learning the data does not specifically mean perfect classification for the test data.

b) There is fast convergence in MSE plots due to Stochastic Gradient Descent. Figure 1 shows fast convergence that it converges around 200-250 epochs, which is a relative small epoch number. There occurs spikes and noise since the weights are characterized in batches and not individually. The training accuracy converges to nearly 100% and the testing accuracy converges around 81%. The training time also decreased in SGD because it does not propagate and update weights at each data point.

MSE may not be preferable for classification problems that there can be cases that classifications for the classes may be accurate but the predictions may be inaccurate and not converging to the ground truth. MSE may not show good succeeded results in these cases. Loss of information may occur after the classification of the predictions. MSE can be useful to measure the network in whether the algorithm is learning or not, but it is not an accurate error measurement for a classification problem.

c) In this part, the same classification problem is considered, but with a higher and lower number of hidden neurons. For the higher number of neurons, I chose the neuron number as 150, and for the lower I chose 5. Plots of MSE and MCE for training and test for the network with high, low and the optimal number of hidden neurons from part a is shown below:

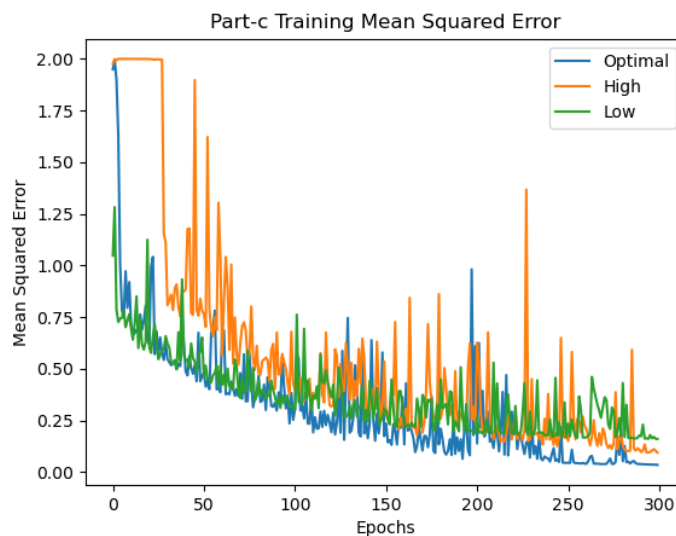


Figure 5: Training Mean Squared Error for Optimal, High and Low Number of Hidden Neurons

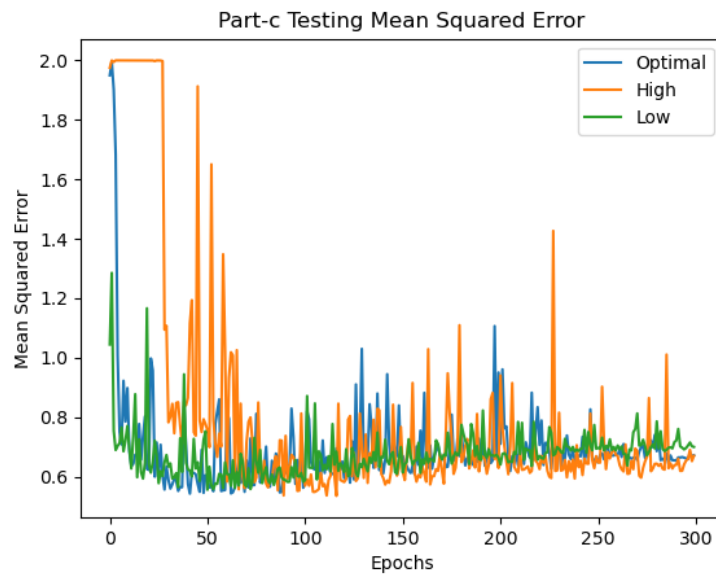


Figure 6: Testing Mean Squared Error for Optimal, High and Low Number of Hidden Neurons

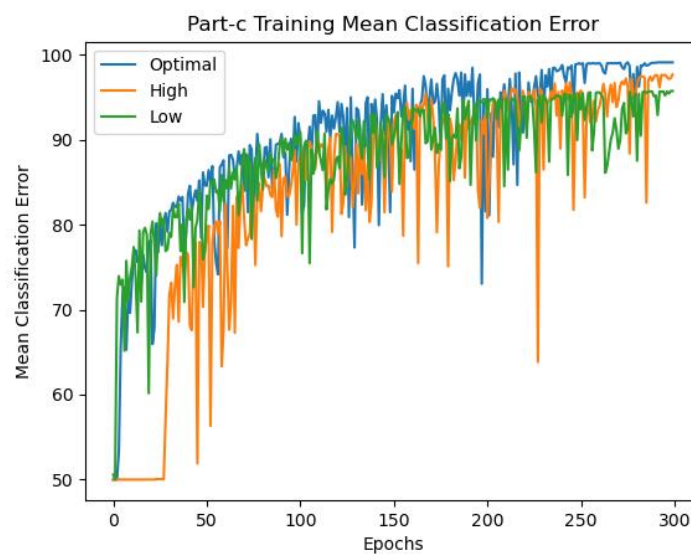


Figure 7: Training Mean Classification Error for Optimal, High and Low Number of Hidden Neurons

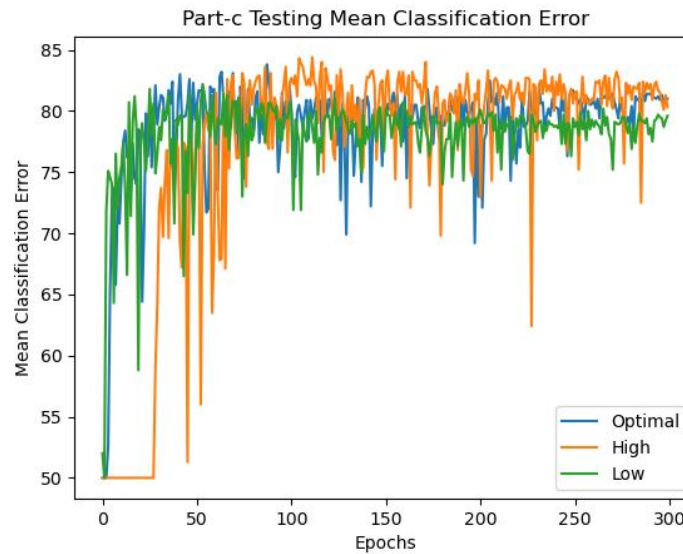


Figure 8: Testing Mean Classification Error for Optimal, High and Low Number of Hidden Neurons

It can be seen that the networks with higher and lower number of hidden units suffer more noise and their convergence is more later compared to the optimal case. Even their converged accuracies are less than the accuracy of optimal network. It can be caused by underfitting and overfitting cases for lower and higher number of hidden unit networks respectively. If the hidden unit number is high, the oscillations increase due to overfitting of the past train. When the hidden unit number is low, the network will struggle classifying the train data and cannot predict the test data correctly. I did not want to choose the higher neuron number much higher because of the increased training time and processing power, therefore I selected it nearly 5 times of the optimal value. When we go right or left from the optimal point, the loss increases.

d) In this part, the network has two hidden layers instead of one. My optimized values for 2 hidden layer network is shown below:

Parameter	Value
Hidden Neuron 1st Layer	300
Hidden Neuron 2nd Layer	30
Output Neuron Number	1
Learning Rate	0.4
Batch Size	50
Std of Weights	0.01
Epochs	300

Table 2: Optimized Parameters of Q1-part d

I calculated and plotted the mean squared and mean classification errors as follows:

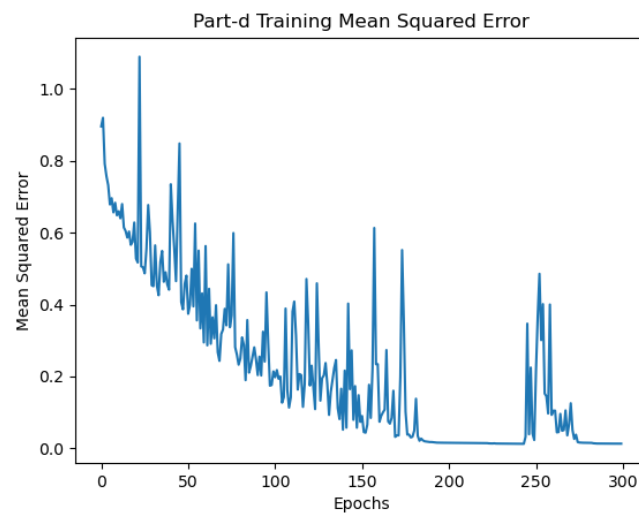


Figure 9: Training Mean Squared Error with Two Hidden Layers

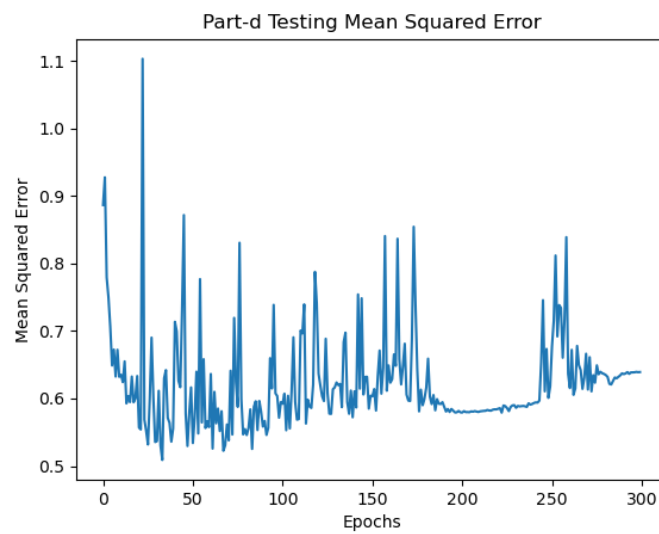


Figure 10: Testing Mean Squared Error with Two Hidden Layers

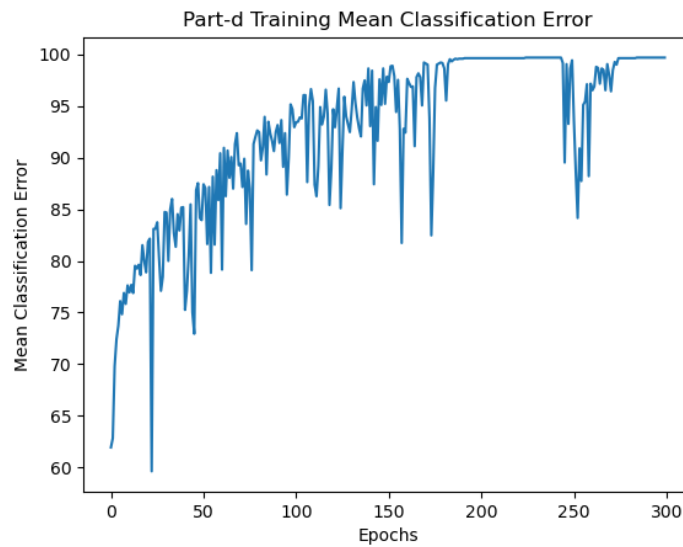


Figure 11: Training Mean Classification Error with Two Hidden Layers

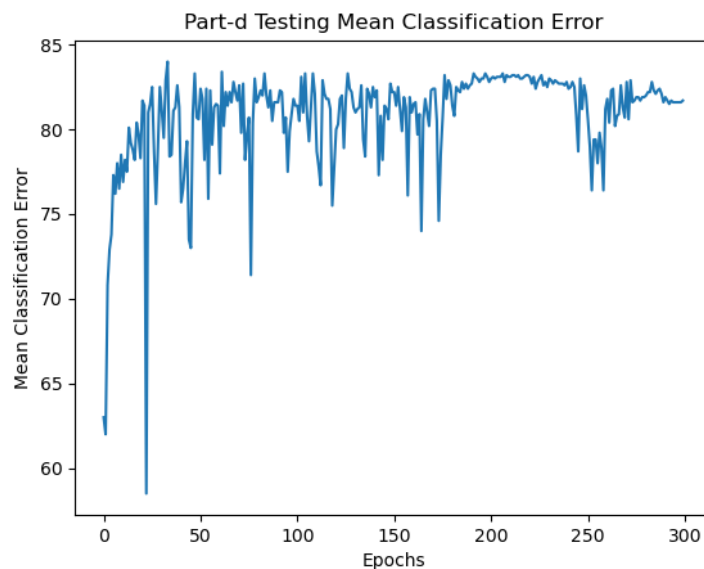


Figure 12: Testing Mean Classification Error with Two Hidden Layers

Comparing the results with part a, there is not a huge change between the learning curves among all metrics. All plots converges to the similar point with the ones in part a. The reason for this is the simplicity of the problem. Since the problem is not complex enough to require a network with 2 hidden layers, the training accuracy still converges nearly 100% but not reaches 100% exactly. Additionally for MSE, there are more spikes in two hidden layer network compared to the network in part a. For more complex problems, a network with

multiple hidden layers is required, but not in this case. Similar to the explanations in part c, if the number of layers increases, there may occur underfitting/overfitting.

e) Similar to the part d, in this part, still a network with 2 hidden layers is used, but additionally a momentum coefficient is used. The iterative momentum formula is shown below:

$$\Delta W(n) = -\eta \frac{\partial E}{\partial W} + \alpha \Delta W(n-1)$$

The parameters are same with the network in part d except the momentum coefficient, which I chose as $\alpha = 0.5$. The results are shown below:

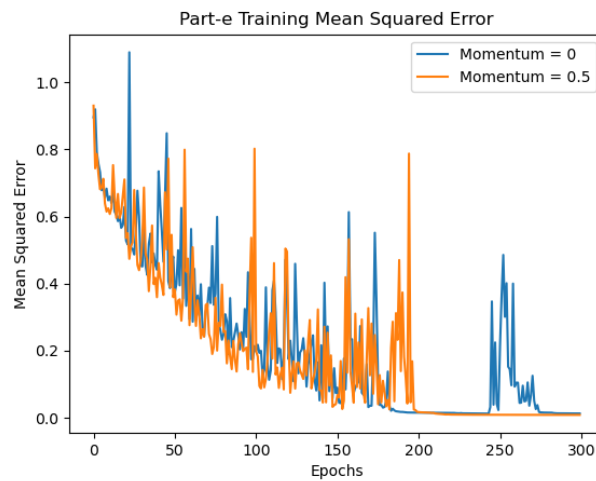


Figure 13: Training Mean Squared Error with Two Hidden Layers with and without Momentum

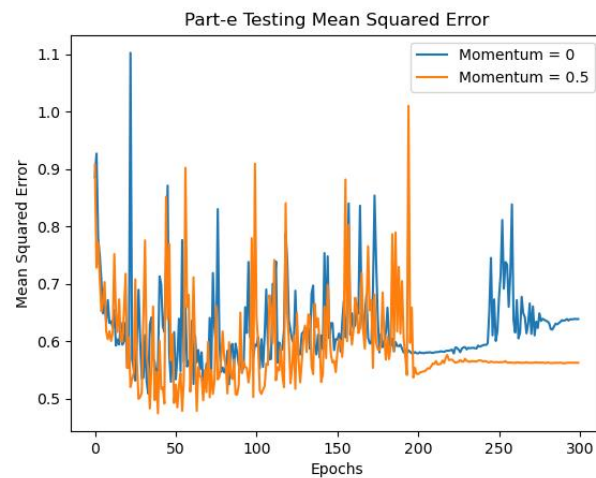


Figure 14: Testing Mean Squared Error with Two Hidden Layers with and without Momentum

