

# ECE C147/247 HW4 Q1: Optimization for Fully Connected Networks

In this notebook, we will implement different optimization rules for gradient descent. We have provided starter code; however, you will need to copy and paste your code from your implementation of the modular fully connected nets in HW #3 to build upon this.

`utils` has built a solid API for building these modular frameworks and training them, and we will use their very well implemented framework as opposed to "reinventing the wheel." This includes using their Solver, various utility functions, and their layer structure. This also includes `nndl.fc_net`, `nndl.layers`, and `nndl.layer_utils`.

```
In [1]: ## Import and setups

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

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

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

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

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

data = get_CIFAR10_data()
for k in data.keys():
    print('{}: {}'.format(k, data[k].shape))
```

```
X_train: (49000, 3, 32, 32)
y_train: (49000,)
X_val: (1000, 3, 32, 32)
y_val: (1000,)
X_test: (1000, 3, 32, 32)
y_test: (1000,)
```

## Building upon your HW #3 implementation

Copy and paste the following functions from your HW #3 implementation of a modular FC net:

- `affine_forward` in `nndl/layers.py`
- `affine_backward` in `nndl/layers.py`

- relu\_forward in nndl/layers.py
- relu\_backward in nndl/layers.py
- affine\_relu\_forward in nndl/layer\_utils.py
- affine\_relu\_backward in nndl/layer\_utils.py
- The FullyConnectedNet class in nndl/fc\_net.py

## Test all functions you copy and pasted

In [3]:

```
from nndl.layer_tests import *

affine_forward_test(); print('\n')
affine_backward_test(); print('\n')
relu_forward_test(); print('\n')
relu_backward_test(); print('\n')
affine_relu_test(); print('\n')
fc_net_test()
```

If affine\_forward function is working, difference should be less than 1e-9:  
difference: 9.769849468192957e-10

If affine\_backward is working, error should be less than 1e-9::  
dx error: 4.517717286767966e-10  
dw error: 2.953555218603546e-10  
db error: 1.3004066660747692e-11

If relu\_forward function is working, difference should be around 1e-8:  
difference: 4.999999798022158e-08

If relu\_backward function is working, error should be less than 1e-9:  
dx error: 3.2756114574567796e-12

If affine\_relu\_forward and affine\_relu\_backward are working, error should be less than 1e-9::  
dx error: 1.797929234036535e-10  
dw error: 1.510465039305725e-10  
db error: 3.2756019865460813e-12

Running check with reg = 0  
Initial loss: 2.300223304820917  
W1 relative error: 4.631040062839554e-08  
W2 relative error: 1.2098517344039791e-06  
W3 relative error: 3.738792456207884e-07  
b1 relative error: 2.0538755808609935e-08  
b2 relative error: 4.805233546659617e-09  
b3 relative error: 1.4364540241108383e-10  
Running check with reg = 3.14  
Initial loss: 6.742053383557732  
W1 relative error: 2.0360662169265662e-08  
W2 relative error: 1.170895493191328e-06  
W3 relative error: 5.248572999147182e-09  
b1 relative error: 2.778938800619773e-07  
b2 relative error: 5.317300046122514e-09  
b3 relative error: 1.296243952406403e-10

## Training a larger model

In general, proceeding with vanilla stochastic gradient descent to optimize models may be fraught with problems and limitations, as discussed in class. Thus, we implement optimizers that improve on SGD.

## SGD + momentum

In the following section, implement SGD with momentum. Read the `nndl/optim.py` API, and be sure you understand it. After, implement `sgd_momentum` in `nndl/optim.py`. Test your implementation of `sgd_momentum` by running the cell below.

```
In [4]: from nndl.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      ]])

print('next_w error: {}'.format(rel_error(next_w, expected_next_w)))
print('velocity error: {}'.format(rel_error(expected_velocity, config['velocity'])))

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

## SGD + Nesterov momentum

Implement `sgd_nesterov_momentum` in `ndl/optim.py`.

```
In [5]: from nndl.optim import sgd_nesterov_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_nesterov_momentum(w, dw, config=config)

expected_next_w = np.asarray([
    [0.08714,      0.15246105,  0.21778211,  0.28310316,  0.34842421],
    [0.41374526,  0.47906632,  0.54438737,  0.60970842,  0.67502947],
    [0.74035053,  0.80567158,  0.87099263,  0.93631368,  1.00163474],
    [1.06695579,  1.13227684,  1.19759789,  1.26291895,  1.32824      ]])
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      ]])
```

```
print('next_w error: {}'.format(rel_error(next_w, expected_next_w)))
print('velocity error: {}'.format(rel_error(expected_velocity, config['velocity'])))
```

```
next_w error: 1.0875186845081027e-08
velocity error: 4.269287743278663e-09
```

## Evaluating SGD, SGD+Momentum, and SGD+NesterovMomentum

Run the following cell to train a 6 layer FC net with SGD, SGD+momentum, and SGD+Nesterov momentum. You should see that SGD+momentum achieves a better loss than SGD, and that SGD+Nesterov momentum achieves a slightly better loss (and training accuracy) than SGD+momentum.

In [6]:

```
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', 'sgd_nesterov_momentum']:
    print('Optimizing with {}'.format(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=False)
    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 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)

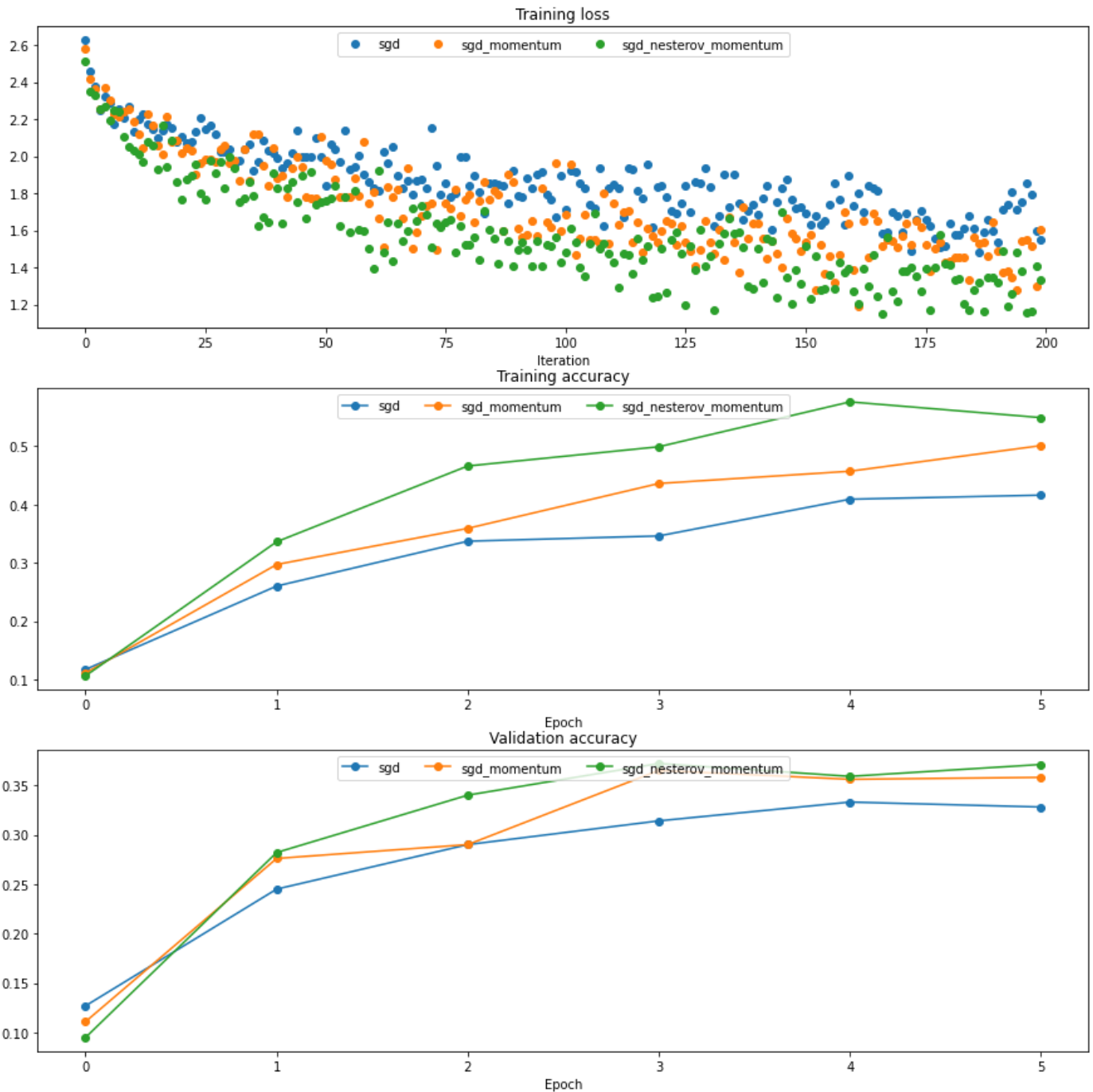
for i in [1, 2, 3]:
```

```
plt.subplot(3, 1, i)
plt.legend(loc='upper center', ncol=4)
plt.gcf().set_size_inches(15, 15)
plt.show()
```

Optimizing with `sgd`

Optimizing with `sgd_momentum`

Optimizing with `sgd_nesterov_momentum`



## RMSProp

Now we go to techniques that adapt the gradient. Implement `rmsprop` in `nndl/optim.py`. Test your implementation by running the cell below.

In [7]:

```
from nndl.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)
a = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)
```

```

config = {'learning_rate': 1e-2, 'a': a}
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 ]])

print('next_w error: {}'.format(rel_error(expected_next_w, next_w)))
print('cache error: {}'.format(rel_error(expected_cache, config['a'])))

```

```

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

```

## Adaptive moments

Now, implement adam in `nndl/optim.py`. Test your implementation by running the cell below.

In [8]:

```

# Test Adam implementation; you should see errors around 1e-7 or less
from nndl.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)
v = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)
a = np.linspace(0.7, 0.5, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-2, 'v': v, 'a': a, '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_a = 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_v = 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 ]])

print('next_w error: {}'.format(rel_error(expected_next_w, next_w)))
print('a error: {}'.format(rel_error(expected_a, config['a'])))
print('v error: {}'.format(rel_error(expected_v, config['v'])))

```

```

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

```

## Comparing SGD, SGD+NesterovMomentum, RMSProp, and

# Adam

The following code will compare optimization with SGD, Momentum, Nesterov Momentum, RMSProp and Adam. In our code, we find that RMSProp, Adam, and SGD + Nesterov Momentum achieve approximately the same training error after a few training epochs.

In [9]:

```
learning_rates = {'rmsprop': 2e-4, 'adam': 1e-3}

for update_rule in ['adam', 'rmsprop']:
    print('Optimizing with {}'.format(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=False)
    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 solvers.items():
    plt.subplot(3, 1, 1)
    plt.plot(solver.loss_history, 'o', label=update_rule)

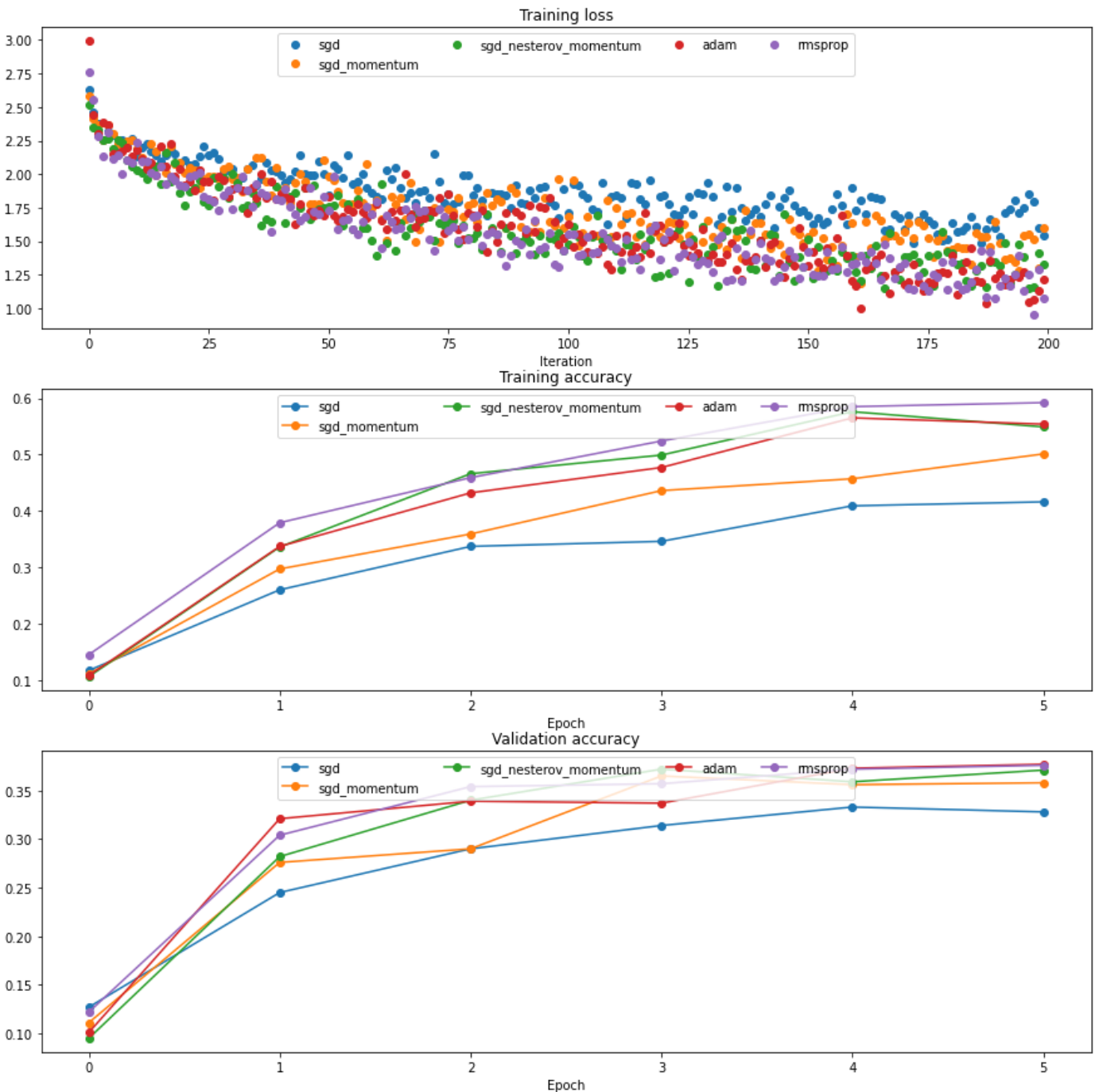
    plt.subplot(3, 1, 2)
    plt.plot(solver.train_acc_history, '-o', label=update_rule)