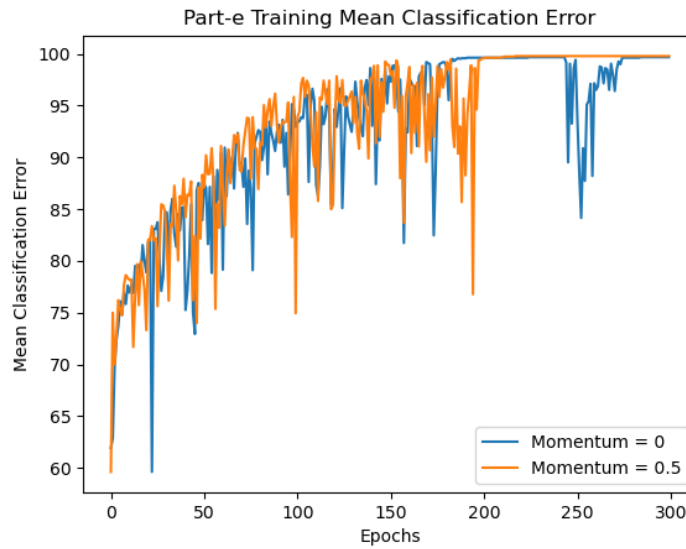


Figure 15: Training Mean Classification Error with Two Hidden Layers with and without Momentum

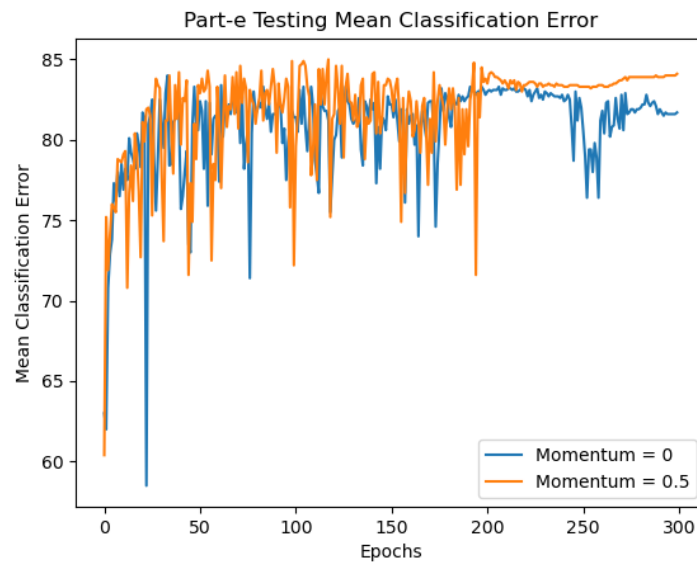


Figure 16: Testing Mean Classification Error with Two Hidden Layers with and without Momentum

In general, the momentum coefficient helps the network to find the global minimum, since the standard gradient descent algorithm may find the local minimum rather. It is prevented for the network to stuck in a local minimum. It helps in faster convergence.

For the results, it can be observed from the plots that with a momentum coefficient, network converges faster compared to the no momentum coefficient case. I chose the momentum coefficient as 0.5 to see the effect of it clearly. It provides a more stable

convergence. Additionally, the network with momentum has a higher classification accuracy, which is nearly 84%, and a lower mean square error for test data compared to the network with no momentum. However, since the task is simple enough as stated before, the effect of momentum in a network can be more clearly seen in a more complex problem. It is relatively less observable in this problem.

Question 2

In question 2, it is asked to implement a network that predicts the 4th Word from a given 3-word sequence. Sigmoid activation is used for the hidden layer and softmax function is used for the output. Their formulas are shown below:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$
$$\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^p e^{z_j}}$$

where, p is the number of neurons in the layer.

Cross entropy error is used for loss function, whose formula is:

$$E_{\text{cross entropy}} = - \sum_{j=1}^M y_j \log(\text{pred}_j)$$

There are 250 words, so the p value is given as 250.

a) The desired parameters, but can be updated for better performance, are listed below:

- Mini batch size = 200
- $\eta = 0.15$
- $\alpha = 0.85$
- Max. Epoch size = 50
- Weight and bias std deviations = 0.01
- (D,P) = (32,256), (16,128), (8,64)

In the structure, a total 3 layers are used, 1 output and 2 hidden layers. I used a linear activation function ($f(x) = x$) for the embedding part and created another layer for it since it simply linearly projects the input to the output.

The methodology is inputting the samples and then concatenating them for the outputs for separate results column-wise, then inputting it to the network.

The optimized parameters are shown at the table below:

Parameter	Value
Momentum Coefficient	0.85
Learning Rate	0.4
Batch Size	250
Std of Weights	0.1
Epochs	50

Table 3: Optimized Parameters of Q2

Results are shown below:

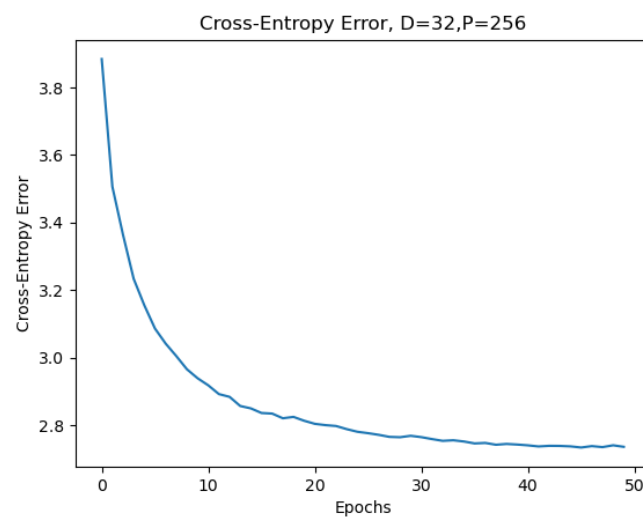


Figure 17: Cross-Entropy Error on the Validation Data when D = 32, P = 256

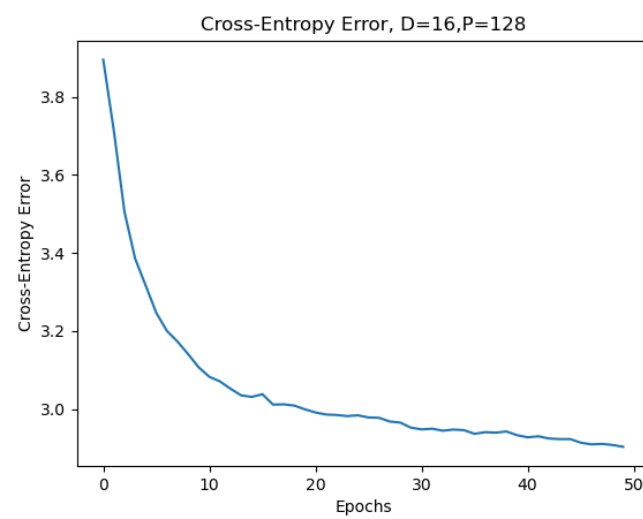


Figure 18: Cross-Entropy Error on the Validation Data when D = 16, P = 128

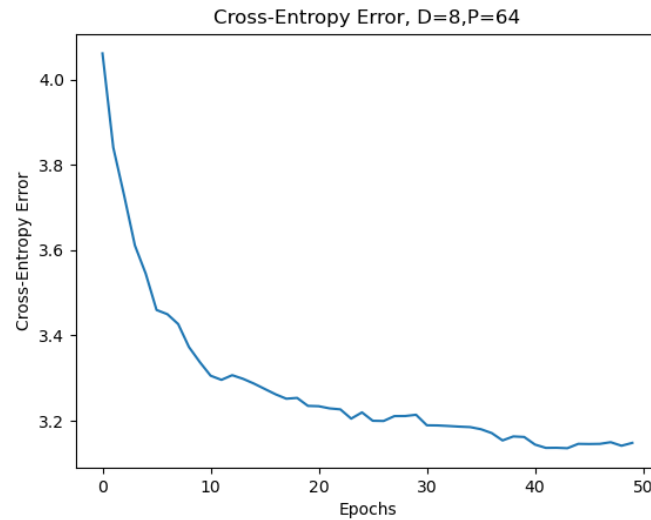


Figure 19: Cross-Entropy Error on the Validation Data when D = 8, P = 64

It can be observed that with more hidden units, the more the network performs better due to the complexity of the problem. The best performance is given at when D=32 and P=256, which has a cross-entropy error nearly 2.75. The errors and accuracies can be seen at the table below:

(D,P) Values	Cross-Entropy Loss	Accuracy
(8,64)	3.15	29.10%
(16,128)	2.90	32.54%
(32,256)	2.74	35.37%

Table 4: Cross-Entropy Losses and Accuracies for different (D,P) values

It is expected to have the better performance since due to the complexity of the problem and large size of the dataset, increasing the hidden unit size increases the performance.

b) In this part, it is asked to take 5 random test samples from the test data and forward it to the trained network to test the network. Afterwards, 10 predictions are executed to test whether they are accurate or not. The highest performed network, which has D = 32 and P = 256, is used for this purpose. The predictions can be seen below:

Phrases	Label	Predictions									
, he said	.	last	.	to	he	today	yesterday	:	,	?	for
put the other	one	in	i	way	,	two	people	day	one	team	.
you get on	with	it	.	your	him	you	with	this	?	the	and
nt want to	see	get	come	say	think	know	.	see	be	do	go
good in those	days	to	that	place	years	school	,	.	days	?	times

Table 5: Top 10 predictions of the 3 word phrases

For predictions, since the accuracy is around 30%, which is not a high accuracy, network slightly struggles while prediction. There are tons of possibilities for the 4th word of the sequence in terms of different areas and contexts. Since the network is trained with one possible outcome, when another possibility of outcome in a different context comes to the network, network fails to predict the correct label and predicts the word that it learned before. However, some predictions make sense and shows that even the accuracy is low, network still learns slightly. Additionally, some predicted words are actually occur as punctuations. Network could not differentiate the words and the punctuations. This is again because of the reason that network does not understand the context clearly and may predict the most repeated word in sequence, that may be the punctuations.

Although the predictions do not look consistent, in all 5 phrase tests, the network succeeded in finding the true label word in one of its 10 predictions. I bolded the true predictions among the 10 predicted values in the table 5. Therefore, as a result, network struggles to find the true labels exactly, but does learn and tries to make accurate predictions. This also may happen due to the input word amount, which is 3. It leaves less context for the network to learn the process.

Question 3

Inline questions for the notebooks are at the end of the report, inside the notebooks.

a) In this FullyConnectedNets Jupyter notebook, a fully connected network is implemented. It is executed by dividing the propagation steps. It leads to implementing the network on different applications. In the process, the forward propagation is the prediction step and outputs activation potentials are outputs. Gradients are found in backpropagation. These two steps are implemented for ReLU function. Afterwards, the sandwich layer is introduced. Then, the softmax and SVM loss functions are tested to confirm that the implementations are accurate and correct. A network with two layers is created by coding these various functions and tested by Solver class. Then, a 3 layer and a 5 layer networks are tested by overfitting 50 samples over 20 epochs. The update rule is implemented with Stochastic Gradient Descent, similar to our problems. Then, the update rules of RMSProp and Adam are introduced, which changes the learning rate for each parameter in the system. Finally, the network is trained on the CIFAR-10 dataset using the Solver class, which results in an accuracy of 51%.

b) In this Dropout notebook, it gives an introduction to dropout, which is a regularization technique used to avoid overfitting and individual learning among neurons. $1-p$ neurons become zero. The implementation of the dropout forward and backward modules is created

by creating a mask that creates a random array with same shape with input matrix. This matrix gives 1 if the matrix elements are less than p or 0 otherwise, then they are divided by p . In the training mode, $1-p$ terms are dropped out from the network, whereas in test mode all nodes are presented. Averages are also shown for different p values and how many terms are 0. Then, two fully connected networks with or without dropout are compared and accuracies are calculated and shown.

Appendix

**Script That is Run in Command Prompt via “python ayberk_yarkin_yildiz_21803386.py x”,
For Both Question 1 and Question 2**

```
import sys  
import numpy as np  
import matplotlib.pyplot as plt  
import h5py
```

```
def graph(x):  
    plt.figure()  
    plt.plot(x)  
    plt.xlabel('Epochs')
```

```
def graphm(x,y,z):  
    plt.figure()  
    plt.plot(x)  
    plt.plot(y)  
    plt.plot(z)  
    plt.legend(['Optimal', 'High', 'Low'])  
    plt.xlabel('Epochs')
```

```
def graphd(x,y):  
    plt.figure()  
    plt.plot(x)  
    plt.plot(y)
```

```
plt.legend(['Momentum = 0','Momentum = 0.5'])  
plt.xlabel('Epochs')
```

```
def checkndim(x):  
    if x.ndim == 1:  
        x = x.reshape(x.shape[0],1)  
    return x
```

```
def vec(x,y):  
    nextt = np.zeros(y)  
    nextt[x-1] = 1  
    return nextt
```

```
def m1(x,y):  
    nextt = np.zeros((x.shape[0],x.shape[1],y))  
    for a in range(x.shape[0]):  
        for b in range(x.shape[1]):  
            nextt[a,b,:]=vec(x[a,b],y)  
    return nextt
```

```
def graph2(x):  
    plt.figure()  
    plt.plot(x)  
    plt.xlabel('Epochs')  
    plt.ylabel('Cross-Entropy Error')
```

```
def m2(x,y):  
    nextt = np.zeros((x.shape[0],y))  
    for a in range(x.shape[0]):
```

```
nextt[a,:]=vec(x[a],y)
return nextt
```

```
def pb(x,y,z,dici,k):
    randoma = np.random.permutation(len(x))[0:5]
    tests = x[randoma]
    tests1 = m1(tests,dici)
    testsl = k[randoma]
    testsl = testsl.reshape(len(testsl),1)
    t10 = z.top10(tests1, 10)
    for i in range(5):
        print('[' + str(i+1) + ']' + str(y[tests[i,0]-1].decode('utf-8')) + ' ' + str(y[tests[i,1]-1].decode('utf-8')) + ' ' + str(y[tests[i,2]-1].decode('utf-8'))))
        print('True label= ' + str(y[testsl[i,0]-1].decode('utf-8'))))
        stri = 'The Top-10 Candidates are: {'
        for j in range(10):
            stri += (str(y[t10[j,i]].decode('utf-8')) + ' '
        print(stri+'}')

```

```
def q1():
```

```
class Network:
    def __init__(self,dim,neuron,std,mean=0):
        self.dim = dim
        self.neuron = neuron
        self.bias = np.random.normal(mean,std,neuron).reshape(neuron,1)
        self.weight = np.random.normal(mean,std,dim*neuron).reshape(neuron,dim)
        self.param = np.concatenate((self.weight,self.bias), axis=1)
```

```
self.chan = None

self.err = None

self.prevAct = None

self.prevchan = 0

def activation(self, x):
    if(x.ndim == 1):
        x = x.reshape(x.shape[0],1)
    self.prevAct = np.tanh(np.matmul(self.param, np.r_[x, [np.ones(x.shape[1])*-1]]))
    return self.prevAct

def derAct(self,a):
    return 1-(a**2)

def __repr__(self):
    return 'Input_Dim: '+str(self.dim)+' , Neuron #: '+str(self.neuron)

class PropagateLayers:
    def __init__(self):
        self.lays = []

    def appendLayers(self,lay):
        self.lays.append(lay)

    def forward(self,imp):
        nextt = imp
        for layers in self.lays:
            nextt = layers.activation(nextt)
        return nextt
```

```
def accur(self,x):
    print('Accuracy:', str(np.sum(x.predict(testimsf.T/255).T == testlbls)/len(testlbls)*100)
+ "%")

def predict(self,imp):
    nextt = self.forward(imp)
    nextt[nextt>=0] = 1
    nextt[nextt<0] = -1
    return nextt

def romando(self,imp,out):
    randindex = np.random.permutation(len(imp))
    imp,out = imp[randindex],out[randindex]
    return imp,out

def BP(self,imp,out,lrate,batch,momentum=0):
    network_output = self.forward(imp)
    for k in reversed(range(len(self.lays))):
        layers = self.lays[k]
        if(layers == self.lays[-1]):
            layers.err = out - network_output
            layers.chan = (layers.derAct(layers.prevAct)) * layers.err

        else:
            nextlayers = self.lays[k+1]
            nextlayers.weight = nextlayers.param[:,0:nextlayers.weight.shape[1]]
            layers.err = np.matmul(nextlayers.weight.T, nextlayers.chan)
            layers.chan = (layers.derAct(layers.prevAct)) * layers.err
```

```
for k in range(len(self.lays)):
    layers = self.lays[k]
    if (k==0):
        checkndim(imp)
        useimp = np.r_[imp,[np.ones(imp.shape[1])*-1]]
    else:
        useimp = np.r_[self.lays[k-1].prevAct,[np.ones(self.lays[k-1].prevAct.shape[1])*-
1]]

    if momentum == 0:
        layers.param = layers.param + (lr*rate*np.matmul(layers.chan, useimp.T))/batch
    else:
        ch = (lr*rate*np.matmul(layers.chan, useimp.T))
        layers.param = layers.param + momentum*layers.prevchan
        layers.prevchan = ch

def mserr(self, imp, out):
    mse = np.mean((out.T - self.forward(imp.T))**2, axis=1)
    return mse

def mcerr(self, imp, out):
    mce = np.sum(self.predict(imp.T) == out.T)/len(out)*100
    return mce

def workout(self, imp, out, impT, outT, lr, epoch, batch, momentum = 0):
    mse_err = []
    mce_err = []
    mset_err = []
```

```
mcet_err = []
for x in range(epoch):
    print('Epoch',x)
    imp,out = self.romando(imp, out)
    bat = int(np.floor(len(imp)/batch))

    for s in range(bat):
        self.BP(imp[batch*s:batch*(s+1)].T, out[batch*s:batch*(s+1)].T, lrate, batch)

    mse_err.append(self.mserr(imp, out))
    mce_err.append(self.mcerr(imp, out))
    mset_err.append(self.mserr(impT, outT))
    mcet_err.append(self.mcerr(impT, outT))

return mse_err, mce_err, mset_err, mcet_err

def __repr__(self):
    strr = ''
    for i, layers in enumerate(self.lays):
        strr += 'Layer ' + str(i) + ': ' + layers.__repr__() + '\n'
    return strr

filename = 'assign2_data1.h5'
f1 = h5py.File(filename,'r+')

testims = np.array(f1['testims'])
testlbls = np.array(f1['testlbls'])
trainims = np.array(f1['trainims'])
```

```
trainlbls = np.array(f1['trainlbls'])

testlbls = testlbls.reshape((len(testlbls),1))
trainlbls = trainlbls.reshape((len(trainlbls),1))

trainimsf = trainims.reshape(trainims.shape[0],(trainims.shape[1])**2)
testimsf = testims.reshape(testims.shape[0], (testims.shape[1])**2)

""" part a """
print('\nQuestion-1 Part-a')
goo = PropagateLayers()
goo.appendLayers(Network((trainims.shape[1])**2, 30, 0.01))
goo.appendLayers(Network(30, 1, 0.01))

print(goo)

trainlbls[trainlbls==0] = -1
testlbls = testlbls.astype(int)
testlbls[testlbls == 0] = -1

get_mse, get_mce, get_mset, get_mcet = goo.workout(trainimsf/255, trainlbls,
testimsf/255, testlbls, 0.4, 300, 50)

goo accur(goo)

""" part c """

print('\n\nQuestion-1 Part-c')
print('High Hidden Neurons Training')
```



```
goohigh = PropagateLayers()
goohigh.appendLayers(Network((trainims.shape[1])**2, 150, 0.01))
goohigh.appendLayers(Network(150, 1, 0.01))
print(goohigh)
```

```
get_mseh, get_mceh, get_mseth, get_mceth = goohigh.workout(trainimsf/255, trainlbls,
testimsf/255, testlbls, 0.4, 300, 50)
```

```
print('Low Hidden Neurons Training')
goolow = PropagateLayers()
goolow.appendLayers(Network((trainims.shape[1])**2, 5, 0.01))
goolow.appendLayers(Network(5, 1, 0.01))
print(goolow)
```

```
get_msel, get_mcel, get_msetl, get_mctl = goolow.workout(trainimsf/255, trainlbls,
testimsf/255, testlbls, 0.4, 300, 50)
```

```
""" part d """
```

```
print('\nQuestion-1 Part-d')
goohid = PropagateLayers()
goohid.appendLayers(Network((trainims.shape[1])**2, 300, 0.01))
goohid.appendLayers(Network(300, 30, 0.01))
goohid.appendLayers(Network(30, 1, 0.01))
print(goohid)
```

```
get_msehid, get_mcehid, get_msethid, get_mcethid = goohid.workout(trainimsf/255,
trainlbls, testimsf/255, testlbls, 0.4, 300, 50)
```

```
goohid accur(goohid)
```

```
""" part e """
```

```
print('\n\nQuestion-1 Part-e')
```

```
goomom = PropagateLayers()
```

```
goomom.appendLayers(Network((trainims.shape[1])**2, 300, 0.01))
```

```
goomom.appendLayers(Network(300, 30, 0.01))
```

```
goomom.appendLayers(Network(30, 1, 0.01))
```

```
print(goomom)
```

```
get_msemom, get_mcemom, get_msetmom, get_mcetmom =  
goomom.workout(trainimsf/255, trainlbls, testimsf/255, testlbls, 0.4, 300, 50, 0.5)
```

```
goomom accur(goomom)
```

```
""" plots """
```

```
graph(get_mse)
```

```
plt.title('Part-a Training Mean Squared Error')
```

```
plt.ylabel('Mean Squared Error')
```

```
plt.show()
```

```
graph(get_mset)
```

```
plt.title('Part-a Testing Mean Squared Error')
```

```
plt.ylabel('Mean Squared Error')
```

```
plt.show()
```

```
graph(get_mce)
```

```
plt.title('Part-a Training Mean Classification Error')
```

```
plt.ylabel('Mean Classification Error')  
plt.show()
```

```
graph(get_mcet)  
plt.title('Part-a Testing Mean Classification Error')  
plt.ylabel('Mean Classification Error')  
plt.show()
```

```
graphm(get_mse,get_mseh,get_msel)  
plt.title('Part-c Training Mean Squared Error')  
plt.ylabel('Mean Squared Error')  
plt.show()
```

```
graphm(get_mset,get_mseth,get_msetl)  
plt.title('Part-c Testing Mean Squared Error')  
plt.ylabel('Mean Squared Error')  
plt.show()
```

```
graphm(get_mce,get_mceh,get_mcel)  
plt.title('Part-c Training Mean Classification Error')  
plt.ylabel('Mean Classification Error')  
plt.show()
```

```
graphm(get_mcet,get_mceth,get_mcetl)  
plt.title('Part-c Testing Mean Classification Error')  
plt.ylabel('Mean Classification Error')  
plt.show()
```

```
graph(get_msehid)
```

```
plt.title('Part-d Training Mean Squared Error')  
plt.ylabel('Mean Squared Error')  
plt.show()
```

```
graph(get_msethid)  
plt.title('Part-d Testing Mean Squared Error')  
plt.ylabel('Mean Squared Error')  
plt.show()
```

```
graph(get_mcehid)  
plt.title('Part-d Training Mean Classification Error')  
plt.ylabel('Mean Classification Error')  
plt.show()
```

```
graph(get_mcethid)  
plt.title('Part-d Testing Mean Classification Error')  
plt.ylabel('Mean Classification Error')  
plt.show()
```

```
graphd(get_msehid,get_msemom)  
plt.title('Part-e Training Mean Squared Error')  
plt.ylabel('Mean Squared Error')  
plt.show()
```

```
graphd(get_msethid,get_msetmom)  
plt.title('Part-e Testing Mean Squared Error')  
plt.ylabel('Mean Squared Error')  
plt.show()
```

```
graphd(get_mcehid,get_mcemom)
plt.title('Part-e Training Mean Classification Error')
plt.ylabel('Mean Classification Error')
plt.show()
```

```
graphd(get_mcethid,get_mcetmom)
plt.title('Part-e Testing Mean Classification Error')
plt.ylabel('Mean Classification Error')
plt.show()
```

```
def q2():
```

```
class Network2:
    def __init__(self, dim, neuron, activ, std, mean=0):
        self.dim = dim
        self.neuron = neuron
        self.activ = activ
        if self.activ == 'sigmoid' or self.activ == 'softmax':
            self.bias = np.random.normal(mean,std,neuron).reshape(neuron,1)
            self.weight = np.random.normal(mean,std,dim*neuron).reshape(neuron,dim)
            self.param = np.concatenate((self.weight,self.bias), axis=1)
        elif self.activ == 'wordembed':
            self.d = neuron
            self.idim = dim
            self.weight = np.random.normal(mean,std,self.d*self.idim).reshape(self.d,self.idim)
        self.chan = None
        self.err = None
```

```
self.prevAct = None

self.prevchan = 0

def activation(self,x):
    if (self.activ == 'sigmoid'):
        return np.exp(2*x)/(1+np.exp(2*x))
    elif (self.activ == 'softmax'):
        return np.exp(x-np.max(x))/np.sum(np.exp(x-np.max(x)),axis=0)
    elif (self.activ == 'wordembed'):
        return x

def activations(self,x):
    if self.activ == 'sigmoid' or self.activ == 'softmax':
        if (x.ndim == 1):
            x=x.reshape(x.shape[0],1)
            self.prevAct = self.activation(np.matmul(self.param,np.r_[x, [np.ones(x.shape[1])*-
1]]))
        elif self.activ == 'wordembed':
            nextlayer = np.zeros((x.shape[0],x.shape[1], self.idim))
            for k in range(nextlayer.shape[0]):
                nextlayer[k,:,:] = self.activation(np.matmul(x[k,:,:], self.weight))
            nextlayer =
nextlayer.reshape((nextlayer.shape[0],nextlayer.shape[1]*nextlayer.shape[2]))
            self.prevAct = nextlayer.T
        return self.prevAct

def derAct(self,a):
    if self.activ == 'sigmoid':
        return 2*(a*(1-a))
```

```
elif self.activ == 'softmax':  
    return a*(1-a)  
elif self.activ == 'wordembed':  
    return np.ones(a.shape)  
  
def __repr__(self):  
    return 'Input_Dim: '+str(self.dim)+' , Neuron #: '+str(self.neuron)+' \n Activation: '+  
self.activ  
  
class PropagateWords:  
    def __init__(self):  
        self.lays = []  
  
    def appendLayers(self,lay):  
        self.lays.append(lay)  
  
    def forward(self,imp):  
        nextt = imp  
        for layers in self.lays:  
            nextt = layers.activations(nextt)  
        return nextt  
  
    def predict(self,imp):  
        nextt = self.forward(imp)  
        if nextt.ndim == 1:  
            return np.argmax(nextt)  
        return np.argmax(nextt, axis=0)
```

```
def romando(self,imp,out):  
    randindex = np.random.permutation(len(imp))  
    imp,out = imp[randindex],out[randindex]  
    return imp,out  
  
def top10(self,imp,n):  
    nextt = self.forward(imp)  
    return np.argpartition(nextt, -n, axis=0)[-n:]  
  
def BP(self,imp,out,lrate,batch,momentum):  
    network_output = self.forward(imp)  
    for k in reversed(range(len(self.lays))):  
        layers = self.lays[k]  
        if layers == self.lays[-1]:  
            layers.chan = out - network_output  
        else:  
            nextlayers = self.lays[k+1]  
            nextlayers.weight = nextlayers.param[:,0:nextlayers.weight.shape[1]]  
            layers.err = np.matmul(nextlayers.weight.T, nextlayers.chan)  
            layers.chan = (layers.derAct(layers.prevAct)) * layers.err  
  
    for k in range(len(self.lays)):  
        layers = self.lays[k]  
        if k==0:  
            checkndim(imp)  
            useimp = imp  
        else:  
            useimp = np.r_[self.lays[k-1].prevAct, [np.ones(self.lays[k-1].prevAct.shape[1])*-  
1]]
```



```
if layers.activ == 'sigmoid' or layers.activ == 'softmax':  
    ch = (lr*rate*np.matmul(layers.chan, useimp.T))/batch  
    layers.param = layers.param + (ch + momentum*layers.prevchan)  
elif layers.activ == 'wordembed':  
    chan3 = layers.chan.reshape((3,batch,layers.idim))  
    useimp = np.transpose(useimp, (1,2,0))  
    ch = np.zeros((useimp.shape[1], chan3.shape[2]))  
    for a in range(chan3.shape[0]):  
        ch = ch + lr*rate*np.matmul(useimp[a,:,:], chan3[a,:,:])  
    ch = ch/batch  
    layers.weight = layers.weight + (ch + momentum*layers.prevchan)  
layers.prevchan = ch
```

```
def workout(self,imp,out,impT,outT,lr*rate,epoch,batch,momentum):  
    entlist = []  
    for x in range(epoch):  
        print('\nEpoch',x)  
  
        imp,out = self.romando(imp, out)  
        bat = int(np.floor(len(imp)/batch))  
  
        for j in range(bat):  
            batimp = m1(imp[batch*j:batch*(j+1)],d)  
            batout = m2(out[batch*j:batch*(j+1)],d).T  
            self.BP(batimp,batout,lr*rate,batch,momentum)  
  
        ov = self.forward(impT)  
        CEerr = -np.sum(np.log(ov)*outT.T)/ov.shape[1]
```

```
print('Cross-Entropy Error: ',CEerr)

entlist.append(CEerr)

av = np.sum(self.predict(impT) == np.argmax(outT.T, axis=0))

print('Correct: ',av)

print('Accuracy: ', (av/ov.shape[1])*100,'%')

return entlist

def __repr__(self):
    strr = ""
    for i, layers in enumerate(self.lays):
        strr += 'Layer ' + str(i) + ': ' + layers.__repr__() + '\n'
    return strr

filename2 = 'assign2_data2.h5'
f2 = h5py.File(filename2,'r+')

words = np.array(f2['words'])
trainin = np.array(f2['trainx'])
trainout = np.array(f2['traind'])
valin = np.array(f2['valx'])
valout = np.array(f2['vald'])
testin = np.array(f2['testx'])
testout = np.array(f2['testd'])
```

```
""" part a """
```

```
print('Question-2 Part-a')
```

```
p = 256
```

```
p2 = 128
```

```
p3 = 64
```

```
idim = 32
```

```
idim2 = 16
```

```
idim3 = 8
```

```
d = 250
```

```
lrate2 = 0.4
```

```
momentum2 = 0.85
```

```
batch2 = 250
```

```
epoch2 = 50
```

```
vin = m1(valin,d)
```

```
vout = m2(valout,d)
```

```
print('\nD = 32, P = 256')
```

```
goo2 = PropagateWords()
```

```
goo2.appendLayers(Network2(idim,d,'wordembed',0.1))
```

```
goo2.appendLayers(Network2(3*idim,p,'sigmoid',0.1))
```

```
goo2.appendLayers(Network2(p,d,'softmax',0.1))
```

```
err1 = goo2.workout(trainin,trainout,vin,vout,lrate2,epoch2,batch2,momentum2)
```

```
print('\nD = 16, P = 128')
```

```
goo22 = PropagateWords()
```

```
goo22.appendLayers(Network2(idim2,d,'wordembed',0.1))
goo22.appendLayers(Network2(3*idim2,p2,'sigmoid',0.1))
goo22.appendLayers(Network2(p2,d,'softmax',0.1))

err2 = goo22.workout(trainin,trainout,vin,vout,lrate2,epoch2,batch2,momentum2)

print('\nD = 8, P = 64')
goo23 = PropagateWords()
goo23.appendLayers(Network2(idim3,d,'wordembed',0.1))
goo23.appendLayers(Network2(3*idim3,p3,'sigmoid',0.1))
goo23.appendLayers(Network2(p3,d,'softmax',0.1))

err3 = goo23.workout(trainin,trainout,vin,vout,lrate2,epoch2,batch2,momentum2)

print('To execute the part b, after observation please close the plots')

graph2(err1)
plt.title('Cross-Entropy Error, D=32,P=256')
plt.show()

graph2(err2)
plt.title('Cross-Entropy Error, D=16,P=128')
plt.show()

graph2(err3)
plt.title('Cross-Entropy Error, D=8,P=64')
plt.show()
```

```
""" part b """
```

```
pb(testin,words,goo2,d,testout)
```

```
question = sys.argv[1]
```

```
def ayberk_yarkin_yildiz_21803386_hw2(question):
```

```
    if question == '1' :
```

```
        q1()
```

```
    elif question == '2' :
```

```
        q2()
```

```
ayberk_yarkin_yildiz_21803386_hw2(question)
```