    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()
```

Optimizing with adam

Optimizing with rmsprop



## Easier optimization

In the following cell, we'll train a 4 layer neural network having 500 units in each hidden layer with the different optimizers, and find that it is far easier to get up to 50+% performance on CIFAR-10. After we implement batchnorm and dropout, we'll ask you to get 55+% on CIFAR-10.

In [10]:

```
optimizer = 'adam'
best_model = None

layer_dims = [500, 500, 500]
weight_scale = 0.01
learning_rate = 1e-3
lr_decay = 0.9

model = FullyConnectedNet(layer_dims, weight_scale=weight_scale,
                           use_batchnorm=True)

solver = Solver(model, data,
```



```
num_epochs=10, batch_size=100,  
update_rule=optimizer,  
optim_config={  
    'learning_rate': learning_rate,  
},  
lr_decay=lr_decay,  
verbose=True, print_every=50)  
solver.train()
```

```
(Iteration 1 / 4900) loss: 2.300475
(Epoch 0 / 10) train acc: 0.165000; val_acc: 0.181000
(Iteration 51 / 4900) loss: 1.932997
(Iteration 101 / 4900) loss: 1.698234
(Iteration 151 / 4900) loss: 1.781077
(Iteration 201 / 4900) loss: 1.744266
(Iteration 251 / 4900) loss: 1.705927
(Iteration 301 / 4900) loss: 1.526387
(Iteration 351 / 4900) loss: 1.702914
(Iteration 401 / 4900) loss: 1.566511
(Iteration 451 / 4900) loss: 1.707089
(Epoch 1 / 10) train acc: 0.399000; val_acc: 0.421000
(Iteration 501 / 4900) loss: 1.556329
(Iteration 551 / 4900) loss: 1.699830
(Iteration 601 / 4900) loss: 1.511032
(Iteration 651 / 4900) loss: 1.479937
(Iteration 701 / 4900) loss: 1.394786
(Iteration 751 / 4900) loss: 1.428443
(Iteration 801 / 4900) loss: 1.461454
(Iteration 851 / 4900) loss: 1.469750
(Iteration 901 / 4900) loss: 1.500166
(Iteration 951 / 4900) loss: 1.564140
(Epoch 2 / 10) train acc: 0.475000; val_acc: 0.470000
(Iteration 1001 / 4900) loss: 1.559267
(Iteration 1051 / 4900) loss: 1.301459
(Iteration 1101 / 4900) loss: 1.494965
(Iteration 1151 / 4900) loss: 1.533820
(Iteration 1201 / 4900) loss: 1.460437
(Iteration 1251 / 4900) loss: 1.357185
(Iteration 1301 / 4900) loss: 1.399734
(Iteration 1351 / 4900) loss: 1.423651
(Iteration 1401 / 4900) loss: 1.289892
(Iteration 1451 / 4900) loss: 1.413948
(Epoch 3 / 10) train acc: 0.492000; val_acc: 0.457000
(Iteration 1501 / 4900) loss: 1.088179
(Iteration 1551 / 4900) loss: 1.374607
(Iteration 1601 / 4900) loss: 1.378255
(Iteration 1651 / 4900) loss: 1.387306
(Iteration 1701 / 4900) loss: 1.296451
(Iteration 1751 / 4900) loss: 1.402440
(Iteration 1801 / 4900) loss: 1.545424
(Iteration 1851 / 4900) loss: 1.205041
(Iteration 1901 / 4900) loss: 1.341622
(Iteration 1951 / 4900) loss: 1.407368
(Epoch 4 / 10) train acc: 0.533000; val_acc: 0.494000
(Iteration 2001 / 4900) loss: 1.243661
(Iteration 2051 / 4900) loss: 1.226753
(Iteration 2101 / 4900) loss: 1.248218
(Iteration 2151 / 4900) loss: 1.168664
(Iteration 2201 / 4900) loss: 1.329747
(Iteration 2251 / 4900) loss: 1.091494
(Iteration 2301 / 4900) loss: 1.298749
(Iteration 2351 / 4900) loss: 1.169799
(Iteration 2401 / 4900) loss: 1.291062
(Epoch 5 / 10) train acc: 0.568000; val_acc: 0.517000
(Iteration 2451 / 4900) loss: 1.090569
(Iteration 2501 / 4900) loss: 1.040105
(Iteration 2551 / 4900) loss: 1.196661
(Iteration 2601 / 4900) loss: 1.184482
(Iteration 2651 / 4900) loss: 1.248305
(Iteration 2701 / 4900) loss: 0.966165
(Iteration 2751 / 4900) loss: 1.073918
(Iteration 2801 / 4900) loss: 1.219766
(Iteration 2851 / 4900) loss: 1.084969
(Iteration 2901 / 4900) loss: 1.214645
(Epoch 6 / 10) train acc: 0.576000; val_acc: 0.527000
```

```
(Iteration 2951 / 4900) loss: 1.353390
(Iteration 3001 / 4900) loss: 1.289739
(Iteration 3051 / 4900) loss: 0.962614
(Iteration 3101 / 4900) loss: 1.228209
(Iteration 3151 / 4900) loss: 1.111609
(Iteration 3201 / 4900) loss: 1.097480
(Iteration 3251 / 4900) loss: 1.163739
(Iteration 3301 / 4900) loss: 0.973220
(Iteration 3351 / 4900) loss: 1.050828
(Iteration 3401 / 4900) loss: 1.176953
(Epoch 7 / 10) train acc: 0.629000; val_acc: 0.522000
(Iteration 3451 / 4900) loss: 1.084450
(Iteration 3501 / 4900) loss: 1.093541
(Iteration 3551 / 4900) loss: 1.387569
(Iteration 3601 / 4900) loss: 1.125746
(Iteration 3651 / 4900) loss: 1.075082
(Iteration 3701 / 4900) loss: 1.066931
(Iteration 3751 / 4900) loss: 0.909745
(Iteration 3801 / 4900) loss: 0.947617
(Iteration 3851 / 4900) loss: 1.034305
(Iteration 3901 / 4900) loss: 1.090769
(Epoch 8 / 10) train acc: 0.667000; val_acc: 0.558000
(Iteration 3951 / 4900) loss: 0.995495
(Iteration 4001 / 4900) loss: 1.008144
(Iteration 4051 / 4900) loss: 0.943819
(Iteration 4101 / 4900) loss: 0.969136
(Iteration 4151 / 4900) loss: 1.047509
(Iteration 4201 / 4900) loss: 0.908712
(Iteration 4251 / 4900) loss: 0.874539
(Iteration 4301 / 4900) loss: 0.944562
(Iteration 4351 / 4900) loss: 0.856182
(Iteration 4401 / 4900) loss: 0.962699
(Epoch 9 / 10) train acc: 0.660000; val_acc: 0.531000
(Iteration 4451 / 4900) loss: 0.945672
(Iteration 4501 / 4900) loss: 0.901853
(Iteration 4551 / 4900) loss: 0.789863
(Iteration 4601 / 4900) loss: 0.997731
(Iteration 4651 / 4900) loss: 0.946927
(Iteration 4701 / 4900) loss: 0.762023
(Iteration 4751 / 4900) loss: 0.799917
(Iteration 4801 / 4900) loss: 0.892969
(Iteration 4851 / 4900) loss: 0.892236
(Epoch 10 / 10) train acc: 0.668000; val_acc: 0.534000
```

In [11]:

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

Validation set accuracy: 0.558

Test set accuracy: 0.533

```
1 import numpy as np
2
3
4 """
5 This file implements various first-order update rules that are commonly used
6 for
7 training neural networks. Each update rule accepts current weights and the
8 gradient of the loss with respect to those weights and produces the next set
9 of
10 weights. Each update rule has the same interface:
11
12 def update(w, dw, config=None):
13
14 Inputs:
15 - w: A numpy array giving the current weights.
16 - dw: A numpy array of the same shape as w giving the gradient of the
17 loss with respect to w.
18 - config: A dictionary containing hyperparameter values such as learning
19 rate,
20 momentum, etc. If the update rule requires caching values over many
21 iterations, then config will also hold these cached values.
22
23 Returns:
24 - next_w: The next point after the update.
25 - config: The config dictionary to be passed to the next iteration of the
26 update rule.
27
28 NOTE: For most update rules, the default learning rate will probably not
29 perform
30 well; however the default values of the other hyperparameters should work
31 well
32 for a variety of different problems.
33
34 For efficiency, update rules may perform in-place updates, mutating w and
35 setting next_w equal to w.
36 """
37
38 def sgd(w, dw, config=None):
39     """
40     Performs vanilla stochastic gradient descent.
41
42     config format:
43     - learning_rate: Scalar learning rate.
44     """
45     if config is None: config = {}
46     config.setdefault('learning_rate', 1e-2)
47
48     w -= config['learning_rate'] * dw
49     return w, config
50
51 def sgd_momentum(w, dw, config=None):
52     """
53     Performs stochastic gradient descent with momentum.
54
55     config format:
56     - learning_rate: Scalar learning rate.
57     - momentum: Scalar between 0 and 1 giving the momentum value.
```

```

55     Setting momentum = 0 reduces to sgd.
56     - velocity: A numpy array of the same shape as w and dw used to store a
moving
57     average of the gradients.
58     """
59     if config is None: config = {}
60     config.setdefault('learning_rate', 1e-2)
61     config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
62     v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets
it to zero.
63
64     # ===== #
65     # YOUR CODE HERE:
66     # Implement the momentum update formula. Return the updated weights
67     # as next_w, and the updated velocity as v.
68     # ===== #
69
70     alpha = config['momentum']
71     eps = config['learning_rate']
72
73     v = alpha * v - eps * dw
74     w += v
75
76     next_w = w
77
78     # ===== #
79     # END YOUR CODE HERE
80     # ===== #
81
82     config['velocity'] = v
83
84     return next_w, config
85
86 def sgd_nesterov_momentum(w, dw, config=None):
87     """
88     Performs stochastic gradient descent with Nesterov momentum.
89
90     config format:
91     - learning_rate: Scalar learning rate.
92     - momentum: Scalar between 0 and 1 giving the momentum value.
93     Setting momentum = 0 reduces to sgd.
94     - velocity: A numpy array of the same shape as w and dw used to store a
moving
95     average of the gradients.
96     """
97     if config is None: config = {}
98     config.setdefault('learning_rate', 1e-2)
99     config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
100    v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets
it to zero.
101
102    # ===== #
103    # YOUR CODE HERE:
104    # Implement the momentum update formula. Return the updated weights
105    # as next_w, and the updated velocity as v.
106    # ===== #
107
108    alpha = config['momentum']
109    eps = config['learning_rate']
110

```