Question 3

Fully-Connected Neural Nets

In the previous homework you implemented a fully-connected two-layer neural network on CIFAR-10. The implementation was simple but not very modular since the loss and gradient were computed in a single monolithic function. This is manageable for a simple two-layer network, but would become impractical as we move to bigger models. Ideally we want to build networks using a more modular design so that we can implement different layer types in isolation and then snap them together into models with different architectures.

In this exercise we will implement fully-connected networks using a more modular approach. For each layer we will implement a `forward` and a `backward` function. The `forward` function will receive inputs, weights, and other parameters and will return both an output and a `cache` object storing data needed for the backward pass, like this:

```
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 dout (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
```

After implementing a bunch of layers this way, we will be able to easily combine them to build classifiers with different architectures.

In addition to implementing fully-connected networks of arbitrary depth, we will also explore different update rules for optimization, and introduce Dropout as a regularizer and Batch/Layer Normalization as a tool to more efficiently optimize deep networks.

```
In [1]: # As usual, a bit of setup
from __future__ import print_function
import time
import numpy as np
```

```

import matplotlib.pyplot as plt
from cs231n.classifiers.fc_net import *
from cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_gradient_a
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))))

```

run the following from the cs231n directory and try again:
python setup.py build_ext --inplace
You may also need to restart your iPython kernel

In [2]: *# Load the (preprocessed) CIFAR10 data.*

```

data = get_CIFAR10_data()
for k, v in list(data.items()):
    print('%s: ' % k, v.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,))

```

Affine layer: forward

Open the file cs231n/layers.py and implement the affine_forward function.

Once you are done you can test your implementaion by running the following:

In [3]: *# 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 e-9 or less.

```

```
print('Testing affine_forward function:')
print('difference: ', rel_error(out, correct_out))
```

Testing affine_forward function:
difference: 9.769847728806635e-10

Affine layer: backward

Now implement the `affine_backward` function and test your implementation using numeric gradient checking.

In [4]:

```
# Test the affine_backward function
np.random.seed(231)
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 e-10 or less
print('Testing affine_backward function:')
print('dx error: ', rel_error(dx_num, dx))
print('dw error: ', rel_error(dw_num, dw))
print('db error: ', rel_error(db_num, db))
```

Testing affine_backward function:
dx error: 5.399100368651805e-11
dw error: 9.904211865398145e-11
db error: 2.4122867568119087e-11

ReLU activation: forward

Implement the forward pass for the ReLU activation function in the `relu_forward` function and test your implementation using the following:

In [5]:

```
# 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 on the order of e-8
print('Testing relu_forward function:')
print('difference: ', rel_error(out, correct_out))
```

Testing relu_forward function:
difference: 4.999999798022158e-08

ReLU activation: backward

Now implement the backward pass for the ReLU activation function in the `relu_backward` function and test your implementation using numeric gradient checking:

```
In [6]: np.random.seed(231)
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 on the order of e-12
print('Testing relu_backward function:')
print('dx error: ', rel_error(dx_num, dx))
```

```
Testing relu_backward function:
dx error: 3.2756349136310288e-12
```

Inline Question 1:

We've only asked you to implement ReLU, but there are a number of different activation functions that one could use in neural networks, each with its pros and cons. In particular, an issue commonly seen with activation functions is getting zero (or close to zero) gradient flow during backpropagation. Which of the following activation functions have this problem? If you consider these functions in the one dimensional case, what types of input would lead to this behaviour?

1. Sigmoid
2. ReLU
3. Leaky ReLU

Answer:

1- For the sigmoid activation function, negative values of the function converge to 0. Therefore, for sigmoid, the negative values cause the activation function to saturate.

2- For ReLU, output of the activation function is 0 for all negative values. Therefore, any negative value will trigger zero gradient flow. Therefore, if the bias is a large negative term, the gradients and the output of the activation function ReLU becomes 0.

To compare, for sigmoid function, the zero gradient flow is recoverable since it gives a small value. But for ReLU, the gradient will always be zero and the zero gradient flow is not recoverable.

3- The Leaky ReLU occurs as a solution to the problem of ReLU, as it does not converge to 0 and it is very small but non-zero for negative values, whereas ReLU is always zero in negative values.

"Sandwich" layers

There are some common patterns of layers that are frequently used in neural nets. For example, affine layers are frequently followed by a ReLU nonlinearity. To make these common patterns easy, we define several convenience layers in the file `cs231n/layer_utils.py`.

For now take a look at the `affine_relu_forward` and `affine_relu_backward` functions, and run the following to numerically gradient check the backward pass:

In [7]:

```
from cs231n.layer_utils import affine_relu_forward, affine_relu_backward
np.random.seed(231)
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,
dw_num = eval_numerical_gradient_array(lambda w: affine_relu_forward(x, w, b)[0], w,
db_num = eval_numerical_gradient_array(lambda b: affine_relu_forward(x, w, b)[0], b,

# Relative error should be around e-10 or less
print('Testing affine_relu_forward and affine_relu_backward:')
print('dx error: ', rel_error(dx_num, dx))
print('dw error: ', rel_error(dw_num, dw))
print('db error: ', rel_error(db_num, db))
```

```
Testing affine_relu_forward and affine_relu_backward:
dx error: 6.750562121603446e-11
dw error: 8.162015570444288e-11
db error: 7.826724021458994e-12
```

Loss layers: Softmax and SVM

You implemented these loss functions in the last assignment, so we'll give them to you for free here. You should still make sure you understand how they work by looking at the implementations in `cs231n/layers.py`.

You can make sure that the implementations are correct by running the following:

In [8]:

```
np.random.seed(231)
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 around the
print('Testing svm_loss:')
print('loss: ', loss)
print('dx error: ', rel_error(dx_num, dx))

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

```
# Test softmax_loss function. Loss should be close to 2.3 and dx error should be around 1e-9
print('\nTesting softmax_loss:')
print('loss: ', loss)
print('dx error: ', rel_error(dx_num, dx))
```

```
Testing svm_loss:
loss: 8.999602749096233
dx error: 1.4021566006651672e-09
```

```
Testing softmax_loss:
loss: 2.302545844500738
dx error: 9.384673161989355e-09
```

Two-layer network

In the previous assignment you implemented a two-layer neural network in a single monolithic class. Now that you have implemented modular versions of the necessary layers, you will reimplement the two layer network using these modular implementations.

Open the file `cs231n/classifiers/fc_net.py` and complete the implementation of the `TwoLayerNet` class. This class will serve as a model for the other networks you will implement in this assignment, so read through it to make sure you understand the API. You can run the cell below to test your implementation.

In [9]:

```
np.random.seed(231)
N, D, H, C = 3, 5, 50, 7
X = np.random.randn(N, D)
y = np.random.randint(C, size=N)

std = 1e-3
model = TwoLayerNet(input_dim=D, hidden_dim=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,
      12.05769098, 12.74614105, 13.43459113, 14.1230412, 14.81149128, 15.49994135,
      12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.05099822, 15.66781506]
    )
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'

# Errors should be around e-7 or Less
for reg in [0.0, 0.7]:
    print('Running numeric gradient check with reg = ', 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('%s relative error: %.2e' % (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: 1.52e-08
W2 relative error: 3.48e-10
b1 relative error: 6.55e-09
b2 relative error: 4.33e-10
Running numeric gradient check with reg = 0.7
W1 relative error: 8.18e-07
W2 relative error: 2.85e-08
b1 relative error: 1.09e-09
b2 relative error: 7.76e-10

```

Solver

In the previous assignment, the logic for training models was coupled to the models themselves. Following a more modular design, for this assignment we have split the logic for training models into a separate class.

Open the file `cs231n/solver.py` and read through it to familiarize yourself with the API. After doing so, use a `Solver` instance to train a `TwoLayerNet` that achieves at least 50% accuracy on the validation set.

In [10]:

```

model = TwoLayerNet()
solver = None

#####
# TODO: Use a Solver instance to train a TwoLayerNet that achieves at least #
# 50% accuracy on the validation set.                                     #
#####

model = TwoLayerNet(hidden_dim=100, reg=0.2)
solver = Solver(model, data, update_rule='sgd',
                optim_config={'learning_rate': 1e-3}, lr_decay=0.95,
                num_epochs=10, batch_size=100, print_every=100)
solver.train()

#####
#                                     END OF YOUR CODE                                     #
#####

```

```

(Iteration 1 / 4900) loss: 2.332096
(Epoch 0 / 10) train acc: 0.164000; val_acc: 0.134000
(Iteration 101 / 4900) loss: 1.857220
(Iteration 201 / 4900) loss: 2.000576
(Iteration 301 / 4900) loss: 1.651815
(Iteration 401 / 4900) loss: 1.538214
(Epoch 1 / 10) train acc: 0.450000; val_acc: 0.454000
(Iteration 501 / 4900) loss: 1.608869
(Iteration 601 / 4900) loss: 1.501398
(Iteration 701 / 4900) loss: 1.615213
(Iteration 801 / 4900) loss: 1.656747
(Iteration 901 / 4900) loss: 1.468052
(Epoch 2 / 10) train acc: 0.484000; val_acc: 0.472000
(Iteration 1001 / 4900) loss: 1.505273
(Iteration 1101 / 4900) loss: 1.503323
(Iteration 1201 / 4900) loss: 1.418404
(Iteration 1301 / 4900) loss: 1.356568
(Iteration 1401 / 4900) loss: 1.507079
(Epoch 3 / 10) train acc: 0.519000; val_acc: 0.475000
(Iteration 1501 / 4900) loss: 1.405298
(Iteration 1601 / 4900) loss: 1.425098
(Iteration 1701 / 4900) loss: 1.388389
(Iteration 1801 / 4900) loss: 1.559448
(Iteration 1901 / 4900) loss: 1.469148
(Epoch 4 / 10) train acc: 0.506000; val_acc: 0.488000
(Iteration 2001 / 4900) loss: 1.521458
(Iteration 2101 / 4900) loss: 1.452836
(Iteration 2201 / 4900) loss: 1.515952
(Iteration 2301 / 4900) loss: 1.253438
(Iteration 2401 / 4900) loss: 1.329813
(Epoch 5 / 10) train acc: 0.546000; val_acc: 0.490000
(Iteration 2501 / 4900) loss: 1.385455
(Iteration 2601 / 4900) loss: 1.380330
(Iteration 2701 / 4900) loss: 1.344157
(Iteration 2801 / 4900) loss: 1.516297
(Iteration 2901 / 4900) loss: 1.373451
(Epoch 6 / 10) train acc: 0.548000; val_acc: 0.519000
(Iteration 3001 / 4900) loss: 1.313017
(Iteration 3101 / 4900) loss: 1.139112
(Iteration 3201 / 4900) loss: 1.596601
(Iteration 3301 / 4900) loss: 1.372248
(Iteration 3401 / 4900) loss: 1.524008
(Epoch 7 / 10) train acc: 0.537000; val_acc: 0.491000
(Iteration 3501 / 4900) loss: 1.325397
(Iteration 3601 / 4900) loss: 1.141724
(Iteration 3701 / 4900) loss: 1.368370
(Iteration 3801 / 4900) loss: 1.319290
(Iteration 3901 / 4900) loss: 1.101957
(Epoch 8 / 10) train acc: 0.545000; val_acc: 0.499000
(Iteration 4001 / 4900) loss: 1.239187
(Iteration 4101 / 4900) loss: 1.346376
(Iteration 4201 / 4900) loss: 1.154919
(Iteration 4301 / 4900) loss: 1.073516
(Iteration 4401 / 4900) loss: 1.577285
(Epoch 9 / 10) train acc: 0.595000; val_acc: 0.503000
(Iteration 4501 / 4900) loss: 1.253220
(Iteration 4601 / 4900) loss: 1.465048
(Iteration 4701 / 4900) loss: 1.484373
(Iteration 4801 / 4900) loss: 1.242994
(Epoch 10 / 10) train acc: 0.564000; val_acc: 0.471000