```

111     v_old = v
112     v = alpha * v - eps * dw
113     w += (v + alpha * (v - v_old))
114
115     next_w = w
116
117     # ===== #
118     # END YOUR CODE HERE
119     # ===== #
120
121     config['velocity'] = v
122
123     return next_w, config
124
125 def rmsprop(w, dw, config=None):
126     """
127     Uses the RMSProp update rule, which uses a moving average of squared
128     gradient
129     values to set adaptive per-parameter learning rates.
130
131     config format:
132     - learning_rate: Scalar learning rate.
133     - decay_rate: Scalar between 0 and 1 giving the decay rate for the squared
134     gradient cache.
135     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
136     - beta: Moving average of second moments of gradients.
137     """
138     if config is None: config = {}
139     config.setdefault('learning_rate', 1e-2)
140     config.setdefault('decay_rate', 0.99)
141     config.setdefault('epsilon', 1e-8)
142     config.setdefault('a', np.zeros_like(w))
143
144     next_w = None
145
146     # ===== #
147     # YOUR CODE HERE:
148     # Implement RMSProp. Store the next value of w as next_w. You need
149     # to also store in config['a'] the moving average of the second
150     # moment gradients, so they can be used for future gradients. Concretely,
151     # config['a'] corresponds to "a" in the lecture notes.
152     # ===== #
153
154     a = config['a']
155     beta = config['decay_rate']
156     eps = config['learning_rate']
157     nu = config['epsilon']
158
159     a = beta * a + (1 - beta) * dw * dw
160     w -= (eps * dw) / (np.sqrt(a) + nu)
161
162     config['a'] = a
163     next_w = w
164
165     # ===== #
166     # END YOUR CODE HERE
167     # ===== #
168
169     return next_w, config

```

```

170
171 def adam(w, dw, config=None):
172     """
173     Uses the Adam update rule, which incorporates moving averages of both the
174     gradient and its square and a bias correction term.
175
176     config format:
177     - learning_rate: Scalar learning rate.
178     - beta1: Decay rate for moving average of first moment of gradient.
179     - beta2: Decay rate for moving average of second moment of gradient.
180     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
181     - m: Moving average of gradient.
182     - v: Moving average of squared gradient.
183     - t: Iteration number.
184     """
185     if config is None: config = {}
186     config.setdefault('learning_rate', 1e-3)
187     config.setdefault('beta1', 0.9)
188     config.setdefault('beta2', 0.999)
189     config.setdefault('epsilon', 1e-8)
190     config.setdefault('v', np.zeros_like(w))
191     config.setdefault('a', np.zeros_like(w))
192     config.setdefault('t', 0)
193
194     next_w = None
195
196     # ===== #
197     # YOUR CODE HERE:
198     # Implement Adam. Store the next value of w as next_w. You need
199     # to also store in config['a'] the moving average of the second
200     # moment gradients, and in config['v'] the moving average of the
201     # first moments. Finally, store in config['t'] the increasing time.
202     # ===== #
203
204     t = config['t']
205     v = config['v']
206     a = config['a']
207     eps = config['learning_rate']
208     nu = config['epsilon']
209     beta1 = config['beta1']
210     beta2 = config['beta2']
211
212     t += 1
213     v = beta1 * v + (1 - beta1) * dw
214     a = beta2 * a + (1 - beta2) * dw * dw
215     v_u = v / (1 - beta1**t)
216     a_u = a / (1 - beta2**t)
217     w -= (eps * v_u) / (np.sqrt(a_u) + nu)
218
219     config['t'] = t
220     config['v'] = v
221     config['a'] = a
222     next_w = w
223
224     # ===== #
225     # END YOUR CODE HERE
226     # ===== #
227
228     return next_w, config
229

```

230  
231  
232  
233  
234



# ECE C147/247 HW4 Q2: Batch Normalization

In this notebook, you will implement the batch normalization layers of a neural network to increase its performance. Please review the details of batch normalization from the lecture notes.

`utils` has built a solid API for building these modular frameworks and training them, and we will use their very well implemented framework as opposed to "reinventing the wheel." This includes using their Solver, various utility functions, and their layer structure. This also includes `nndl.fc_net` , `nndl.layers` , and `nndl.layer_utils` .

```
In [1]: ## Import and setups

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

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

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

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

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

data = get_CIFAR10_data()
for k in data.keys():
    print('{}: {} '.format(k, data[k].shape))

x_train: (49000, 3, 32, 32)
y_train: (49000,)
x_val: (1000, 3, 32, 32)
y_val: (1000,)
x_test: (1000, 3, 32, 32)
y_test: (1000,)
```

## Batchnorm forward pass

Implement the training time batchnorm forward pass, `batchnorm_forward` , in `nndl/layers.py` . After that, test your implementation by running the following cell.

```
In [3]: # Check the training-time forward pass by checking means and variances
# of features both before and after batch normalization
```

```

# Simulate the forward pass for a two-layer network
N, D1, D2, D3 = 200, 50, 60, 3
X = np.random.randn(N, D1)
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)
a = np.maximum(0, X.dot(W1)).dot(W2)

print('Before batch normalization:')
print('  means: ', a.mean(axis=0))
print('  stds: ', a.std(axis=0))

# Means should be close to zero and stds close to one
print('After batch normalization (gamma=1, beta=0)')
a_norm, _ = batchnorm_forward(a, np.ones(D3), np.zeros(D3), {'mode': 'train'})
print('  mean: ', a_norm.mean(axis=0))
print('  std: ', a_norm.std(axis=0))

# Now means should be close to beta and stds close to gamma
gamma = np.asarray([1.0, 2.0, 3.0])
beta = np.asarray([11.0, 12.0, 13.0])
a_norm, _ = batchnorm_forward(a, gamma, beta, {'mode': 'train'})
print('After batch normalization (nontrivial gamma, beta)')
print('  means: ', a_norm.mean(axis=0))
print('  stds: ', a_norm.std(axis=0))

```

```

Before batch normalization:
  means: [-27.52804285  3.23108626 -39.84281267]
  stds:  [32.45802149 34.16563075 41.61402063]
After batch normalization (gamma=1, beta=0)
  mean:  [6.63913369e-16 2.27595720e-17 4.26603197e-16]
  std:   [1. 1. 1.]
After batch normalization (nontrivial gamma, beta)
  means: [11. 12. 13.]
  stds:  [1.          1.99999999 2.99999999]

```

Implement the testing time batchnorm forward pass, `batchnorm_forward`, in `nndl/layers.py`. After that, test your implementation by running the following cell.

In [4]:

```

# Check the test-time forward pass by running the training-time
# forward pass many times to warm up the running averages, and then
# checking the means and variances of activations after a test-time
# forward pass.

N, D1, D2, D3 = 200, 50, 60, 3
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)

bn_param = {'mode': 'train'}
gamma = np.ones(D3)
beta = np.zeros(D3)
for t in np.arange(50):
    X = np.random.randn(N, D1)
    a = np.maximum(0, X.dot(W1)).dot(W2)
    batchnorm_forward(a, gamma, beta, bn_param)
bn_param['mode'] = 'test'
X = np.random.randn(N, D1)
a = np.maximum(0, X.dot(W1)).dot(W2)
a_norm, _ = batchnorm_forward(a, gamma, beta, bn_param)

# Means should be close to zero and stds close to one, but will be
# noisier than training-time forward passes.
print('After batch normalization (test-time):')
print('  means: ', a_norm.mean(axis=0))
print('  stds: ', a_norm.std(axis=0))

```

```
After batch normalization (test-time):
means: [-0.07597434  0.00739562  0.07691946]
stds:  [0.88788487  1.04432885  0.99089243]
```

## Batchnorm backward pass

Implement the backward pass for the batchnorm layer, `batchnorm_backward` in `nndl/layers.py`. Check your implementation by running the following cell.

```
In [5]: # Gradient check batchnorm backward pass

N, D = 4, 5
x = 5 * np.random.randn(N, D) + 12
gamma = np.random.randn(D)
beta = np.random.randn(D)
dout = np.random.randn(N, D)

bn_param = {'mode': 'train'}
fx = lambda x: batchnorm_forward(x, gamma, beta, bn_param)[0]
fg = lambda a: batchnorm_forward(x, gamma, beta, bn_param)[0]
fb = lambda b: batchnorm_forward(x, gamma, beta, bn_param)[0]

dx_num = eval_numerical_gradient_array(fx, x, dout)
da_num = eval_numerical_gradient_array(fg, gamma, dout)
db_num = eval_numerical_gradient_array(fb, beta, dout)

_, cache = batchnorm_forward(x, gamma, beta, bn_param)
dx, dgamma, dbeta = batchnorm_backward(dout, cache)
print('dx error: ', rel_error(dx_num, dx))
print('dgamma error: ', rel_error(da_num, dgamma))
print('dbeta error: ', rel_error(db_num, dbeta))

dx error:  1.3987316119808057e-09
dgamma error:  8.256068052448124e-12
dbeta error:  3.2755634106925108e-12
```

## Implement a fully connected neural network with batchnorm layers

Modify the `FullyConnectedNet()` class in `nndl/fc_net.py` to incorporate batchnorm layers. You will need to modify the class in the following areas:

- (1) The gammas and betas need to be initialized to 1's and 0's respectively in `__init__`.
- (2) The `batchnorm_forward` layer needs to be inserted between each affine and relu layer (except in the output layer) in a forward pass computation in `loss`. You may find it helpful to write an `affine_batchnorm_relu()` layer in `nndl/layer_utils.py` although this is not necessary.
- (3) The `batchnorm_backward` layer has to be appropriately inserted when calculating gradients.

After you have done the appropriate modifications, check your implementation by running the following cell.

Note, while the relative error for W3 should be small, as we backprop gradients more, you may find the relative error increases. Our relative error for W1 is on the order of  $1e-4$ .

In [6]:

```

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,
                              use_batchnorm=True)

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

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

```

```

Running check with reg = 0
Initial loss: 2.191057485443417
W1 relative error: 8.308090754249485e-05
W2 relative error: 1.0217928713160541e-05
W3 relative error: 5.235349587857187e-10
b1 relative error: 1.7763568394002505e-07
b2 relative error: 2.220446049250313e-08
b3 relative error: 1.0585275403998414e-10
beta1 relative error: 7.306782539667682e-08
beta2 relative error: 8.096766141771687e-09
gamma1 relative error: 7.579108849046012e-08
gamma2 relative error: 3.08976181884722e-09

```

```

Running check with reg = 3.14
Initial loss: 6.988257449407779
W1 relative error: 1.8170001423524735e-06
W2 relative error: 1.5313005259248896e-06
W3 relative error: 1.3481894007180487e-08
b1 relative error: 1.7763568394002505e-07
b2 relative error: 8.881784197001252e-08
b3 relative error: 2.573882924097516e-10
beta1 relative error: 2.520345435204669e-09
beta2 relative error: 7.583300904467152e-09
gamma1 relative error: 2.485653322350572e-09
gamma2 relative error: 1.4850195551907493e-08

```

## Training a deep fully connected network with batch normalization.

To see if batchnorm helps, let's train a deep neural network with and without batch normalization.

```

In [7]: # Try training a very deep net with batchnorm
hidden_dims = [100, 100, 100, 100, 100]

num_train = 1000
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 = 2e-2

```

```

bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=True)
model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=False)

bn_solver = Solver(bn_model, small_data,
                    num_epochs=10, batch_size=50,
                    update_rule='adam',
                    optim_config={
                        'learning_rate': 1e-3,
                    },
                    verbose=True, print_every=200)
bn_solver.train()

solver = Solver(model, small_data,
                num_epochs=10, batch_size=50,
                update_rule='adam',
                optim_config={
                    'learning_rate': 1e-3,
                },
                verbose=True, print_every=200)
solver.train()

```

```

(Iteration 1 / 200) loss: 2.325430
(Epoch 0 / 10) train acc: 0.100000; val_acc: 0.099000
(Epoch 1 / 10) train acc: 0.322000; val_acc: 0.270000
(Epoch 2 / 10) train acc: 0.392000; val_acc: 0.294000
(Epoch 3 / 10) train acc: 0.429000; val_acc: 0.273000
(Epoch 4 / 10) train acc: 0.503000; val_acc: 0.295000
(Epoch 5 / 10) train acc: 0.629000; val_acc: 0.330000
(Epoch 6 / 10) train acc: 0.678000; val_acc: 0.325000
(Epoch 7 / 10) train acc: 0.689000; val_acc: 0.332000
(Epoch 8 / 10) train acc: 0.726000; val_acc: 0.308000
(Epoch 9 / 10) train acc: 0.790000; val_acc: 0.329000
(Epoch 10 / 10) train acc: 0.828000; val_acc: 0.332000
(Iteration 1 / 200) loss: 2.301877
(Epoch 0 / 10) train acc: 0.127000; val_acc: 0.140000
(Epoch 1 / 10) train acc: 0.231000; val_acc: 0.221000
(Epoch 2 / 10) train acc: 0.305000; val_acc: 0.278000
(Epoch 3 / 10) train acc: 0.315000; val_acc: 0.263000
(Epoch 4 / 10) train acc: 0.360000; val_acc: 0.260000
(Epoch 5 / 10) train acc: 0.415000; val_acc: 0.292000
(Epoch 6 / 10) train acc: 0.457000; val_acc: 0.298000
(Epoch 7 / 10) train acc: 0.509000; val_acc: 0.306000
(Epoch 8 / 10) train acc: 0.523000; val_acc: 0.289000
(Epoch 9 / 10) train acc: 0.548000; val_acc: 0.313000
(Epoch 10 / 10) train acc: 0.586000; val_acc: 0.321000

```

In [8]:

```

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')

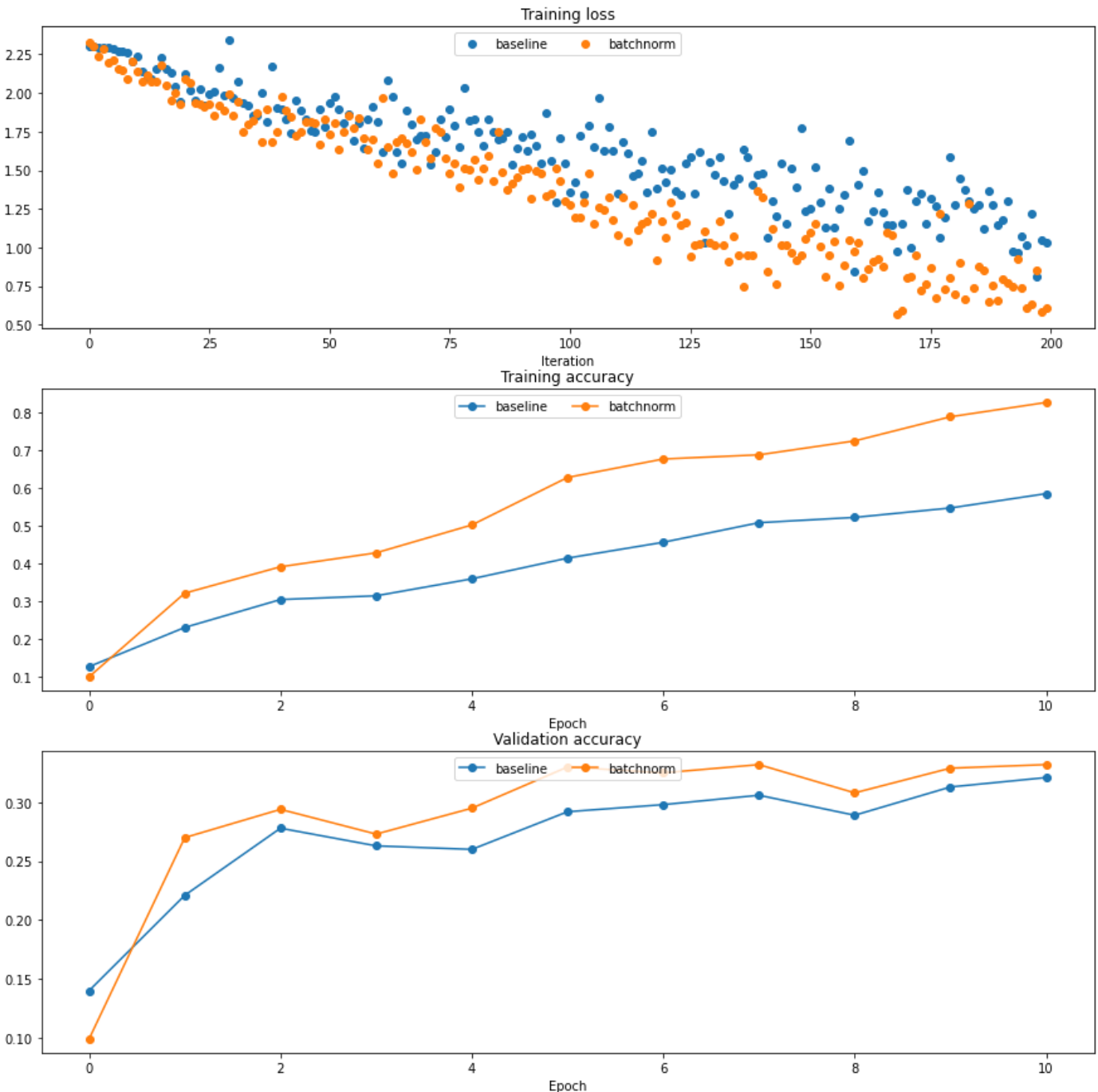
plt.subplot(3, 1, 1)
plt.plot(solver.loss_history, 'o', label='baseline')
plt.plot(bn_solver.loss_history, 'o', label='batchnorm')

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

```

```
plt.subplot(3, 1, 3)
plt.plot(solver.val_acc_history, '-o', label='baseline')
plt.plot(bn_solver.val_acc_history, '-o', label='batchnorm')

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()
```



## Batchnorm and initialization

The following cells run an experiment where for a deep network, the initialization is varied. We do training for when batchnorm layers are and are not included.

```
In [9]: # Try training a very deep net with batchnorm
hidden_dims = [50, 50, 50, 50, 50, 50, 50]
```

```

num_train = 1000
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'],
}

bn_solvers = {}
solvers = {}
weight_scales = np.logspace(-4, 0, num=20)
for i, weight_scale in enumerate(weight_scales):
    print('Running weight scale {} / {}'.format(i + 1, len(weight_scales)))
    bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=True)
    model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=False)

    bn_solver = Solver(bn_model, small_data,
                       num_epochs=10, batch_size=50,
                       update_rule='adam',
                       optim_config={
                           'learning_rate': 1e-3,
                       },
                       verbose=False, print_every=200)
    bn_solver.train()
    bn_solvers[weight_scale] = bn_solver

    solver = Solver(model, small_data,
                    num_epochs=10, batch_size=50,
                    update_rule='adam',
                    optim_config={
                        'learning_rate': 1e-3,
                    },
                    verbose=False, print_every=200)
    solver.train()
    solvers[weight_scale] = solver

```

```

Running weight scale 1 / 20
Running weight scale 2 / 20
Running weight scale 3 / 20
Running weight scale 4 / 20
Running weight scale 5 / 20
Running weight scale 6 / 20
Running weight scale 7 / 20
Running weight scale 8 / 20
Running weight scale 9 / 20
Running weight scale 10 / 20
Running weight scale 11 / 20
Running weight scale 12 / 20
Running weight scale 13 / 20
Running weight scale 14 / 20
Running weight scale 15 / 20
Running weight scale 16 / 20
Running weight scale 17 / 20
Running weight scale 18 / 20
Running weight scale 19 / 20
Running weight scale 20 / 20

```

In [10]:

```

# Plot results of weight scale experiment
best_train_accs, bn_best_train_accs = [], []
best_val_accs, bn_best_val_accs = [], []
final_train_loss, bn_final_train_loss = [], []

for ws in weight_scales:
    best_train_accs.append(max(solvers[ws].train_acc_history))
    bn_best_train_accs.append(max(bn_solvers[ws].train_acc_history))

```

```

best_val_accs.append(max(solvers[ws].val_acc_history))
bn_best_val_accs.append(max(bn_solvers[ws].val_acc_history))

final_train_loss.append(np.mean(solvers[ws].loss_history[-100:]))
bn_final_train_loss.append(np.mean(bn_solvers[ws].loss_history[-100:]))

plt.subplot(3, 1, 1)
plt.title('Best val accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Best val accuracy')
plt.semilogx(weight_scales, best_val_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_val_accs, '-o', label='batchnorm')
plt.legend(ncol=2, loc='lower right')

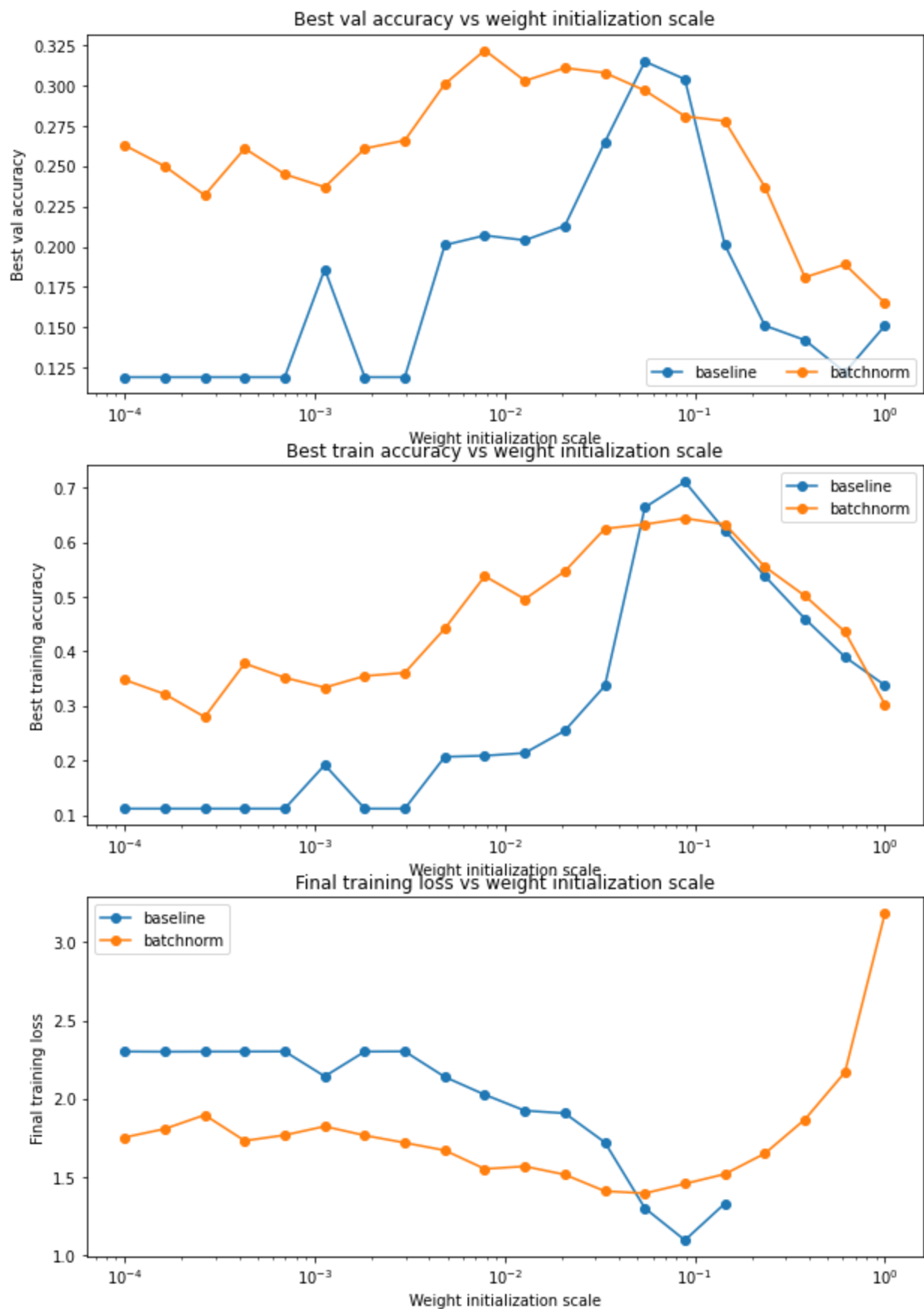
plt.subplot(3, 1, 2)
plt.title('Best train accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Best training accuracy')
plt.semilogx(weight_scales, best_train_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_train_accs, '-o', label='batchnorm')
plt.legend()

plt.subplot(3, 1, 3)
plt.title('Final training loss vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Final training loss')
plt.semilogx(weight_scales, final_train_loss, '-o', label='baseline')
plt.semilogx(weight_scales, bn_final_train_loss, '-o', label='batchnorm')
plt.legend()

plt.gcf().set_size_inches(10, 15)
plt.show()

```





Question:

In the cell below, summarize the findings of this experiment, and WHY these results make sense.

Answer:

- From the plots, we can know that batchnorm can achieve higher accuracy. When the weight initialization scale is smaller than  $10^{-1}$ , the model with batchnorm performs better accuracies and has lower training loss than the model without batchnorm.
- The model with batchnorm (the orange curve) has wider ranges on the weight initialization scale in the training loss plot. And, it is stable and consistent with the change of the weight initialization scale. It means that batchnorm makes the model less sensitive to the weight initialization scale.
- This is because batchnorm is used in regularization, and makes results more stable. Therefore, from the results shown above, the model with batchnorm (the orange curve) is less sensitive and less affected by the weight initialization scale.

# ECE C147/247 HW4 Q3: Dropout

In this notebook, you will implement dropout. Then we will ask you to train a network with batchnorm and dropout, and achieve over 55% accuracy on CIFAR-10.

`utils` has built a solid API for building these modular frameworks and training them, and we will use their very well implemented framework as opposed to "reinventing the wheel." This includes using their Solver, various utility functions, and their layer structure. This also includes `nndl.fc_net`, `nndl.layers`, and `nndl.layer_utils`.