```

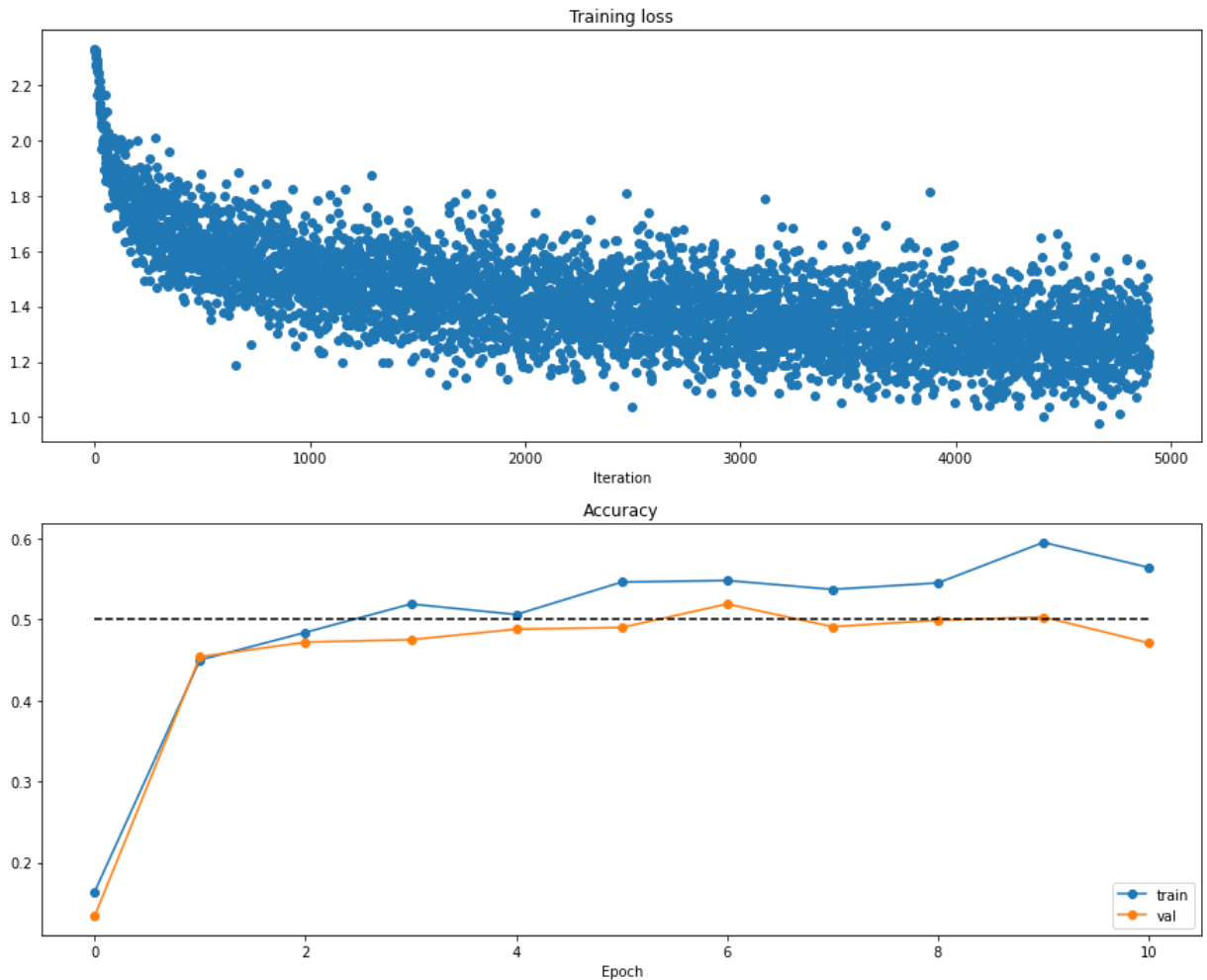
```

In [11]: # 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 network

Next you will implement a fully-connected network with an arbitrary number of hidden layers.

Read through the `FullyConnectedNet` class in the file `cs231n/classifiers/fc_net.py`.

Implement the initialization, the forward pass, and the backward pass. For the moment don't worry about implementing dropout or batch/layer normalization; we will add those features soon.

Initial loss and gradient check

As a sanity check, run the following to check the initial loss and to gradient check the network both with and without regularization. Do the initial losses seem reasonable?

For gradient checking, you should expect to see errors around $1e-7$ or less.

```
In [12]: np.random.seed(231)
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 = ', 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: ', loss)

    # Most of the errors should be on the order of e-7 or smaller.
    # NOTE: It is fine however to see an error for W2 on the order of e-5
    # for the check when reg = 0.0
    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('%s relative error: %.2e' % (name, rel_error(grad_num, grads[name])))
```

```
Running check with reg = 0
Initial loss: 2.3004790897684924
W1 relative error: 1.48e-07
W2 relative error: 2.21e-05
W3 relative error: 3.53e-07
b1 relative error: 5.38e-09
b2 relative error: 2.09e-09
b3 relative error: 5.80e-11
Running check with reg = 3.14
Initial loss: 7.052114776533016
W1 relative error: 7.36e-09
W2 relative error: 6.87e-08
W3 relative error: 3.48e-08
b1 relative error: 1.48e-08
b2 relative error: 1.72e-09
b3 relative error: 1.80e-10
```

As another sanity check, make sure you can overfit a small dataset of 50 images. First we will try a three-layer network with 100 units in each hidden layer. In the following cell, tweak the learning rate and initialization scale to overfit and achieve 100% training accuracy within 20 epochs.

```
In [13]: # TODO: Use a three-layer Net to overfit 50 training examples by
#         tweaking just the learning rate and initialization scale.

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'],
}

weight_scale = 1e-1
learning_rate = 1e-3
model = FullyConnectedNet([100, 100],
                          weight_scale=weight_scale, dtype=np.float64)
solver = Solver(model, small_data,
                print_every=10, num_epochs=20, batch_size=25,
                update_rule='sgd',
                optim_config={
                    'learning_rate': learning_rate,
                })
```

```

    )
    solver.train()

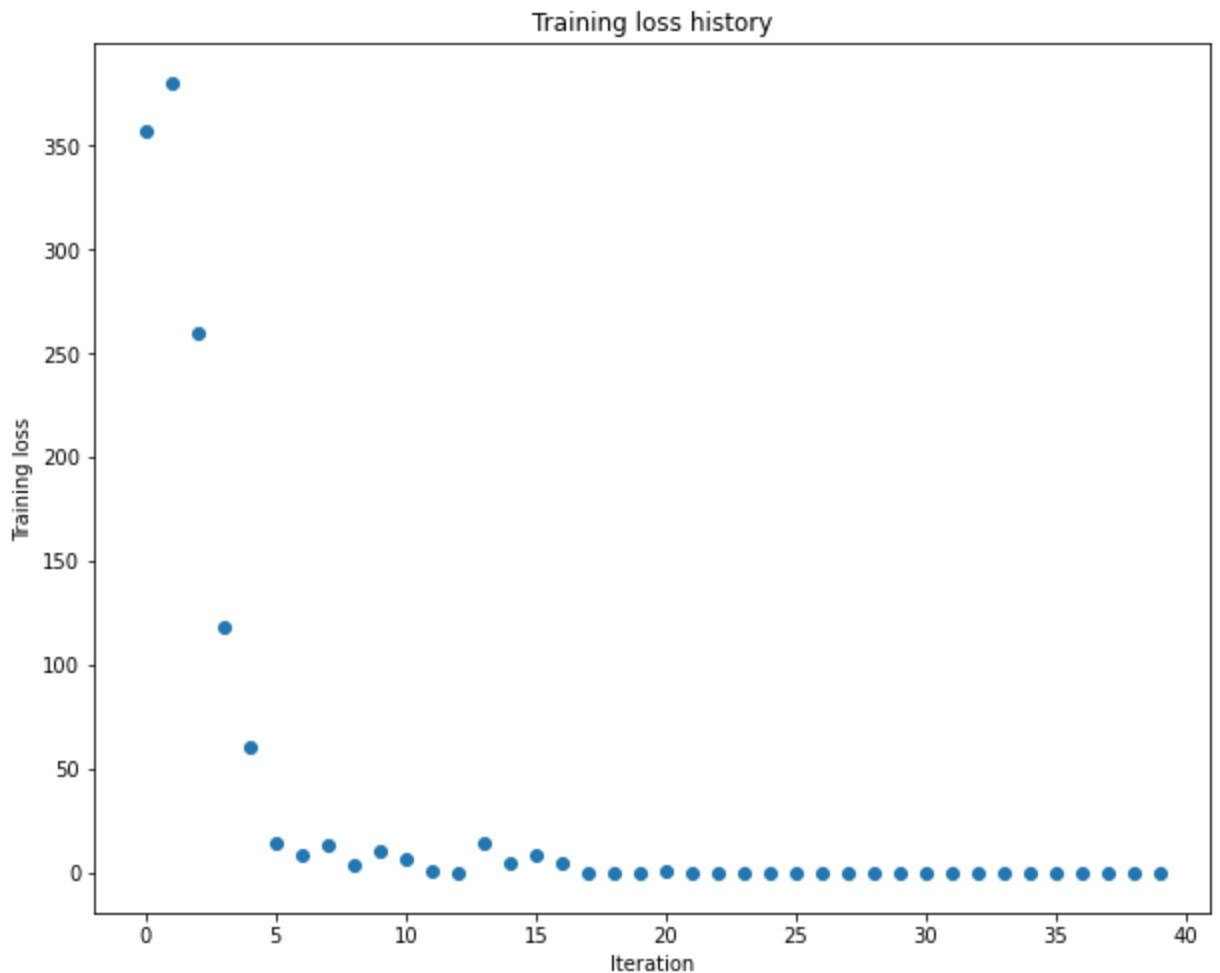
    plt.plot(solver.loss_history, 'o')
    plt.title('Training loss history')
    plt.xlabel('Iteration')
    plt.ylabel('Training loss')
    plt.show()

```

```

(Iteration 1 / 40) loss: 357.428290
(Epoch 0 / 20) train acc: 0.220000; val_acc: 0.111000
(Epoch 1 / 20) train acc: 0.380000; val_acc: 0.141000
(Epoch 2 / 20) train acc: 0.520000; val_acc: 0.138000
(Epoch 3 / 20) train acc: 0.740000; val_acc: 0.130000
(Epoch 4 / 20) train acc: 0.820000; val_acc: 0.153000
(Epoch 5 / 20) train acc: 0.860000; val_acc: 0.175000
(Iteration 11 / 40) loss: 6.726589
(Epoch 6 / 20) train acc: 0.940000; val_acc: 0.163000
(Epoch 7 / 20) train acc: 0.960000; val_acc: 0.166000
(Epoch 8 / 20) train acc: 0.960000; val_acc: 0.164000
(Epoch 9 / 20) train acc: 0.980000; val_acc: 0.162000
(Epoch 10 / 20) train acc: 0.980000; val_acc: 0.162000
(Iteration 21 / 40) loss: 0.800243
(Epoch 11 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 12 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 13 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 14 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.158000
(Iteration 31 / 40) loss: 0.000000
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.158000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.158000

```



Now try to use a five-layer network with 100 units on each layer to overfit 50 training examples. Again you will have to adjust the learning rate and weight initialization, but you should be able to achieve 100% training accuracy within 20 epochs.

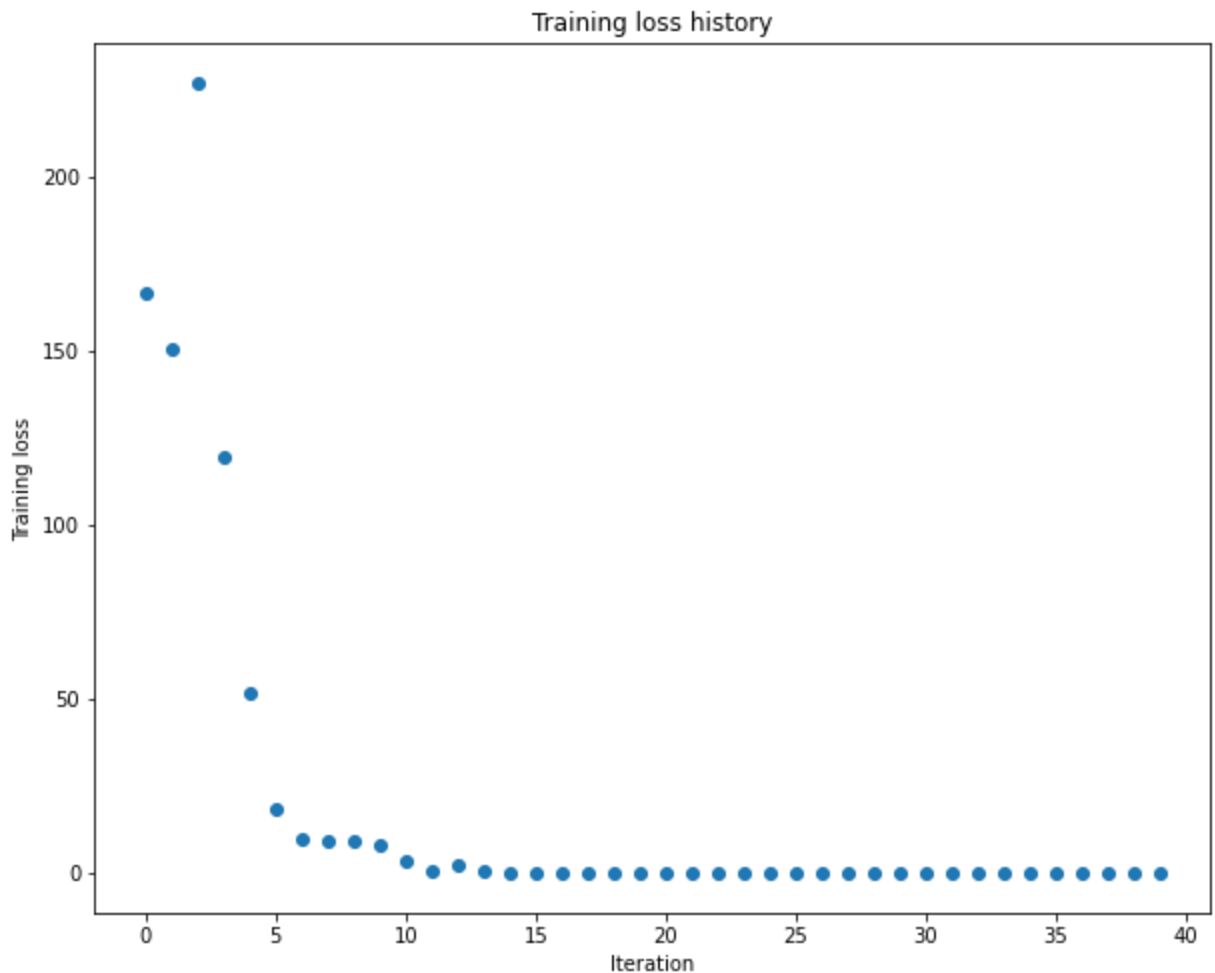
```
In [14]: # TODO: Use a five-layer Net to overfit 50 training examples by
# tweaking just the learning rate and initialization scale.

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'],
}

learning_rate = 2e-3
weight_scale = 1e-1
model = FullyConnectedNet([100, 100, 100, 100],
                           weight_scale=weight_scale, dtype=np.float64)
solver = Solver(model, small_data,
                 print_every=10, num_epochs=20, batch_size=25,
                 update_rule='sgd',
                 optim_config={
                     'learning_rate': learning_rate,
                 })
solver.train()

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

```
(Iteration 1 / 40) loss: 166.501707
(Epoch 0 / 20) train acc: 0.100000; val_acc: 0.107000
(Epoch 1 / 20) train acc: 0.320000; val_acc: 0.101000
(Epoch 2 / 20) train acc: 0.160000; val_acc: 0.122000
(Epoch 3 / 20) train acc: 0.380000; val_acc: 0.106000
(Epoch 4 / 20) train acc: 0.520000; val_acc: 0.111000
(Epoch 5 / 20) train acc: 0.760000; val_acc: 0.113000
(Iteration 11 / 40) loss: 3.343141
(Epoch 6 / 20) train acc: 0.840000; val_acc: 0.122000
(Epoch 7 / 20) train acc: 0.920000; val_acc: 0.113000
(Epoch 8 / 20) train acc: 0.940000; val_acc: 0.125000
(Epoch 9 / 20) train acc: 0.960000; val_acc: 0.125000
(Epoch 10 / 20) train acc: 0.980000; val_acc: 0.121000
(Iteration 21 / 40) loss: 0.039138
(Epoch 11 / 20) train acc: 0.980000; val_acc: 0.123000
(Epoch 12 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 13 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 14 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.121000
(Iteration 31 / 40) loss: 0.000644
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.121000
```



Inline Question 2:

Did you notice anything about the comparative difficulty of training the three-layer net vs training the five layer net? In particular, based on your experience, which network seemed more sensitive to the initialization scale? Why do you think that is the case?

Answer:

The larger number of layers in a network can cause difficulties in controlling the network. The network becomes way more sensitive to the parameters, such as the learning rate. If any initialization parameter is changed, the network behavior is affected extensively. Therefore, the larger layered networks is harder to train. When number of hidden layers increases, more complex the network becomes in the hidden part, and it becomes more unpredictable. In this hidden part, since the complexity increases, any certain change in weight initialization can change the behavior of the hidden part dramatically, if there is bigger number of layers.

Update rules

So far we have used vanilla stochastic gradient descent (SGD) as our update rule. More sophisticated update rules can make it easier to train deep networks. We will implement a few of the most commonly used update rules and compare them to vanilla SGD.

SGD+Momentum

Stochastic gradient descent with momentum is a widely used update rule that tends to make deep networks converge faster than vanilla stochastic gradient descent. See the Momentum Update section at <http://cs231n.github.io/neural-networks-3/#sgd> for more information.

Open the file `cs231n/optim.py` and read the documentation at the top of the file to make sure you understand the API. Implement the SGD+momentum update rule in the function `sgd_momentum` and run the following to check your implementation. You should see errors less than $e-8$.

```
In [15]: from cs231n.optim import sgd_momentum

N, D = 4, 5
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
v = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-3, 'velocity': v}
next_w, _ = sgd_momentum(w, dw, config=config)

expected_next_w = np.asarray([
    [ 0.1406,      0.20738947,  0.27417895,  0.34096842,  0.40775789],
    [ 0.47454737,  0.54133684,  0.60812632,  0.67491579,  0.74170526],
    [ 0.80849474,  0.87528421,  0.94207368,  1.00886316,  1.07565263],
    [ 1.14244211,  1.20923158,  1.27602105,  1.34281053,  1.4096      ]])
expected_velocity = np.asarray([
    [ 0.5406,      0.55475789,  0.56891579,  0.58307368,  0.59723158],
    [ 0.61138947,  0.62554737,  0.63970526,  0.65386316,  0.66802105],
    [ 0.68217895,  0.69633684,  0.71049474,  0.72465263,  0.73881053],
    [ 0.75296842,  0.76712632,  0.78128421,  0.79544211,  0.8096      ]])

# Should see relative errors around e-8 or less
print('next_w error: ', rel_error(next_w, expected_next_w))
print('velocity error: ', rel_error(expected_velocity, config['velocity']))
```

```
next_w error: 8.882347033505819e-09
velocity error: 4.269287743278663e-09
```

Once you have done so, run the following to train a six-layer network with both SGD and SGD+momentum. You should see the SGD+momentum update rule converge faster.

```
In [16]: num_train = 4000
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'],
}

solvers = {}

for update_rule in ['sgd', 'sgd_momentum']:
    print('running with ', update_rule)
    model = FullyConnectedNet([100, 100, 100, 100, 100], weight_scale=5e-2)

    solver = Solver(model, small_data,
                    num_epochs=5, batch_size=100,
                    update_rule=update_rule,
```

```

        optim_config={
            'learning_rate': 1e-2,
        },
        verbose=True)
    solvers[update_rule] = solver
    solver.train()
    print()

    plt.subplot(3, 1, 1)
    plt.title('Training loss')
    plt.xlabel('Iteration')

    plt.subplot(3, 1, 2)
    plt.title('Training accuracy')
    plt.xlabel('Epoch')

    plt.subplot(3, 1, 3)
    plt.title('Validation accuracy')
    plt.xlabel('Epoch')

    for update_rule, solver in list(solvers.items()):
        plt.subplot(3, 1, 1)
        plt.plot(solver.loss_history, 'o', label=update_rule)

        plt.subplot(3, 1, 2)
        plt.plot(solver.train_acc_history, '-o', label=update_rule)

        plt.subplot(3, 1, 3)
        plt.plot(solver.val_acc_history, '-o', label=update_rule)

    for i in [1, 2, 3]:
        plt.subplot(3, 1, i)
        plt.legend(loc='upper center', ncol=4)
    plt.gcf().set_size_inches(15, 15)
    plt.show()

```

```

running with  sgd
(Iteration 1 / 200) loss: 2.559978
(Epoch 0 / 5) train acc: 0.103000; val_acc: 0.108000
(Iteration 11 / 200) loss: 2.291086
(Iteration 21 / 200) loss: 2.153591
(Iteration 31 / 200) loss: 2.082693
(Epoch 1 / 5) train acc: 0.277000; val_acc: 0.242000
(Iteration 41 / 200) loss: 2.004171
(Iteration 51 / 200) loss: 2.010409
(Iteration 61 / 200) loss: 2.023753
(Iteration 71 / 200) loss: 2.026621
(Epoch 2 / 5) train acc: 0.352000; val_acc: 0.312000
(Iteration 81 / 200) loss: 1.807163
(Iteration 91 / 200) loss: 1.914256
(Iteration 101 / 200) loss: 1.917177
(Iteration 111 / 200) loss: 1.706193
(Epoch 3 / 5) train acc: 0.405000; val_acc: 0.322000
(Iteration 121 / 200) loss: 1.697994
(Iteration 131 / 200) loss: 1.768837
(Iteration 141 / 200) loss: 1.784967
(Iteration 151 / 200) loss: 1.823291
(Epoch 4 / 5) train acc: 0.431000; val_acc: 0.324000
(Iteration 161 / 200) loss: 1.626499
(Iteration 171 / 200) loss: 1.901366
(Iteration 181 / 200) loss: 1.550534
(Iteration 191 / 200) loss: 1.716921
(Epoch 5 / 5) train acc: 0.436000; val_acc: 0.330000

```

```

running with  sgd_momentum
(Iteration 1 / 200) loss: 3.153778

```

```
(Epoch 0 / 5) train acc: 0.105000; val_acc: 0.093000
(Iteration 11 / 200) loss: 2.145874
(Iteration 21 / 200) loss: 2.032563
(Iteration 31 / 200) loss: 1.985848
(Epoch 1 / 5) train acc: 0.311000; val_acc: 0.281000
(Iteration 41 / 200) loss: 1.882354
(Iteration 51 / 200) loss: 1.855372
(Iteration 61 / 200) loss: 1.649133
(Iteration 71 / 200) loss: 1.806432
(Epoch 2 / 5) train acc: 0.415000; val_acc: 0.324000
(Iteration 81 / 200) loss: 1.907840
(Iteration 91 / 200) loss: 1.510681
(Iteration 101 / 200) loss: 1.546872
(Iteration 111 / 200) loss: 1.512047
(Epoch 3 / 5) train acc: 0.434000; val_acc: 0.321000
(Iteration 121 / 200) loss: 1.677301
(Iteration 131 / 200) loss: 1.504686
(Iteration 141 / 200) loss: 1.633253
(Iteration 151 / 200) loss: 1.745081
(Epoch 4 / 5) train acc: 0.460000; val_acc: 0.353000
(Iteration 161 / 200) loss: 1.485411
(Iteration 171 / 200) loss: 1.610416
(Iteration 181 / 200) loss: 1.528331
(Iteration 191 / 200) loss: 1.447238
(Epoch 5 / 5) train acc: 0.515000; val_acc: 0.384000
```

<ipython-input-16-239314d269f9>:39: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

```
plt.subplot(3, 1, 1)
```

<ipython-input-16-239314d269f9>:42: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

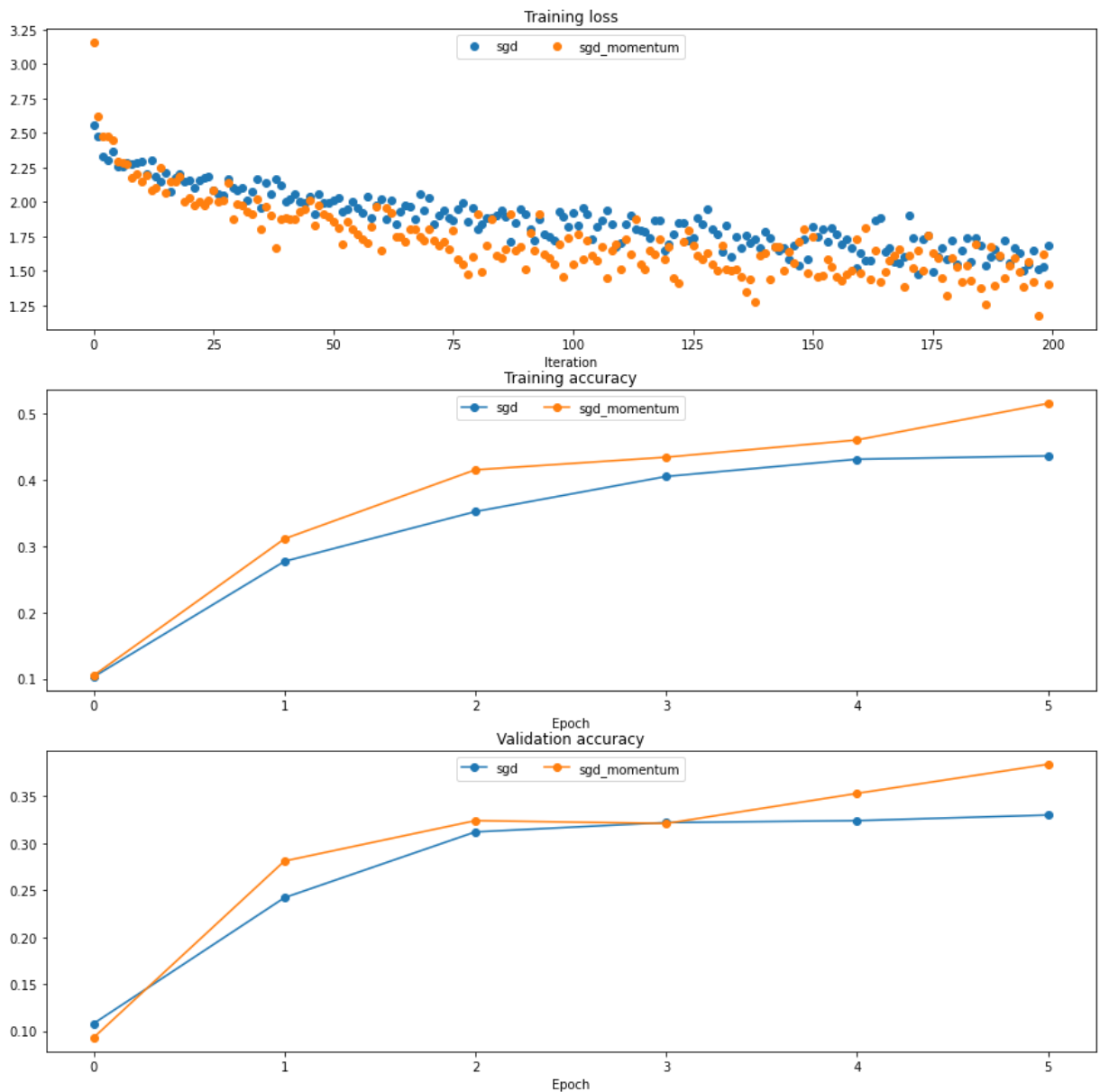
```
plt.subplot(3, 1, 2)
```

<ipython-input-16-239314d269f9>:45: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

```
plt.subplot(3, 1, 3)
```

<ipython-input-16-239314d269f9>:49: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

```
plt.subplot(3, 1, i)
```



RMSProp and Adam

RMSProp [1] and Adam [2] are update rules that set per-parameter learning rates by using a running average of the second moments of gradients.

In the file `cs231n/optim.py`, implement the RMSProp update rule in the `rmsprop` function and implement the Adam update rule in the `adam` function, and check your implementations using the tests below.

NOTE: Please implement the *complete* Adam update rule (with the bias correction mechanism), not the first simplified version mentioned in the course notes.

[1] Tijmen Tieleman and Geoffrey Hinton. "Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude." COURSE: Neural Networks for Machine Learning 4 (2012).