```
In [1]: ## Import and setups

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

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

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

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

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

data = get_CIFAR10_data()
for k in data.keys():
    print('{}: {}'.format(k, data[k].shape))

x_train: (49000, 3, 32, 32)
y_train: (49000,)
x_val: (1000, 3, 32, 32)
y_val: (1000,)
x_test: (1000, 3, 32, 32)
y_test: (1000,)
```

## Dropout forward pass

Implement the training and test time dropout forward pass, `dropout_forward`, in `nndl/layers.py`. After that, test your implementation by running the following cell.

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

```

for p in [0.3, 0.6, 0.75]:
    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())

```

```

Running tests with p = 0.3
Mean of input: 10.001329958949242
Mean of train-time output: 10.000677630103738
Mean of test-time output: 10.001329958949242
Fraction of train-time output set to zero: 0.300204
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.6
Mean of input: 10.001329958949242
Mean of train-time output: 9.995013082981687
Mean of test-time output: 10.001329958949242
Fraction of train-time output set to zero: 0.600344
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.75
Mean of input: 10.001329958949242
Mean of train-time output: 9.99724972751136
Mean of test-time output: 10.001329958949242
Fraction of train-time output set to zero: 0.750256
Fraction of test-time output set to zero: 0.0

```

## Dropout backward pass

Implement the backward pass, `dropout_backward`, in `nndl/layers.py`. After that, test your gradients by running the following cell:

In [4]:

```

x = np.random.randn(10, 10) + 10
dout = np.random.randn(*x.shape)

dropout_param = {'mode': 'train', 'p': 0.8, '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)[0],

print('dx relative error: ', rel_error(dx, dx_num))

```

```
dx relative error: 1.8929063206183813e-11
```

## Implement a fully connected neural network with dropout layers

Modify the `FullyConnectedNet()` class in `nndl/fc_net.py` to incorporate dropout. A dropout layer should be incorporated after every ReLU layer. Concretely, there shouldn't be a dropout at the output layer since there is no ReLU at the output layer. You will need to modify the class in the following areas:

- (1) In the forward pass, you will need to incorporate a dropout layer after every relu layer.
- (2) In the backward pass, you will need to incorporate a dropout backward pass layer.

Check your implementation by running the following code. Our W1 gradient relative error is on the order of  $1e-6$  (the largest of all the relative errors).

In [5]:

```
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 [0, 0.25, 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)

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

```
Running check with dropout = 0
Initial loss: 2.3051948273987857
W1 relative error: 2.5272575344376073e-07
W2 relative error: 1.8875602996866792e-05
W3 relative error: 2.3446001188064074e-07
b1 relative error: 1.3413524910372306e-07
b2 relative error: 5.051339805546953e-08
b3 relative error: 1.4926760614954125e-10
```

```
Running check with dropout = 0.25
Initial loss: 2.29898614757146
W1 relative error: 9.737733598900162e-07
W2 relative error: 5.073657932383196e-08
W3 relative error: 2.8870225716091813e-08
b1 relative error: 9.618747087456625e-09
b2 relative error: 1.897778283870511e-09
b3 relative error: 8.933554101600298e-11
```

```
Running check with dropout = 0.5
Initial loss: 2.302437587710995
W1 relative error: 4.553387957138422e-08
W2 relative error: 1.8475263967389784e-08
W3 relative error: 4.3413247403122424e-07
b1 relative error: 1.872462967441693e-08
b2 relative error: 5.045591219274328e-09
b3 relative error: 1.2715825087959627e-10
```

## Dropout as a regularizer

In class, we claimed that dropout acts as a regularizer by effectively bagging. To check this, we will train two small networks, one with dropout and one without dropout.

In [6]:

```
# Train two identical nets, one with dropout and one without

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 = [0, 0.6]
for dropout in dropout_choices:
    model = FullyConnectedNet([100, 100, 100], dropout=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

```

```

(Iteration 1 / 125) loss: 2.300804
(Epoch 0 / 25) train acc: 0.220000; val_acc: 0.168000
(Epoch 1 / 25) train acc: 0.188000; val_acc: 0.147000
(Epoch 2 / 25) train acc: 0.266000; val_acc: 0.200000
(Epoch 3 / 25) train acc: 0.338000; val_acc: 0.262000
(Epoch 4 / 25) train acc: 0.378000; val_acc: 0.278000
(Epoch 5 / 25) train acc: 0.428000; val_acc: 0.297000
(Epoch 6 / 25) train acc: 0.468000; val_acc: 0.323000
(Epoch 7 / 25) train acc: 0.494000; val_acc: 0.287000
(Epoch 8 / 25) train acc: 0.566000; val_acc: 0.328000
(Epoch 9 / 25) train acc: 0.572000; val_acc: 0.322000
(Epoch 10 / 25) train acc: 0.622000; val_acc: 0.324000
(Epoch 11 / 25) train acc: 0.670000; val_acc: 0.279000
(Epoch 12 / 25) train acc: 0.710000; val_acc: 0.338000
(Epoch 13 / 25) train acc: 0.746000; val_acc: 0.319000
(Epoch 14 / 25) train acc: 0.792000; val_acc: 0.307000
(Epoch 15 / 25) train acc: 0.834000; val_acc: 0.297000
(Epoch 16 / 25) train acc: 0.876000; val_acc: 0.327000
(Epoch 17 / 25) train acc: 0.886000; val_acc: 0.320000
(Epoch 18 / 25) train acc: 0.918000; val_acc: 0.314000
(Epoch 19 / 25) train acc: 0.922000; val_acc: 0.290000
(Epoch 20 / 25) train acc: 0.944000; val_acc: 0.306000
(Iteration 101 / 125) loss: 0.156105
(Epoch 21 / 25) train acc: 0.968000; val_acc: 0.302000
(Epoch 22 / 25) train acc: 0.978000; val_acc: 0.302000
(Epoch 23 / 25) train acc: 0.976000; val_acc: 0.289000
(Epoch 24 / 25) train acc: 0.986000; val_acc: 0.285000
(Epoch 25 / 25) train acc: 0.978000; val_acc: 0.311000
(Iteration 1 / 125) loss: 2.306395
(Epoch 0 / 25) train acc: 0.120000; val_acc: 0.131000
(Epoch 1 / 25) train acc: 0.170000; val_acc: 0.166000
(Epoch 2 / 25) train acc: 0.246000; val_acc: 0.208000
(Epoch 3 / 25) train acc: 0.240000; val_acc: 0.193000
(Epoch 4 / 25) train acc: 0.234000; val_acc: 0.203000
(Epoch 5 / 25) train acc: 0.234000; val_acc: 0.207000
(Epoch 6 / 25) train acc: 0.238000; val_acc: 0.202000
(Epoch 7 / 25) train acc: 0.276000; val_acc: 0.224000
(Epoch 8 / 25) train acc: 0.288000; val_acc: 0.249000
(Epoch 9 / 25) train acc: 0.314000; val_acc: 0.250000
(Epoch 10 / 25) train acc: 0.324000; val_acc: 0.267000
(Epoch 11 / 25) train acc: 0.360000; val_acc: 0.263000
(Epoch 12 / 25) train acc: 0.360000; val_acc: 0.293000
(Epoch 13 / 25) train acc: 0.352000; val_acc: 0.266000
(Epoch 14 / 25) train acc: 0.366000; val_acc: 0.273000
(Epoch 15 / 25) train acc: 0.390000; val_acc: 0.280000
(Epoch 16 / 25) train acc: 0.438000; val_acc: 0.294000
(Epoch 17 / 25) train acc: 0.442000; val_acc: 0.293000
(Epoch 18 / 25) train acc: 0.416000; val_acc: 0.303000
(Epoch 19 / 25) train acc: 0.400000; val_acc: 0.271000
(Epoch 20 / 25) train acc: 0.384000; val_acc: 0.280000
(Iteration 101 / 125) loss: 1.881304
(Epoch 21 / 25) train acc: 0.412000; val_acc: 0.290000
(Epoch 22 / 25) train acc: 0.460000; val_acc: 0.303000
(Epoch 23 / 25) train acc: 0.494000; val_acc: 0.309000
(Epoch 24 / 25) train acc: 0.474000; val_acc: 0.312000
(Epoch 25 / 25) train acc: 0.488000; val_acc: 0.311000

```

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])

```

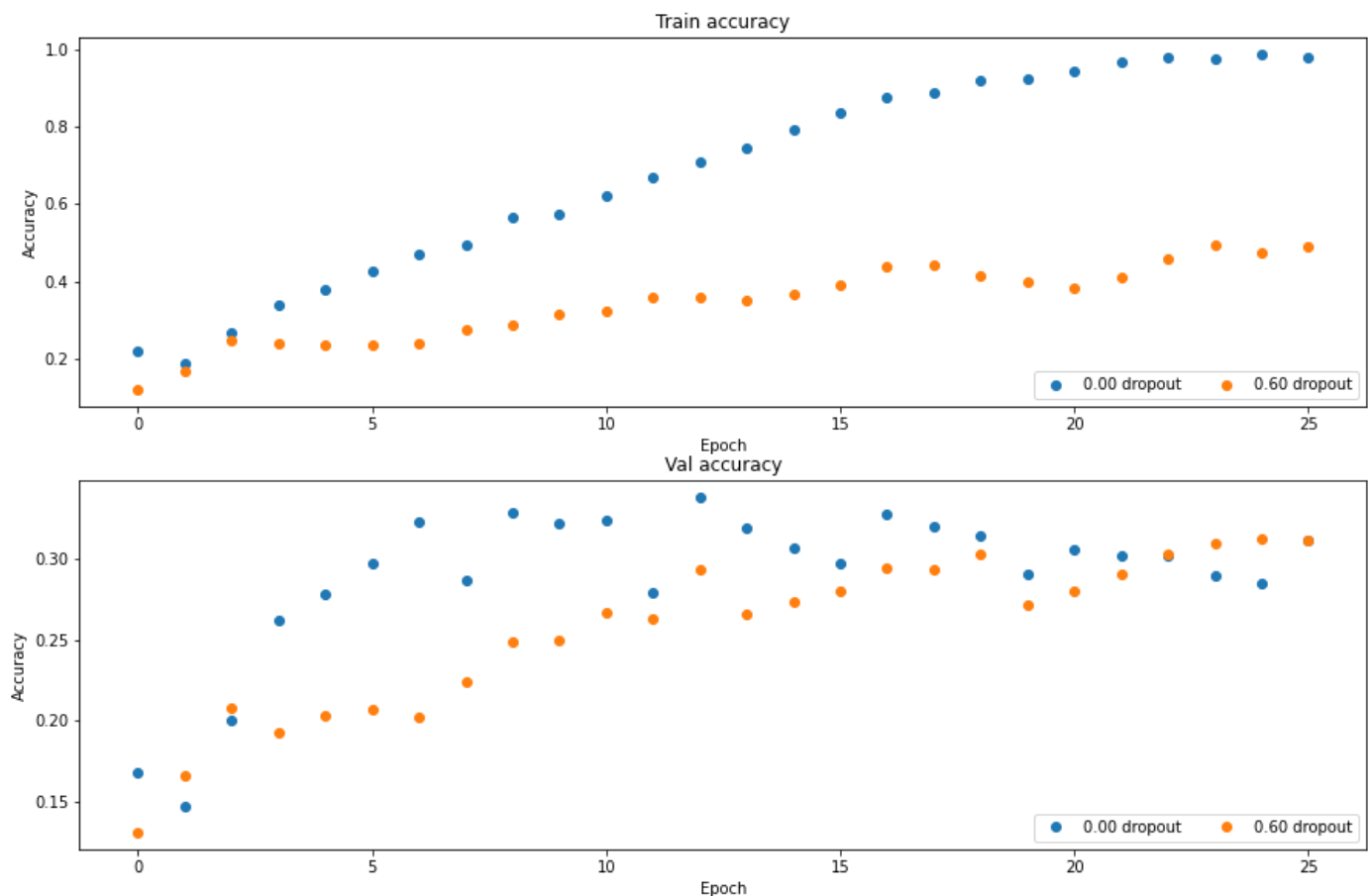
```

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()

```



## Question

Based off the results of this experiment, is dropout performing regularization? Explain your answer.

## Answer:

- Yes, it is performing regularization. From the second plot, we can find that the model with or without dropout have similar validation accuracies. While in the first plot, the model without dropout has significantly higher training accuracies than the model with dropout. It means that the additional training accuracies of the model without dropout is overfitting, and the dropout regularized it.



# Final part of the assignment

Get over 55% validation accuracy on CIFAR-10 by using the layers you have implemented. You will be graded according to the following equation:

$$\min(\text{floor}((X - 32\%) / 28\%, 1)$$

where if you get 60% or higher validation accuracy, you get full points.

In [8]:

```
# ===== #
# YOUR CODE HERE:
#   Implement a FC-net that achieves at least 55% validation accuracy
#   on CIFAR-10.
# ===== #

hidden_dims = [600, 600, 600, 600]
weight_scale = 0.04
dropout = 0.2
learning_rate = 3e-3
lr_decay = 0.95
update_rule = 'adam'

model = FullyConnectedNet(hidden_dims = hidden_dims, weight_scale = weight_scale, dropout
                           use_batchnorm = True, reg = 0.0)

solver = Solver(model, data,
                 num_epochs = 100, batch_size = 500,
                 update_rule = update_rule,
                 optim_config = {
                     'learning_rate': learning_rate,
                 },
                 lr_decay = lr_decay,
                 verbose=True, print_every = 100)

solver.train()

y_test_pred = np.argmax(model.loss(data['X_test']), axis = 1)
y_val_pred = np.argmax(model.loss(data['X_val']), axis = 1)
print("\nTest set accuracy: {}".format(np.mean(np.equal(y_test_pred, data['y_test']))))
print("Validation set accuracy: {}".format(np.mean(np.equal(y_val_pred, data['y_val']))))