[2] Diederik Kingma and Jimmy Ba, "Adam: A Method for Stochastic Optimization", ICLR 2015.

```
In [17]: # Test RMSProp implementation
         from cs231n.optim import rmsprop
```

```

N, D = 4, 5
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
cache = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-2, 'cache': cache}
next_w, _ = rmsprop(w, dw, config=config)

expected_next_w = np.asarray([
    [-0.39223849, -0.34037513, -0.28849239, -0.23659121, -0.18467247],
    [-0.132737, -0.08078555, -0.02881884, 0.02316247, 0.07515774],
    [0.12716641, 0.17918792, 0.23122175, 0.28326742, 0.33532447],
    [0.38739248, 0.43947102, 0.49155973, 0.54365823, 0.59576619]])
expected_cache = np.asarray([
    [0.5976, 0.6126277, 0.6277108, 0.64284931, 0.65804321],
    [0.67329252, 0.68859723, 0.70395734, 0.71937285, 0.73484377],
    [0.75037008, 0.7659518, 0.78158892, 0.79728144, 0.81302936],
    [0.82883269, 0.84469141, 0.86060554, 0.87657507, 0.8926 ]])

# You should see relative errors around e-7 or less
print('next_w error: ', rel_error(expected_next_w, next_w))
print('cache error: ', rel_error(expected_cache, config['cache']))

```

```

next_w error: 9.524687511038133e-08
cache error: 2.6477955807156126e-09

```

In [18]:

```

# Test Adam implementation
from cs231n.optim import adam

N, D = 4, 5
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
m = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)
v = np.linspace(0.7, 0.5, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-2, 'm': m, 'v': v, 't': 5}
next_w, _ = adam(w, dw, config=config)

expected_next_w = np.asarray([
    [-0.40094747, -0.34836187, -0.29577703, -0.24319299, -0.19060977],
    [-0.1380274, -0.08544591, -0.03286534, 0.01971428, 0.0722929],
    [0.1248705, 0.17744702, 0.23002243, 0.28259667, 0.33516969],
    [0.38774145, 0.44031188, 0.49288093, 0.54544852, 0.59801459]])
expected_v = np.asarray([
    [0.69966, 0.68908382, 0.67851319, 0.66794809, 0.65738853],
    [0.64683452, 0.63628604, 0.6257431, 0.61520571, 0.60467385],
    [0.59414753, 0.58362676, 0.57311152, 0.56260183, 0.55209767],
    [0.54159906, 0.53110598, 0.52061845, 0.51013645, 0.49966, ]])
expected_m = np.asarray([
    [0.48, 0.49947368, 0.51894737, 0.53842105, 0.55789474],
    [0.57736842, 0.59684211, 0.61631579, 0.63578947, 0.65526316],
    [0.67473684, 0.69421053, 0.71368421, 0.73315789, 0.75263158],
    [0.77210526, 0.79157895, 0.81105263, 0.83052632, 0.85 ]])

# You should see relative errors around e-7 or less
print('next_w error: ', rel_error(expected_next_w, next_w))
print('v error: ', rel_error(expected_v, config['v']))
print('m error: ', rel_error(expected_m, config['m']))

```

```

next_w error: 1.1395691798535431e-07
v error: 4.208314038113071e-09
m error: 4.214963193114416e-09

```

Once you have debugged your RMSProp and Adam implementations, run the following to train

a pair of deep networks using these new update rules:

```
In [19]: learning_rates = {'rmsprop': 1e-4, 'adam': 1e-3}
for update_rule in ['adam', 'rmsprop']:
    print('running with ', update_rule)
    model = FullyConnectedNet([100, 100, 100, 100, 100], weight_scale=5e-2)

    solver = Solver(model, small_data,
                    num_epochs=5, batch_size=100,
                    update_rule=update_rule,
                    optim_config={
                        'learning_rate': learning_rates[update_rule]
                    },
                    verbose=True)
    solvers[update_rule] = solver
    solver.train()
    print()

plt.subplot(3, 1, 1)
plt.title('Training loss')
plt.xlabel('Iteration')

plt.subplot(3, 1, 2)
plt.title('Training accuracy')
plt.xlabel('Epoch')

plt.subplot(3, 1, 3)
plt.title('Validation accuracy')
plt.xlabel('Epoch')

for update_rule, solver in list(solvers.items()):
    plt.subplot(3, 1, 1)
    plt.plot(solver.loss_history, 'o', label=update_rule)

    plt.subplot(3, 1, 2)
    plt.plot(solver.train_acc_history, '-o', label=update_rule)

    plt.subplot(3, 1, 3)
    plt.plot(solver.val_acc_history, '-o', label=update_rule)

for i in [1, 2, 3]:
    plt.subplot(3, 1, i)
    plt.legend(loc='upper center', ncol=4)
plt.gcf().set_size_inches(15, 15)
plt.show()
```

```
running with adam
(Iteration 1 / 200) loss: 3.476928
(Epoch 0 / 5) train acc: 0.126000; val_acc: 0.110000
(Iteration 11 / 200) loss: 2.027712
(Iteration 21 / 200) loss: 2.183358
(Iteration 31 / 200) loss: 1.744257
(Epoch 1 / 5) train acc: 0.363000; val_acc: 0.330000
(Iteration 41 / 200) loss: 1.707951
(Iteration 51 / 200) loss: 1.703835
(Iteration 61 / 200) loss: 2.094758
(Iteration 71 / 200) loss: 1.505558
(Epoch 2 / 5) train acc: 0.419000; val_acc: 0.362000
(Iteration 81 / 200) loss: 1.594429
(Iteration 91 / 200) loss: 1.519016
(Iteration 101 / 200) loss: 1.368522
(Iteration 111 / 200) loss: 1.470400
(Epoch 3 / 5) train acc: 0.460000; val_acc: 0.378000
(Iteration 121 / 200) loss: 1.199064
```



```
(Iteration 131 / 200) loss: 1.464705
(Iteration 141 / 200) loss: 1.359863
(Iteration 151 / 200) loss: 1.415068
(Epoch 4 / 5) train acc: 0.521000; val_acc: 0.374000
(Iteration 161 / 200) loss: 1.382818
(Iteration 171 / 200) loss: 1.359900
(Iteration 181 / 200) loss: 1.095948
(Iteration 191 / 200) loss: 1.243088
(Epoch 5 / 5) train acc: 0.572000; val_acc: 0.382000
```

running with rmsprop

```
(Iteration 1 / 200) loss: 2.589166
(Epoch 0 / 5) train acc: 0.119000; val_acc: 0.146000
(Iteration 11 / 200) loss: 2.032921
(Iteration 21 / 200) loss: 1.897278
(Iteration 31 / 200) loss: 1.770793
(Epoch 1 / 5) train acc: 0.381000; val_acc: 0.320000
(Iteration 41 / 200) loss: 1.895732
(Iteration 51 / 200) loss: 1.681091
(Iteration 61 / 200) loss: 1.487204
(Iteration 71 / 200) loss: 1.629973
(Epoch 2 / 5) train acc: 0.429000; val_acc: 0.350000
(Iteration 81 / 200) loss: 1.506686
(Iteration 91 / 200) loss: 1.610742
(Iteration 101 / 200) loss: 1.486124
(Iteration 111 / 200) loss: 1.559454
(Epoch 3 / 5) train acc: 0.492000; val_acc: 0.359000
(Iteration 121 / 200) loss: 1.496860
(Iteration 131 / 200) loss: 1.531552
(Iteration 141 / 200) loss: 1.550195
(Iteration 151 / 200) loss: 1.657838
(Epoch 4 / 5) train acc: 0.533000; val_acc: 0.354000
(Iteration 161 / 200) loss: 1.603105
(Iteration 171 / 200) loss: 1.405372
(Iteration 181 / 200) loss: 1.503740
(Iteration 191 / 200) loss: 1.385278
(Epoch 5 / 5) train acc: 0.531000; val_acc: 0.374000
```

<ipython-input-19-c31f2247ce3b>:30: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

```
plt.subplot(3, 1, 1)
```

<ipython-input-19-c31f2247ce3b>:33: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

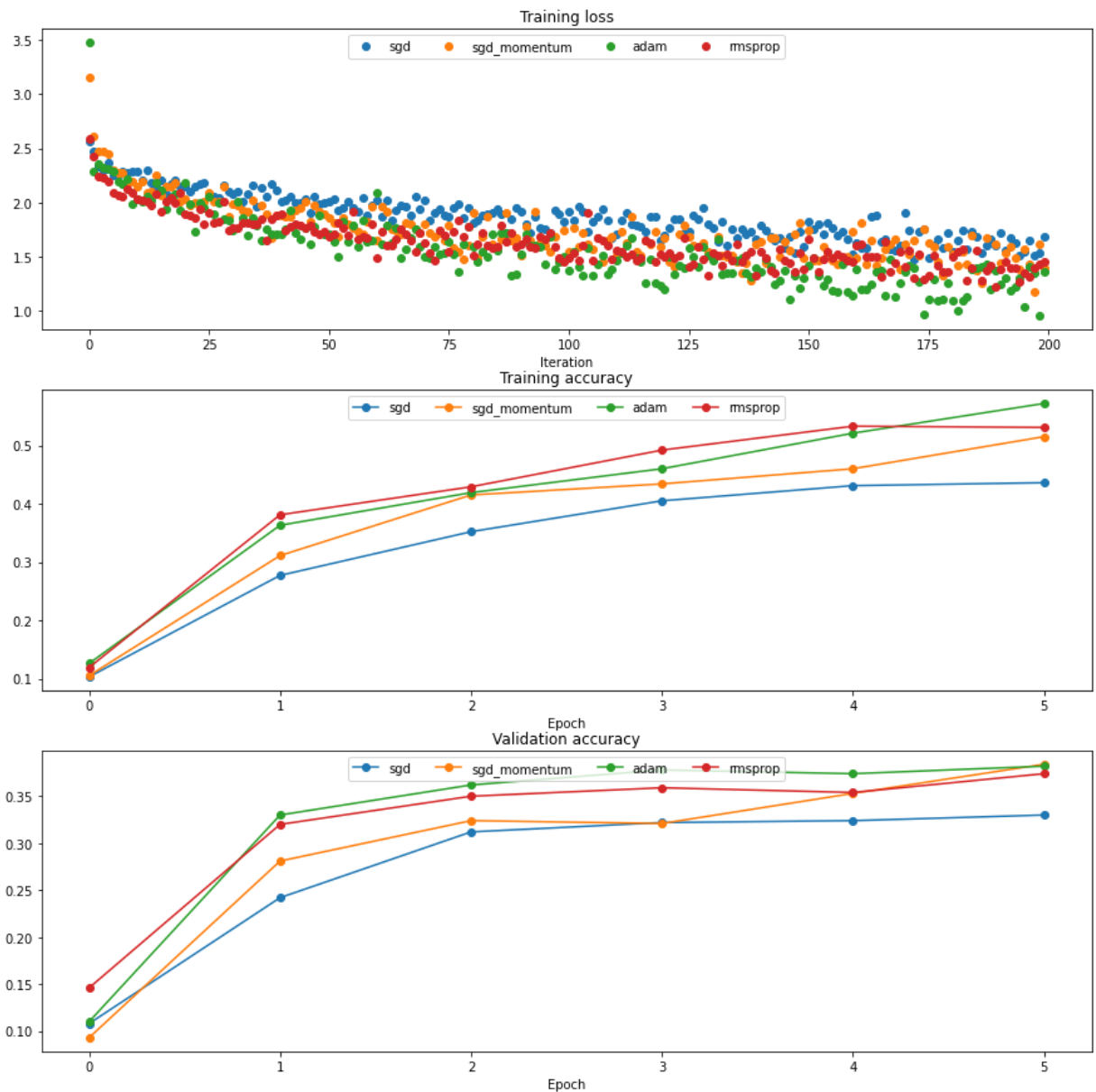
```
plt.subplot(3, 1, 2)
```

<ipython-input-19-c31f2247ce3b>:36: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

```
plt.subplot(3, 1, 3)
```

<ipython-input-19-c31f2247ce3b>:40: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

```
plt.subplot(3, 1, i)
```



Inline Question 3:

AdaGrad, like Adam, is a per-parameter optimization method that uses the following update rule:

```
cache += dw**2
w += - learning_rate * dw / (np.sqrt(cache) + eps)
```

John notices that when he was training a network with AdaGrad that the updates became very small, and that his network was learning slowly. Using your knowledge of the AdaGrad update rule, why do you think the updates would become very small? Would Adam have the same issue?

Answer:

In the update weight formula, in the division part, for the small gradient values, the square root of the cache becomes smaller and therefore the update value increases. For the higher gradient values, cache increases and makes the update value increasingly smaller. As a result, the cache

gets massively growing and makes the update value decrease. It is unpreventable in AdaGrad, but it is preventable in Adam because there is a decaying term in Adam.

Train a good model!

Train the best fully-connected model that you can on CIFAR-10, storing your best model in the `best_model` variable. We require you to get at least 50% accuracy on the validation set using a fully-connected net.

If you are careful it should be possible to get accuracies above 55%, but we don't require it for this part and won't assign extra credit for doing so. Later in the assignment we will ask you to train the best convolutional network that you can on CIFAR-10, and we would prefer that you spend your effort working on convolutional nets rather than fully-connected nets.

You might find it useful to complete the `BatchNormalization.ipynb` and `Dropout.ipynb` notebooks before completing this part, since those techniques can help you train powerful models.

In [20]:

```
best_model = None
#####
# TODO: Train the best FullyConnectedNet that you can on CIFAR-10. You might #
# find batch/layer normalization and dropout useful. Store your best model in #
# the best_model variable. #####

hidden_dims = [100] * 4

range_weight_scale = [1e-2, 2e-2, 5e-3]
range_lr = [1e-5, 5e-4, 1e-5]

best_val_acc = -1
best_weight_scale = 0
best_lr = 0

print("Training...")

for weight_scale in range_weight_scale:
    for lr in range_lr:
        model = FullyConnectedNet(hidden_dims=hidden_dims, reg=0.0,
                                   weight_scale=weight_scale)
        solver = Solver(model, data, update_rule='adam',
                         optim_config={'learning_rate': lr},
                         batch_size=100, num_epochs=5,
                         verbose=False)
        solver.train()
        val_acc = solver.best_val_acc

        print('Weight_scale: %f, lr: %f, val_acc: %f' % (weight_scale, lr, val_acc))

        if val_acc > best_val_acc:
            best_val_acc = val_acc
            best_weight_scale = weight_scale
            best_lr = lr
            best_model = model

print("Best val_acc: %f" % best_val_acc)
print("Best weight_scale: %f" % best_weight_scale)
```

```
print("Best lr: %f" % best_lr)
```

```
#####
#                                     END OF YOUR CODE                                #
#####
```

Training...

```
Weight_scale: 0.010000, lr: 0.000010, val_acc: 0.341000
Weight_scale: 0.010000, lr: 0.000500, val_acc: 0.492000
Weight_scale: 0.010000, lr: 0.000010, val_acc: 0.341000
Weight_scale: 0.020000, lr: 0.000010, val_acc: 0.426000
Weight_scale: 0.020000, lr: 0.000500, val_acc: 0.512000
Weight_scale: 0.020000, lr: 0.000010, val_acc: 0.421000
Weight_scale: 0.005000, lr: 0.000010, val_acc: 0.272000
Weight_scale: 0.005000, lr: 0.000500, val_acc: 0.517000
Weight_scale: 0.005000, lr: 0.000010, val_acc: 0.260000
Best val_acc: 0.517000
Best weight_scale: 0.005000
Best lr: 0.000500
```

Test your model!

Run your best model on the validation and test sets. You should achieve above 50% accuracy on the validation set.

In [21]:

```
y_test_pred = np.argmax(best_model.loss(data['X_test']), axis=1)
y_val_pred = np.argmax(best_model.loss(data['X_val']), axis=1)
print('Validation set accuracy: ', (y_val_pred == data['y_val']).mean())
print('Test set accuracy: ', (y_test_pred == data['y_test']).mean())
```

```
Validation set accuracy: 0.517
Test set accuracy: 0.507
```

Dropout

Dropout [1] is a technique for regularizing neural networks by randomly setting some features to zero during the forward pass. In this exercise you will implement a dropout layer and modify your fully-connected network to optionally use dropout.

[1] [Geoffrey E. Hinton et al, "Improving neural networks by preventing co-adaptation of feature detectors", arXiv 2012](#)

In [1]:

```
# As usual, a bit of setup
from __future__ import print_function
import time
import numpy as np
import matplotlib.pyplot as plt
from cs231n.classifiers.fc_net import *
from cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_gradient_a
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))))
```

run the following from the cs231n directory and try again:
python setup.py build_ext --inplace
You may also need to restart your iPython kernel

In [2]:

```
# Load the (preprocessed) CIFAR10 data.

data = get_CIFAR10_data()
for k, v in data.items():
    print('%s: ' % k, v.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,)
```

Dropout forward pass

In the file `cs231n/layers.py`, implement the forward pass for dropout. Since dropout behaves differently during training and testing, make sure to implement the operation for both modes.