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

```
(Iteration 1 / 9800) loss: 2.489383
(Epoch 0 / 100) train acc: 0.186000; val_acc: 0.200000
(Epoch 1 / 100) train acc: 0.436000; val_acc: 0.443000
(Iteration 101 / 9800) loss: 1.597815
(Epoch 2 / 100) train acc: 0.484000; val_acc: 0.498000
(Iteration 201 / 9800) loss: 1.366692
(Epoch 3 / 100) train acc: 0.537000; val_acc: 0.515000
(Iteration 301 / 9800) loss: 1.398289
(Epoch 4 / 100) train acc: 0.577000; val_acc: 0.566000
(Iteration 401 / 9800) loss: 1.248255
(Epoch 5 / 100) train acc: 0.599000; val_acc: 0.538000
(Iteration 501 / 9800) loss: 1.140919
(Epoch 6 / 100) train acc: 0.664000; val_acc: 0.548000
(Iteration 601 / 9800) loss: 1.117880
(Epoch 7 / 100) train acc: 0.613000; val_acc: 0.548000
(Iteration 701 / 9800) loss: 1.022194
(Epoch 8 / 100) train acc: 0.691000; val_acc: 0.580000
(Iteration 801 / 9800) loss: 1.053228
(Epoch 9 / 100) train acc: 0.702000; val_acc: 0.565000
(Iteration 901 / 9800) loss: 0.900552
(Epoch 10 / 100) train acc: 0.715000; val_acc: 0.581000
(Iteration 1001 / 9800) loss: 0.955007
(Epoch 11 / 100) train acc: 0.746000; val_acc: 0.593000
(Iteration 1101 / 9800) loss: 0.873861
(Epoch 12 / 100) train acc: 0.752000; val_acc: 0.579000
(Iteration 1201 / 9800) loss: 0.802197
(Epoch 13 / 100) train acc: 0.757000; val_acc: 0.585000
(Iteration 1301 / 9800) loss: 0.784768
(Epoch 14 / 100) train acc: 0.802000; val_acc: 0.597000
(Iteration 1401 / 9800) loss: 0.775939
(Epoch 15 / 100) train acc: 0.789000; val_acc: 0.583000
(Iteration 1501 / 9800) loss: 0.742718
(Epoch 16 / 100) train acc: 0.838000; val_acc: 0.596000
(Iteration 1601 / 9800) loss: 0.663400
(Epoch 17 / 100) train acc: 0.833000; val_acc: 0.585000
(Iteration 1701 / 9800) loss: 0.644984
(Epoch 18 / 100) train acc: 0.845000; val_acc: 0.597000
(Iteration 1801 / 9800) loss: 0.589513
(Epoch 19 / 100) train acc: 0.838000; val_acc: 0.597000
(Iteration 1901 / 9800) loss: 0.619806
(Epoch 20 / 100) train acc: 0.876000; val_acc: 0.595000
(Iteration 2001 / 9800) loss: 0.576475
(Epoch 21 / 100) train acc: 0.903000; val_acc: 0.586000
(Iteration 2101 / 9800) loss: 0.573623
(Epoch 22 / 100) train acc: 0.899000; val_acc: 0.600000
(Iteration 2201 / 9800) loss: 0.478417
(Epoch 23 / 100) train acc: 0.903000; val_acc: 0.585000
(Iteration 2301 / 9800) loss: 0.470228
(Epoch 24 / 100) train acc: 0.916000; val_acc: 0.599000
(Iteration 2401 / 9800) loss: 0.435533
(Epoch 25 / 100) train acc: 0.932000; val_acc: 0.588000
(Iteration 2501 / 9800) loss: 0.453185
(Epoch 26 / 100) train acc: 0.933000; val_acc: 0.596000
(Iteration 2601 / 9800) loss: 0.473915
(Epoch 27 / 100) train acc: 0.950000; val_acc: 0.603000
(Iteration 2701 / 9800) loss: 0.383685
(Epoch 28 / 100) train acc: 0.944000; val_acc: 0.606000
(Iteration 2801 / 9800) loss: 0.369992
(Epoch 29 / 100) train acc: 0.946000; val_acc: 0.595000
(Iteration 2901 / 9800) loss: 0.412054
(Epoch 30 / 100) train acc: 0.944000; val_acc: 0.594000
(Iteration 3001 / 9800) loss: 0.320266
(Epoch 31 / 100) train acc: 0.969000; val_acc: 0.600000
(Iteration 3101 / 9800) loss: 0.291068
(Epoch 32 / 100) train acc: 0.962000; val_acc: 0.598000
(Iteration 3201 / 9800) loss: 0.294330
```

(Epoch 33 / 100) train acc: 0.963000; val\_acc: 0.608000  
(Iteration 3301 / 9800) loss: 0.261495  
(Epoch 34 / 100) train acc: 0.978000; val\_acc: 0.597000  
(Iteration 3401 / 9800) loss: 0.317429  
(Epoch 35 / 100) train acc: 0.981000; val\_acc: 0.609000  
(Iteration 3501 / 9800) loss: 0.266649  
(Epoch 36 / 100) train acc: 0.974000; val\_acc: 0.598000  
(Iteration 3601 / 9800) loss: 0.309518  
(Epoch 37 / 100) train acc: 0.981000; val\_acc: 0.596000  
(Iteration 3701 / 9800) loss: 0.271904  
(Epoch 38 / 100) train acc: 0.984000; val\_acc: 0.597000  
(Iteration 3801 / 9800) loss: 0.221919  
(Epoch 39 / 100) train acc: 0.980000; val\_acc: 0.610000  
(Iteration 3901 / 9800) loss: 0.247744  
(Epoch 40 / 100) train acc: 0.990000; val\_acc: 0.601000  
(Iteration 4001 / 9800) loss: 0.211057  
(Epoch 41 / 100) train acc: 0.986000; val\_acc: 0.602000  
(Iteration 4101 / 9800) loss: 0.215980  
(Epoch 42 / 100) train acc: 0.989000; val\_acc: 0.602000  
(Iteration 4201 / 9800) loss: 0.189695  
(Epoch 43 / 100) train acc: 0.990000; val\_acc: 0.600000  
(Iteration 4301 / 9800) loss: 0.237263  
(Epoch 44 / 100) train acc: 0.989000; val\_acc: 0.593000  
(Iteration 4401 / 9800) loss: 0.149792  
(Epoch 45 / 100) train acc: 0.981000; val\_acc: 0.603000  
(Iteration 4501 / 9800) loss: 0.195684  
(Epoch 46 / 100) train acc: 0.994000; val\_acc: 0.605000  
(Iteration 4601 / 9800) loss: 0.225804  
(Epoch 47 / 100) train acc: 0.990000; val\_acc: 0.608000  
(Iteration 4701 / 9800) loss: 0.227298  
(Epoch 48 / 100) train acc: 0.993000; val\_acc: 0.604000  
(Iteration 4801 / 9800) loss: 0.155038  
(Epoch 49 / 100) train acc: 0.997000; val\_acc: 0.601000  
(Epoch 50 / 100) train acc: 0.995000; val\_acc: 0.611000  
(Iteration 4901 / 9800) loss: 0.179200  
(Epoch 51 / 100) train acc: 0.994000; val\_acc: 0.603000  
(Iteration 5001 / 9800) loss: 0.148308  
(Epoch 52 / 100) train acc: 0.996000; val\_acc: 0.601000  
(Iteration 5101 / 9800) loss: 0.203314  
(Epoch 53 / 100) train acc: 0.998000; val\_acc: 0.602000  
(Iteration 5201 / 9800) loss: 0.133153  
(Epoch 54 / 100) train acc: 0.998000; val\_acc: 0.613000  
(Iteration 5301 / 9800) loss: 0.179897  
(Epoch 55 / 100) train acc: 0.997000; val\_acc: 0.603000  
(Iteration 5401 / 9800) loss: 0.138111  
(Epoch 56 / 100) train acc: 0.998000; val\_acc: 0.606000  
(Iteration 5501 / 9800) loss: 0.120119  
(Epoch 57 / 100) train acc: 0.999000; val\_acc: 0.607000  
(Iteration 5601 / 9800) loss: 0.167840  
(Epoch 58 / 100) train acc: 0.996000; val\_acc: 0.610000  
(Iteration 5701 / 9800) loss: 0.188360  
(Epoch 59 / 100) train acc: 0.996000; val\_acc: 0.603000  
(Iteration 5801 / 9800) loss: 0.190754  
(Epoch 60 / 100) train acc: 0.999000; val\_acc: 0.612000  
(Iteration 5901 / 9800) loss: 0.189417  
(Epoch 61 / 100) train acc: 0.998000; val\_acc: 0.603000  
(Iteration 6001 / 9800) loss: 0.206282  
(Epoch 62 / 100) train acc: 0.997000; val\_acc: 0.597000  
(Iteration 6101 / 9800) loss: 0.140797  
(Epoch 63 / 100) train acc: 0.997000; val\_acc: 0.604000  
(Iteration 6201 / 9800) loss: 0.153129  
(Epoch 64 / 100) train acc: 0.999000; val\_acc: 0.613000  
(Iteration 6301 / 9800) loss: 0.166696  
(Epoch 65 / 100) train acc: 0.997000; val\_acc: 0.605000  
(Iteration 6401 / 9800) loss: 0.113103  
(Epoch 66 / 100) train acc: 0.999000; val\_acc: 0.604000

(Iteration 6501 / 9800) loss: 0.103705  
(Epoch 67 / 100) train acc: 0.998000; val\_acc: 0.613000  
(Iteration 6601 / 9800) loss: 0.100708  
(Epoch 68 / 100) train acc: 0.999000; val\_acc: 0.607000  
(Iteration 6701 / 9800) loss: 0.108514  
(Epoch 69 / 100) train acc: 0.999000; val\_acc: 0.608000  
(Iteration 6801 / 9800) loss: 0.097242  
(Epoch 70 / 100) train acc: 0.999000; val\_acc: 0.608000  
(Iteration 6901 / 9800) loss: 0.119275  
(Epoch 71 / 100) train acc: 0.999000; val\_acc: 0.606000  
(Iteration 7001 / 9800) loss: 0.092603  
(Epoch 72 / 100) train acc: 0.999000; val\_acc: 0.612000  
(Iteration 7101 / 9800) loss: 0.140257  
(Epoch 73 / 100) train acc: 0.999000; val\_acc: 0.611000  
(Iteration 7201 / 9800) loss: 0.101904  
(Epoch 74 / 100) train acc: 0.999000; val\_acc: 0.614000  
(Iteration 7301 / 9800) loss: 0.109674  
(Epoch 75 / 100) train acc: 0.999000; val\_acc: 0.611000  
(Iteration 7401 / 9800) loss: 0.111727  
(Epoch 76 / 100) train acc: 0.998000; val\_acc: 0.615000  
(Iteration 7501 / 9800) loss: 0.128696  
(Epoch 77 / 100) train acc: 1.000000; val\_acc: 0.613000  
(Iteration 7601 / 9800) loss: 0.145574  
(Epoch 78 / 100) train acc: 0.999000; val\_acc: 0.607000  
(Iteration 7701 / 9800) loss: 0.134654  
(Epoch 79 / 100) train acc: 0.997000; val\_acc: 0.604000  
(Iteration 7801 / 9800) loss: 0.114184  
(Epoch 80 / 100) train acc: 0.999000; val\_acc: 0.612000  
(Iteration 7901 / 9800) loss: 0.173852  
(Epoch 81 / 100) train acc: 0.999000; val\_acc: 0.613000  
(Iteration 8001 / 9800) loss: 0.126730  
(Epoch 82 / 100) train acc: 1.000000; val\_acc: 0.612000  
(Iteration 8101 / 9800) loss: 0.085527  
(Epoch 83 / 100) train acc: 0.997000; val\_acc: 0.613000  
(Iteration 8201 / 9800) loss: 0.090590  
(Epoch 84 / 100) train acc: 0.999000; val\_acc: 0.613000  
(Iteration 8301 / 9800) loss: 0.109501  
(Epoch 85 / 100) train acc: 0.999000; val\_acc: 0.616000  
(Iteration 8401 / 9800) loss: 0.103534  
(Epoch 86 / 100) train acc: 0.999000; val\_acc: 0.613000  
(Iteration 8501 / 9800) loss: 0.096822  
(Epoch 87 / 100) train acc: 1.000000; val\_acc: 0.615000  
(Iteration 8601 / 9800) loss: 0.095079  
(Epoch 88 / 100) train acc: 0.998000; val\_acc: 0.614000  
(Iteration 8701 / 9800) loss: 0.125121  
(Epoch 89 / 100) train acc: 0.999000; val\_acc: 0.611000  
(Iteration 8801 / 9800) loss: 0.127577  
(Epoch 90 / 100) train acc: 0.999000; val\_acc: 0.614000  
(Iteration 8901 / 9800) loss: 0.101643  
(Epoch 91 / 100) train acc: 0.999000; val\_acc: 0.619000  
(Iteration 9001 / 9800) loss: 0.125605  
(Epoch 92 / 100) train acc: 0.999000; val\_acc: 0.611000  
(Iteration 9101 / 9800) loss: 0.133782  
(Epoch 93 / 100) train acc: 0.999000; val\_acc: 0.609000  
(Iteration 9201 / 9800) loss: 0.095340  
(Epoch 94 / 100) train acc: 0.999000; val\_acc: 0.615000  
(Iteration 9301 / 9800) loss: 0.129781  
(Epoch 95 / 100) train acc: 0.997000; val\_acc: 0.614000  
(Iteration 9401 / 9800) loss: 0.121475  
(Epoch 96 / 100) train acc: 0.999000; val\_acc: 0.612000  
(Iteration 9501 / 9800) loss: 0.133388  
(Epoch 97 / 100) train acc: 0.999000; val\_acc: 0.615000  
(Iteration 9601 / 9800) loss: 0.118243  
(Epoch 98 / 100) train acc: 0.998000; val\_acc: 0.615000  
(Iteration 9701 / 9800) loss: 0.100336  
(Epoch 99 / 100) train acc: 1.000000; val\_acc: 0.611000

(Epoch 100 / 100) train acc: 1.000000; val\_acc: 0.613000

Test set accuracy: 0.593

Validation set accuracy: 0.615

```

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

```

```

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

```

```

118 - dx: Gradient with respect to x
119 """
120 x = cache
121
122 # ===== #
123 # YOUR CODE HERE:
124 #   Implement the ReLU backward pass
125 # ===== #
126
127 dx = dout * (x > 0)
128
129 # ===== #
130 # END YOUR CODE HERE
131 # ===== #
132
133 return dx
134
135 def batchnorm_forward(x, gamma, beta, bn_param):
136     """
137     Forward pass for batch normalization.
138
139     During training the sample mean and (uncorrected) sample variance are
140     computed from minibatch statistics and used to normalize the incoming data.
141     During training we also keep an exponentially decaying running mean of the
142     mean
143     and variance of each feature, and these averages are used to normalize data
144     at test-time.
145
146     At each timestep we update the running averages for mean and variance using
147     an exponential decay based on the momentum parameter:
148
149     running_mean = momentum * running_mean + (1 - momentum) * sample_mean
150     running_var = momentum * running_var + (1 - momentum) * sample_var
151
152     Note that the batch normalization paper suggests a different test-time
153     behavior: they compute sample mean and variance for each feature using a
154     large number of training images rather than using a running average. For
155     this implementation we have chosen to use running averages instead since
156     they do not require an additional estimation step; the torch7
157     implementation
158     of batch normalization also uses running averages.
159
160     Input:
161     - x: Data of shape (N, D)
162     - gamma: Scale parameter of shape (D,)
163     - beta: Shift parameter of shape (D,)
164     - bn_param: Dictionary with the following keys:
165       - mode: 'train' or 'test'; required
166       - eps: Constant for numeric stability
167       - momentum: Constant for running mean / variance.
168       - running_mean: Array of shape (D,) giving running mean of features
169       - running_var: Array of shape (D,) giving running variance of features
170
171     Returns a tuple of:
172     - out: of shape (N, D)
173     - cache: A tuple of values needed in the backward pass
174     """
175     mode = bn_param['mode']
176     eps = bn_param.get('eps', 1e-5)
177     momentum = bn_param.get('momentum', 0.9)

```



```

176
177 N, D = x.shape
178 running_mean = bn_param.get('running_mean', np.zeros(D, dtype=x.dtype))
179 running_var = bn_param.get('running_var', np.zeros(D, dtype=x.dtype))
180
181 out, cache = None, None
182 if mode == 'train':
183
184     # ===== #
185     # YOUR CODE HERE:
186     #   A few steps here:
187     #   (1) Calculate the running mean and variance of the minibatch.
188     #   (2) Normalize the activations with the running mean and variance.
189     #   (3) Scale and shift the normalized activations. Store this
190     #       as the variable 'out'
191     #   (4) Store any variables you may need for the backward pass in
192     #       the 'cache' variable.
193     # ===== #
194
195     mean = np.mean(x, axis = 0)
196     var = np.var(x, axis = 0)
197     normalize_x = (x - mean) / np.sqrt(var + eps)
198
199     running_mean = momentum * running_mean + (1 - momentum) * mean
200     running_var = momentum * running_var + (1 - momentum) * var
201
202     out = gamma * normalize_x + beta
203
204     cache = {'normalize_x': normalize_x,
205             'x_minus_mean': (x - mean),
206             'sqrt_var_eps': np.sqrt(var + eps),
207             'gamma': gamma
208             }
209
210     # ===== #
211     # END YOUR CODE HERE
212     # ===== #
213
214 elif mode == 'test':
215
216     # ===== #
217     # YOUR CODE HERE:
218     #   Calculate the testing time normalized activation. Normalize using
219     #   the running mean and variance, and then scale and shift
220     #   appropriately.
221     #   Store the output as 'out'.
222     # ===== #
223
224     normalize_x = (x - running_mean) / np.sqrt(running_var + eps)
225     out = gamma * normalize_x + beta
226
227     # ===== #
228     # END YOUR CODE HERE
229     # ===== #
230
231 else:
232     raise ValueError('Invalid forward batchnorm mode "%s"' % mode)
233
234 # Store the updated running means back into bn_param
235 bn_param['running_mean'] = running_mean

```

```

235     bn_param['running_var'] = running_var
236
237     return out, cache
238
239 def batchnorm_backward(dout, cache):
240     """
241     Backward pass for batch normalization.
242
243     For this implementation, you should write out a computation graph for
244     batch normalization on paper and propagate gradients backward through
245     intermediate nodes.
246
247     Inputs:
248     - dout: Upstream derivatives, of shape (N, D)
249     - cache: Variable of intermediates from batchnorm_forward.
250
251     Returns a tuple of:
252     - dx: Gradient with respect to inputs x, of shape (N, D)
253     - dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
254     - dbeta: Gradient with respect to shift parameter beta, of shape (D,)
255     """
256     dx, dgamma, dbeta = None, None, None
257
258     # ===== #
259     # YOUR CODE HERE:
260     # Implement the batchnorm backward pass, calculating dx, dgamma, and
261     # dbeta.
262     # ===== #
263
264     normalize_x = cache.get('normalize_x')
265     x_minus_mean = cache.get('x_minus_mean')
266     sqrt_var_eps = cache.get('sqrt_var_eps')
267     gamma = cache.get('gamma')
268     N = dout.shape[0]
269
270     dbeta = np.sum(dout, axis = 0)
271     dgamma = np.sum(dout * normalize_x, axis = 0)
272
273     dnormalize_x = dout * gamma
274
275     db = x_minus_mean * dnormalize_x # b = 1 / sqrt_var_eps
276     dc = (-1 / (sqrt_var_eps * sqrt_var_eps)) * db # c = sqrt_var_eps
277     de = (1 / (2 * sqrt_var_eps)) * dc # e = sqrt_var_eps * sqrt_var_eps
278     dvar = np.sum(de, axis = 0)
279
280     da = dnormalize_x / sqrt_var_eps # a = x - mu
281     dmu = -np.sum(da, axis = 0) - 2 * np.sum(x_minus_mean, axis = 0) * dvar / N
282
283     dx = da + 2 * x_minus_mean * dvar / N + dmu / N
284
285     # ===== #
286     # END YOUR CODE HERE
287     # ===== #
288
289     return dx, dgamma, dbeta
290
291 def dropout_forward(x, dropout_param):
292     """
293     Performs the forward pass for (inverted) dropout.

```

```

294 Inputs:
295 - x: Input data, of any shape
296 - dropout_param: A dictionary with the following keys:
297   - p: Dropout parameter. We drop each neuron output with probability p.
298   - mode: 'test' or 'train'. If the mode is train, then perform dropout;
299     if the mode is test, then just return the input.
300   - seed: Seed for the random number generator. Passing seed makes this
301     function deterministic, which is needed for gradient checking but not
in
302     real networks.
303
304 Outputs:
305 - out: Array of the same shape as x.
306 - cache: A tuple (dropout_param, mask). In training mode, mask is the
dropout
307     mask that was used to multiply the input; in test mode, mask is None.
308 """
309 p, mode = dropout_param['p'], dropout_param['mode']
310 if 'seed' in dropout_param:
311     np.random.seed(dropout_param['seed'])
312
313 mask = None
314 out = None
315
316 if mode == 'train':
317     # ===== #
318     # YOUR CODE HERE:
319     # Implement the inverted dropout forward pass during training time.
320     # Store the masked and scaled activations in out, and store the
321     # dropout mask as the variable mask.
322     # ===== #
323
324     mask = (np.random.rand(*x.shape) < (1 - p)) / (1 - p)
325     out = x * mask
326
327     # ===== #
328     # END YOUR CODE HERE
329     # ===== #
330
331 elif mode == 'test':
332
333     # ===== #
334     # YOUR CODE HERE:
335     # Implement the inverted dropout forward pass during test time.
336     # ===== #
337
338     out = x
339
340     # ===== #
341     # END YOUR CODE HERE
342     # ===== #
343
344 cache = (dropout_param, mask)
345 out = out.astype(x.dtype, copy=False)
346
347 return out, cache
348
349 def dropout_backward(dout, cache):
350     """
351     Perform the backward pass for (inverted) dropout.

```

```

352
353 Inputs:
354 - dout: Upstream derivatives, of any shape
355 - cache: (dropout_param, mask) from dropout_forward.
356 """
357 dropout_param, mask = cache
358 mode = dropout_param['mode']
359
360 dx = None
361 if mode == 'train':
362     # ===== #
363     # YOUR CODE HERE:
364     # Implement the inverted dropout backward pass during training time.
365     # ===== #
366
367     dx = dout * mask
368
369     # ===== #
370     # END YOUR CODE HERE
371     # ===== #
372 elif mode == 'test':
373     # ===== #
374     # YOUR CODE HERE:
375     # Implement the inverted dropout backward pass during test time.
376     # ===== #
377
378     dx = dout
379
380     # ===== #
381     # END YOUR CODE HERE
382     # ===== #
383 return dx
384
385 def svm_loss(x, y):
386     """
387     Computes the loss and gradient using for multiclass SVM classification.
388
389     Inputs:
390     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth
class
391     for the ith input.
392     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
393     0 <= y[i] < C
394
395     Returns a tuple of:
396     - loss: Scalar giving the loss
397     - dx: Gradient of the loss with respect to x
398     """
399     N = x.shape[0]
400     correct_class_scores = x[np.arange(N), y]
401     margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
402     margins[np.arange(N), y] = 0
403     loss = np.sum(margins) / N
404     num_pos = np.sum(margins > 0, axis=1)
405     dx = np.zeros_like(x)
406     dx[margins > 0] = 1
407     dx[np.arange(N), y] -= num_pos
408     dx /= N
409     return loss, dx
410

```

```
411
412 def softmax_loss(x, y):
413     """
414     Computes the loss and gradient for softmax classification.
415
416     Inputs:
417     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth
class
418     for the ith input.
419     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
420     0 <= y[i] < C
421
422     Returns a tuple of:
423     - loss: Scalar giving the loss
424     - dx: Gradient of the loss with respect to x
425     """
426
427     probs = np.exp(x - np.max(x, axis=1, keepdims=True))
428     probs /= np.sum(probs, axis=1, keepdims=True)
429     N = x.shape[0]
430     loss = -np.sum(np.log(probs[np.arange(N), y])) / N
431     dx = probs.copy()
432     dx[np.arange(N), y] -= 1
433     dx /= N
434     return loss, dx
435
```

```
1 from .layers import *
2
3
4 def affine_relu_forward(x, w, b):
5     """
6     Convenience layer that performs an affine transform followed by a ReLU
7
8     Inputs:
9     - x: Input to the affine layer
10    - w, b: Weights for the affine layer
11
12    Returns a tuple of:
13    - out: Output from the ReLU
14    - cache: Object to give to the backward pass
15    """
16    a, fc_cache = affine_forward(x, w, b)
17    out, relu_cache = relu_forward(a)
18    cache = (fc_cache, relu_cache)
19    return out, cache
20
21
22 def affine_relu_backward(dout, cache):
23     """
24     Backward pass for the affine-relu convenience layer
25     """
26    fc_cache, relu_cache = cache
27    da = relu_backward(dout, relu_cache)
28    dx, dw, db = affine_backward(da, fc_cache)
29    return dx, dw, db
30
```

```

1 import numpy as np
2 from .layers import *
3 from .layer_utils import *
4
5
6 class TwoLayerNet(object):
7     """
8     A two-layer fully-connected neural network with ReLU nonlinearity and
9     softmax loss that uses a modular layer design. We assume an input dimension
10    of D, a hidden dimension of H, and perform classification over C classes.
11
12    The architecture should be affine - relu - affine - softmax.
13
14    Note that this class does not implement gradient descent; instead, it
15    will interact with a separate Solver object that is responsible for running
16    optimization.
17
18    The learnable parameters of the model are stored in the dictionary
19    self.params that maps parameter names to numpy arrays.
20    """
21
22    def __init__(self, input_dim=3*32*32, hidden_dims=100, num_classes=10,
23                  dropout=0, weight_scale=1e-3, reg=0.0):
24        """
25        Initialize a new network.
26
27        Inputs:
28        - input_dim: An integer giving the size of the input
29        - hidden_dims: An integer giving the size of the hidden layer
30        - num_classes: An integer giving the number of classes to classify
31        - dropout: Scalar between 0 and 1 giving dropout strength.
32        - weight_scale: Scalar giving the standard deviation for random
33          initialization of the weights.
34        - reg: Scalar giving L2 regularization strength.
35        """
36        self.params = {}
37        self.reg = reg
38
39        # ===== #
40        # YOUR CODE HERE:
41        # Initialize W1, W2, b1, and b2. Store these as self.params['W1'],
42        # self.params['W2'], self.params['b1'] and self.params['b2']. The
43        # biases are initialized to zero and the weights are initialized
44        # so that each parameter has mean 0 and standard deviation
45        weight_scale.
46        # The dimensions of W1 should be (input_dim, hidden_dim) and the
47        # dimensions of W2 should be (hidden_dims, num_classes)
48        # ===== #
49        self.params['W1'] = weight_scale * np.random.randn(input_dim,
50        hidden_dims) + 0
51        self.params['W2'] = weight_scale * np.random.randn(hidden_dims,
52        num_classes) + 0
53
54        self.params['b1'] = np.zeros((hidden_dims, 1))
55        self.params['b2'] = np.zeros((num_classes, 1))
56
57        # ===== #
58        # END YOUR CODE HERE