Once you have done so, run the cell below to test your implementation.

```
In [3]: np.random.seed(231)
x = np.random.randn(500, 500) + 10

for p in [0.25, 0.4, 0.7]:
    out, _ = dropout_forward(x, {'mode': 'train', 'p': p})
    out_test, _ = dropout_forward(x, {'mode': 'test', 'p': p})

    print('Running tests with p = ', p)
    print('Mean of input: ', x.mean())
    print('Mean of train-time output: ', out.mean())
    print('Mean of test-time output: ', out_test.mean())
    print('Fraction of train-time output set to zero: ', (out == 0).mean())
    print('Fraction of test-time output set to zero: ', (out_test == 0).mean())
    print()
```

```
Running tests with p = 0.25
Mean of input: 10.000207878477502
Mean of train-time output: 10.014059116977283
Mean of test-time output: 10.000207878477502
Fraction of train-time output set to zero: 0.749784
Fraction of test-time output set to zero: 0.0
```

```
Running tests with p = 0.4
Mean of input: 10.000207878477502
Mean of train-time output: 9.977917658761159
Mean of test-time output: 10.000207878477502
Fraction of train-time output set to zero: 0.600796
Fraction of test-time output set to zero: 0.0
```

```
Running tests with p = 0.7
Mean of input: 10.000207878477502
Mean of train-time output: 9.987811912159426
Mean of test-time output: 10.000207878477502
Fraction of train-time output set to zero: 0.30074
Fraction of test-time output set to zero: 0.0
```

Dropout backward pass

In the file `cs231n/layers.py`, implement the backward pass for dropout. After doing so, run the following cell to numerically gradient-check your implementation.

```
In [4]: np.random.seed(231)
x = np.random.randn(10, 10) + 10
dout = np.random.randn(*x.shape)

dropout_param = {'mode': 'train', 'p': 0.2, 'seed': 123}
out, cache = dropout_forward(x, dropout_param)
dx = dropout_backward(dout, cache)
dx_num = eval_numerical_gradient_array(lambda xx: dropout_forward(xx, dropout_param)

# Error should be around e-10 or less
print('dx relative error: ', rel_error(dx, dx_num))
```

```
dx relative error: 5.44560814873387e-11
```

Inline Question 1:

What happens if we do not divide the values being passed through inverse dropout by p in the dropout layer? Why does that happen?

Answer:

It is wanted that the dropout function to output on a similar scale to the input to have a consistent network. The mask is divided by p to keep the expected values same. If the input and output averages of the dropout function are compared, the reason of the division process by p can be understood. When dividing by p , the average of the output terms becomes 10. If not dividing by p , the average of the output terms is 3, whereas the average of the input terms is 10. The output average terms when dividing by p becomes similar to the input average terms but they become inconsistent if we do not divide by p . Even if the output neurons are dropping by dividing, the expected value of the dropout output is kept same. Therefore, the network becomes consistent if the dropout is used.

Fully-connected nets with Dropout

In the file `cs231n/classifiers/fc_net.py`, modify your implementation to use dropout. Specifically, if the constructor of the net receives a value that is not 1 for the `dropout` parameter, then the net should add dropout immediately after every ReLU nonlinearity. After doing so, run the following to numerically gradient-check your implementation.

In [5]:

```
np.random.seed(231)
N, D, H1, H2, C = 2, 15, 20, 30, 10
X = np.random.randn(N, D)
y = np.random.randint(C, size=(N,))

for dropout in [1, 0.75, 0.5]:
    print('Running check with dropout = ', dropout)
    model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
                              weight_scale=5e-2, dtype=np.float64,
                              dropout=dropout, seed=123)

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

    # Relative errors should be around e-6 or less; Note that it's fine
    # if for dropout=1 you have W2 error be on the order of e-5.
    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('%s relative error: %.2e' % (name, rel_error(grad_num, grads[name])))
    print()
```

```
Running check with dropout = 1
Initial loss: 2.3004790897684924
W1 relative error: 1.48e-07
W2 relative error: 2.21e-05
W3 relative error: 3.53e-07
b1 relative error: 5.38e-09
b2 relative error: 2.09e-09
b3 relative error: 5.80e-11
```

```
Running check with dropout = 0.75
Initial loss: 2.302371489704412
W1 relative error: 1.90e-07
```

```

W2 relative error: 4.76e-06
W3 relative error: 2.60e-08
b1 relative error: 4.73e-09
b2 relative error: 1.82e-09
b3 relative error: 1.70e-10

```

```

Running check with dropout = 0.5
Initial loss: 2.3042759220785896
W1 relative error: 3.11e-07
W2 relative error: 1.84e-08
W3 relative error: 5.35e-08
b1 relative error: 2.58e-08
b2 relative error: 2.99e-09
b3 relative error: 1.13e-10

```

Regularization experiment

As an experiment, we will train a pair of two-layer networks on 500 training examples: one will use no dropout, and one will use a keep probability of 0.25. We will then visualize the training and validation accuracies of the two networks over time.

```

In [6]: # Train two identical nets, one with dropout and one without
np.random.seed(231)
num_train = 500
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'],
}

solvers = {}
dropout_choices = [1, 0.25]
for dropout in dropout_choices:
    model = FullyConnectedNet([500], dropout=dropout)
    print(dropout)

    solver = Solver(model, small_data,
                    num_epochs=25, batch_size=100,
                    update_rule='adam',
                    optim_config={
                        'learning_rate': 5e-4,
                    },
                    verbose=True, print_every=100)
    solver.train()
    solvers[dropout] = solver

```

```

1
(Iteration 1 / 125) loss: 7.856644
(Epoch 0 / 25) train acc: 0.260000; val_acc: 0.184000
(Epoch 1 / 25) train acc: 0.416000; val_acc: 0.258000
(Epoch 2 / 25) train acc: 0.482000; val_acc: 0.276000
(Epoch 3 / 25) train acc: 0.532000; val_acc: 0.277000
(Epoch 4 / 25) train acc: 0.600000; val_acc: 0.271000
(Epoch 5 / 25) train acc: 0.708000; val_acc: 0.299000
(Epoch 6 / 25) train acc: 0.722000; val_acc: 0.282000
(Epoch 7 / 25) train acc: 0.832000; val_acc: 0.255000
(Epoch 8 / 25) train acc: 0.878000; val_acc: 0.269000
(Epoch 9 / 25) train acc: 0.902000; val_acc: 0.275000
(Epoch 10 / 25) train acc: 0.888000; val_acc: 0.261000
(Epoch 11 / 25) train acc: 0.926000; val_acc: 0.278000
(Epoch 12 / 25) train acc: 0.960000; val_acc: 0.302000

```



```
(Epoch 13 / 25) train acc: 0.964000; val_acc: 0.306000
(Epoch 14 / 25) train acc: 0.966000; val_acc: 0.309000
(Epoch 15 / 25) train acc: 0.976000; val_acc: 0.289000
(Epoch 16 / 25) train acc: 0.988000; val_acc: 0.302000
(Epoch 17 / 25) train acc: 0.988000; val_acc: 0.309000
(Epoch 18 / 25) train acc: 0.990000; val_acc: 0.312000
(Epoch 19 / 25) train acc: 0.990000; val_acc: 0.316000
(Epoch 20 / 25) train acc: 0.988000; val_acc: 0.313000
(Iteration 101 / 125) loss: 0.117108
(Epoch 21 / 25) train acc: 0.992000; val_acc: 0.297000
(Epoch 22 / 25) train acc: 0.980000; val_acc: 0.302000
(Epoch 23 / 25) train acc: 0.984000; val_acc: 0.294000
(Epoch 24 / 25) train acc: 0.994000; val_acc: 0.303000
(Epoch 25 / 25) train acc: 0.992000; val_acc: 0.289000
0.25
(Iteration 1 / 125) loss: 17.318479
(Epoch 0 / 25) train acc: 0.230000; val_acc: 0.177000
(Epoch 1 / 25) train acc: 0.378000; val_acc: 0.243000
(Epoch 2 / 25) train acc: 0.402000; val_acc: 0.254000
(Epoch 3 / 25) train acc: 0.502000; val_acc: 0.276000
(Epoch 4 / 25) train acc: 0.528000; val_acc: 0.298000
(Epoch 5 / 25) train acc: 0.562000; val_acc: 0.297000
(Epoch 6 / 25) train acc: 0.628000; val_acc: 0.291000
(Epoch 7 / 25) train acc: 0.622000; val_acc: 0.299000
(Epoch 8 / 25) train acc: 0.684000; val_acc: 0.311000
(Epoch 9 / 25) train acc: 0.714000; val_acc: 0.291000
(Epoch 10 / 25) train acc: 0.724000; val_acc: 0.297000
(Epoch 11 / 25) train acc: 0.762000; val_acc: 0.303000
(Epoch 12 / 25) train acc: 0.786000; val_acc: 0.284000
(Epoch 13 / 25) train acc: 0.824000; val_acc: 0.303000
(Epoch 14 / 25) train acc: 0.826000; val_acc: 0.347000
(Epoch 15 / 25) train acc: 0.862000; val_acc: 0.349000
(Epoch 16 / 25) train acc: 0.862000; val_acc: 0.315000
(Epoch 17 / 25) train acc: 0.870000; val_acc: 0.307000
(Epoch 18 / 25) train acc: 0.876000; val_acc: 0.340000
(Epoch 19 / 25) train acc: 0.886000; val_acc: 0.327000
(Epoch 20 / 25) train acc: 0.880000; val_acc: 0.322000
(Iteration 101 / 125) loss: 3.870114
(Epoch 21 / 25) train acc: 0.894000; val_acc: 0.325000
(Epoch 22 / 25) train acc: 0.898000; val_acc: 0.306000
(Epoch 23 / 25) train acc: 0.878000; val_acc: 0.306000
(Epoch 24 / 25) train acc: 0.914000; val_acc: 0.321000
(Epoch 25 / 25) train acc: 0.910000; val_acc: 0.328000
```

In [7]:

```
# Plot train and validation accuracies of the two models

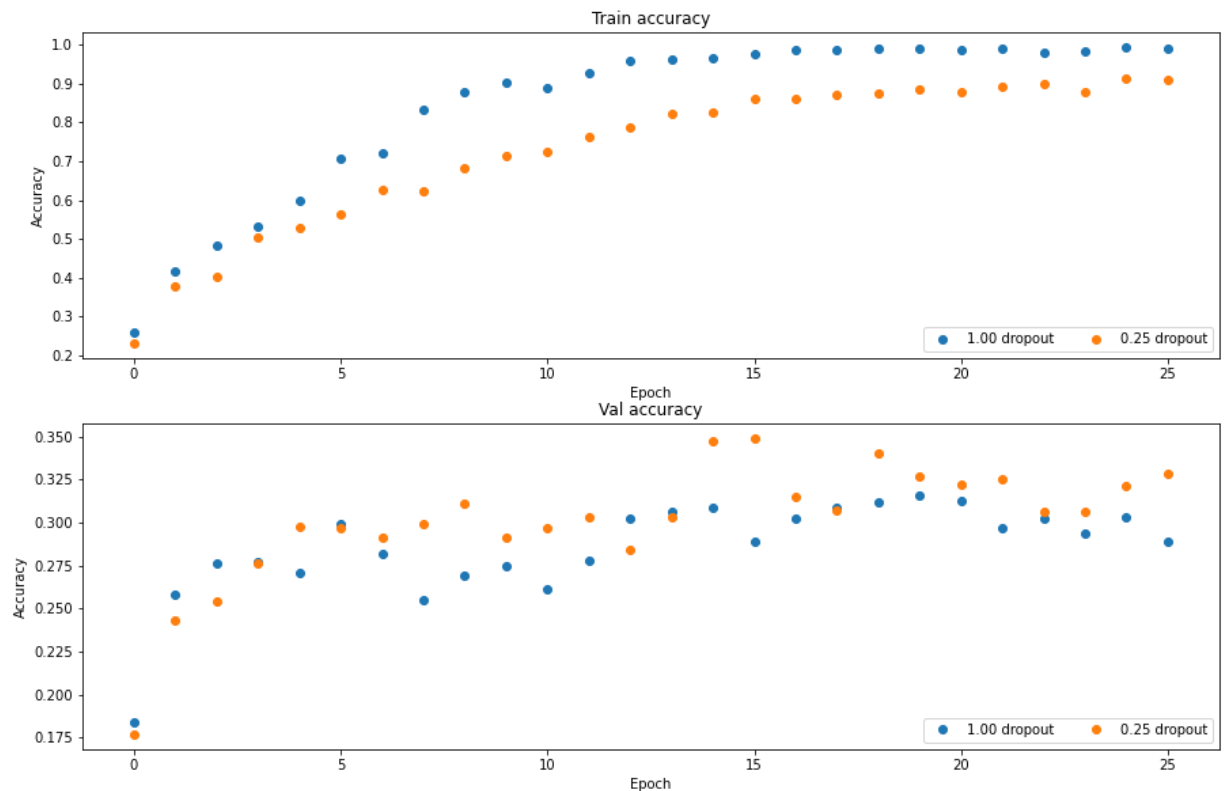
train_accs = []
val_accs = []
for dropout in dropout_choices:
    solver = solvers[dropout]
    train_accs.append(solver.train_acc_history[-1])
    val_accs.append(solver.val_acc_history[-1])

plt.subplot(3, 1, 1)
for dropout in dropout_choices:
    plt.plot(solvers[dropout].train_acc_history, 'o', label='%.2f dropout' % dropout)
plt.title('Train accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend(ncol=2, loc='lower right')

plt.subplot(3, 1, 2)
for dropout in dropout_choices:
    plt.plot(solvers[dropout].val_acc_history, 'o', label='%.2f dropout' % dropout)
plt.title('Val accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
```

```
plt.legend(ncol=2, loc='lower right')

plt.gcf().set_size_inches(15, 15)
plt.show()
```



Inline Question 2:

Compare the validation and training accuracies with and without dropout -- what do your results suggest about dropout as a regularizer?

Answer:

The results show that the dropout is preventing the overfitting of the training data. If the dropout is used in the network, the network will have a training accuracy of 91%. If there is no dropout in the network, the training accuracy becomes nearly 99%. Therefore, the network overfits the data when there is no dropout in the network. Accordingly, if we measure the validation accuracy in network with dropout, it becomes 32.8%. Without dropout, the validation accuracy becomes nearly 29%. Therefore, it can be said that the network with dropout is able to predict the validation data better. Lastly, plots also show these facts that validation accuracy with dropout is usually higher than the network without dropout.

Inline Question 3:

Suppose we are training a deep fully-connected network for image classification, with dropout after hidden layers (parameterized by keep probability p). How should we modify p , if at all, if we decide to decrease the size of the hidden layers (that is, the number of nodes in each layer)?

Answer:

The keep probability p value should be increased. As in first question of the assignment, it is seen that if there are more hidden layers, the learning becomes better. However, if there are more layers in the network, the overfitting fact can be occur. Training data may result in overfitting as the layer amount increases. So, there is a trade-off. We can control the p so that if there are many hidden layers that can cause overfitting, the p value should be chosen small. However, if the amount of hidden layer is decreased that can result in decreasing probability of overfitting, p value can be increased to make the learning process stay in a steady speed.