```

```

57     # ===== #
58
59     def loss(self, X, y=None):
60         """
61         Compute loss and gradient for a minibatch of data.
62
63         Inputs:
64         - X: Array of input data of shape (N, d_1, ..., d_k)
65         - y: Array of labels, of shape (N,). y[i] gives the label for X[i].
66
67         Returns:
68         If y is None, then run a test-time forward pass of the model and return:
69         - scores: Array of shape (N, C) giving classification scores, where
70           scores[i, c] is the classification score for X[i] and class c.
71
72         If y is not None, then run a training-time forward and backward pass and
73         return a tuple of:
74         - loss: Scalar value giving the loss
75         - grads: Dictionary with the same keys as self.params, mapping parameter
76           names to gradients of the loss with respect to those parameters.
77         """
78         scores = None
79
80         # ===== #
81         # YOUR CODE HERE:
82         #   Implement the forward pass of the two-layer neural network. Store
83         #   the class scores as the variable 'scores'. Be sure to use the layers
84         #   you prior implemented.
85         # ===== #
86
87         out_l1, cache_l1 = affine_forward(X, self.params['W1'],
self.params['b1'])
88         out_relu, cache_relu = relu_forward(out_l1)
89         scores, cache_l2 = affine_forward(out_relu, self.params['W2'],
self.params['b2'])
90
91         # ===== #
92         # END YOUR CODE HERE
93         # ===== #
94
95         # If y is None then we are in test mode so just return scores
96         if y is None:
97             return scores
98
99         loss, grads = 0, {}
100        # ===== #
101        # YOUR CODE HERE:
102        #   Implement the backward pass of the two-layer neural net. Store
103        #   the loss as the variable 'loss' and store the gradients in the
104        #   'grads' dictionary. For the grads dictionary, grads['W1'] holds
105        #   the gradient for W1, grads['b1'] holds the gradient for b1, etc.
106        #   i.e., grads[k] holds the gradient for self.params[k].
107        #
108        #   Add L2 regularization, where there is an added cost 0.5*self.reg*W^2
109        #   for each W. Be sure to include the 0.5 multiplying factor to
110        #   match our implementation.
111        #
112        #   And be sure to use the layers you prior implemented.
113        # ===== #
114

```



```

115     loss, dx2 = softmax_loss(scores, y)
116     reg_loss = 0.5 * self.reg * (np.linalg.norm(self.params['W1'], 'fro')**2
+ np.linalg.norm(self.params['W2'], 'fro')**2)
117     loss += reg_loss
118
119     dh1, dw2, db2 = affine_backward(dx2, cache_l2)
120     da = relu_backward(dh1, cache_relu)
121     dx1, dw1, db1 = affine_backward(da, cache_l1)
122
123     grads['W1'] = dw1 + self.reg * self.params['W1']
124     grads['b1'] = db1.T
125
126     grads['W2'] = dw2 + self.reg * self.params['W2']
127     grads['b2'] = db2.T
128
129     # ===== #
130     # END YOUR CODE HERE
131     # ===== #
132
133     return loss, grads
134
135
136 class FullyConnectedNet(object):
137     """
138     A fully-connected neural network with an arbitrary number of hidden layers,
139     ReLU nonlinearities, and a softmax loss function. This will also implement
140     dropout and batch normalization as options. For a network with L layers,
141     the architecture will be
142
143     {affine - [batch norm] - relu - [dropout]} x (L - 1) - affine - softmax
144
145     where batch normalization and dropout are optional, and the {...} block is
146     repeated L - 1 times.
147
148     Similar to the TwoLayerNet above, learnable parameters are stored in the
149     self.params dictionary and will be learned using the Solver class.
150     """
151
152     def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,
153                 dropout=0, use_batchnorm=False, reg=0.0,
154                 weight_scale=1e-2, dtype=np.float32, seed=None):
155         """
156         Initialize a new FullyConnectedNet.
157
158         Inputs:
159         - hidden_dims: A list of integers giving the size of each hidden layer.
160         - input_dim: An integer giving the size of the input.
161         - num_classes: An integer giving the number of classes to classify.
162         - dropout: Scalar between 0 and 1 giving dropout strength. If dropout=0
163         then
164             the network should not use dropout at all.
165         - use_batchnorm: Whether or not the network should use batch
166         normalization.
167         - reg: Scalar giving L2 regularization strength.
168         - weight_scale: Scalar giving the standard deviation for random
169         initialization of the weights.
170         - dtype: A numpy datatype object; all computations will be performed
171         using
172             this datatype. float32 is faster but less accurate, so you should use
173             float64 for numeric gradient checking.

```

```

171     - seed: If not None, then pass this random seed to the dropout layers.
This
172     will make the dropout layers deterministic so we can gradient check the
173     model.
174     """
175     self.use_batchnorm = use_batchnorm
176     self.use_dropout = dropout > 0
177     self.reg = reg
178     self.num_layers = 1 + len(hidden_dims)
179     self.dtype = dtype
180     self.params = {}
181
182     # ===== #
183     # YOUR CODE HERE:
184     # Initialize all parameters of the network in the self.params
dictionary.
185     # The weights and biases of layer 1 are W1 and b1; and in general the
186     # weights and biases of layer i are Wi and bi. The
187     # biases are initialized to zero and the weights are initialized
188     # so that each parameter has mean 0 and standard deviation
weight_scale.
189     #
190     # BATCHNORM: Initialize the gammas of each layer to 1 and the beta
191     # parameters to zero. The gamma and beta parameters for layer 1 should
192     # be self.params['gamma1'] and self.params['beta1']. For layer 2, they
193     # should be gamma2 and beta2, etc. Only use batchnorm if
self.use_batchnorm
194     # is true and DO NOT do batch normalize the output scores.
195     # ===== #
196
197     dims = []
198     dims.append(input_dim)
199     dims.extend(hidden_dims)
200     dims.append(num_classes)
201
202     for i in np.arange(self.num_layers):
203         num = str(i+1)
204         self.params['W'+num] = weight_scale * np.random.randn(dims[i],
dims[i+1]) + 0
205         self.params['b'+num] = np.zeros((dims[i+1]))
206
207         if i == (self.num_layers - 1):
208             break
209
210         if self.use_batchnorm:
211             self.params['gamma'+num] = np.ones((dims[i+1]))
212             self.params['beta'+num] = np.zeros((dims[i+1]))
213
214     # ===== #
215     # END YOUR CODE HERE
216     # ===== #
217
218     # When using dropout we need to pass a dropout_param dictionary to each
219     # dropout layer so that the layer knows the dropout probability and the
mode
220     # (train / test). You can pass the same dropout_param to each dropout
layer.
221     self.dropout_param = {}
222     if self.use_dropout:
223         self.dropout_param = {'mode': 'train', 'p': dropout}

```

```

224         if seed is not None:
225             self.dropout_param['seed'] = seed
226
227         # With batch normalization we need to keep track of running means and
228         # variances, so we need to pass a special bn_param object to each batch
229         # normalization layer. You should pass self.bn_params[0] to the forward
pass
230         # of the first batch normalization layer, self.bn_params[1] to the
forward
231         # pass of the second batch normalization layer, etc.
232         self.bn_params = []
233         if self.use_batchnorm:
234             self.bn_params = [{'mode': 'train'} for i in np.arange(self.num_layers
- 1)]
235
236         # Cast all parameters to the correct datatype
237         for k, v in self.params.items():
238             self.params[k] = v.astype(dtype)
239
240
241     def loss(self, X, y=None):
242         """
243         Compute loss and gradient for the fully-connected net.
244
245         Input / output: Same as TwoLayerNet above.
246         """
247         X = X.astype(self.dtype)
248         mode = 'test' if y is None else 'train'
249
250         # Set train/test mode for batchnorm params and dropout param since they
251         # behave differently during training and testing.
252         if self.dropout_param is not None:
253             self.dropout_param['mode'] = mode
254         if self.use_batchnorm:
255             for bn_param in self.bn_params:
256                 bn_param['mode'] = mode
257
258         scores = None
259
260         # ===== #
261         # YOUR CODE HERE:
262         # Implement the forward pass of the FC net and store the output
263         # scores as the variable "scores".
264         #
265         # BATCHNORM: If self.use_batchnorm is true, insert a bathnorm layer
266         # between the affine_forward and relu_forward layers. You may
267         # also write an affine_batchnorm_relu() function in layer_utils.py.
268         #
269         # DROPOUT: If dropout is non-zero, insert a dropout layer after
270         # every ReLU layer.
271         # ===== #
272
273         fc_outs = {}
274         fc_caches = {}
275
276         batchnorm_outs = {}
277         batchnorm_caches = {}
278
279         h = {}
280         h[0] = X

```

```

281     relu_caches = {}
282
283     dropout_caches = {}
284
285     for i in np.arange(self.num_layers):
286         num = str(i+1)
287         # fc
288         fc_outs[i+1], fc_caches[i+1] = affine_forward(h[i],
self.params['W'+num], self.params['b'+num])
289         if i == (self.num_layers - 1):
290             break
291         relu_input = fc_outs[i+1]
292
293         # batch-norm
294         if self.use_batchnorm:
295             batchnorm_outs[i+1], batchnorm_caches[i+1] =
batchnorm_forward(fc_outs[i+1], self.params['gamma'+num],
self.params['beta'+num], self.bn_params[i])
296             relu_input = batchnorm_outs[i+1]
297
298         # relu
299         h[i+1], relu_caches[i+1] = relu_forward(relu_input)
300
301         # dropout
302         if self.use_dropout:
303             h[i+1], dropout_caches[i+1] = dropout_forward(h[i+1],
self.dropout_param)
304
305     scores = fc_outs[self.num_layers]
306
307     # ===== #
308     # END YOUR CODE HERE
309     # ===== #
310
311     # If test mode return early
312     if mode == 'test':
313         return scores
314
315     loss, grads = 0.0, {}
316     # ===== #
317     # YOUR CODE HERE:
318     # Implement the backwards pass of the FC net and store the gradients
319     # in the grads dict, so that grads[k] is the gradient of self.params[k]
320     # Be sure your L2 regularization includes a 0.5 factor.
321     #
322     # BATCHNORM: Incorporate the backward pass of the batchnorm.
323     #
324     # DROPOUT: Incorporate the backward pass of dropout.
325     # ===== #
326
327     loss, dx = softmax_loss(scores, y)
328     reg_loss_sum = 0
329     for i in np.arange(self.num_layers):
330         num = str(i+1)
331         reg_loss_sum += np.linalg.norm(self.params['W'+num], 'fro')**2
332     loss += 0.5 * self.reg * reg_loss_sum
333
334     dict_dW = {}
335     dict_db = {}
336

```

```

337     dict_dgamma = {}
338     dict_dbeta = {}
339
340     # fc - last layer
341     dh, dW, db = affine_backward(dx, fc_caches[self.num_layers])
342     dict_dW[self.num_layers] = dW
343     dict_db[self.num_layers] = db
344
345     for i in np.arange(self.num_layers - 1, 0, -1):
346         # dropout
347         if self.use_dropout:
348             dh = dropout_backward(dh, dropout_caches[i])
349
350         # relu
351         dx_relu = relu_backward(dh, relu_caches[i])
352         fc_input = dx_relu
353
354         # batch-norm
355         if self.use_batchnorm:
356             dx_batchnorm, dgamma, dbeta = batchnorm_backward(dx_relu,
batchnorm_caches[i])
357             dict_dgamma[i] = dgamma
358             dict_dbeta[i] = dbeta
359             fc_input = dx_batchnorm
360
361         # fc
362         dh, dW, db = affine_backward(fc_input, fc_caches[i])
363         dict_dW[i] = dW
364         dict_db[i] = db
365
366     for i in np.arange(self.num_layers):
367         num = str(i+1)
368         grads['W'+num] = dict_dW[i+1] + self.reg * self.params['W'+num]
369         grads['b'+num] = dict_db[i+1]
370         if i == (self.num_layers - 1):
371             break
372         if self.use_batchnorm:
373             grads['gamma'+num] = dict_dgamma[i+1]
374             grads['beta'+num] = dict_dbeta[i+1]
375
376     # ===== #
377     # END YOUR CODE HERE
378     # ===== #
379
380     return loss, grads
381

```