

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.

If you did not complete affine forward and backwards passes, or relu forward and backward passes from HW #3 correctly, you may use another classmate's implementation of these functions for this assignment, or contact us at ece239as.w18@gmail.com.

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

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

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

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

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

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

```
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: 9.7698500479884e-10

If affine_backward is working, error should be less than 1e-9::
dx error: 9.284406992232215e-11
dw error: 1.6325273442881732e-10
db error: 6.057786314879854e-12

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: 3.2756345400232475e-12

If affine_relu_forward and affine_relu_backward are working, error should be less than 1e-9::
dx error: 5.871903409121762e-11
dw error: 1.47644176702134e-10
db error: 3.2755739481167644e-12

Running check with reg = 0
Initial loss: 2.3030445959139376
W1 relative error: 5.031470609833072e-08
W2 relative error: 2.8639211742447507e-07
W3 relative error: 3.8931780465848915e-07
b1 relative error: 1.3580922768177407e-08
b2 relative error: 1.1767511109917437e-09
b3 relative error: 1.735940605963334e-10
Running check with reg = 3.14
Initial loss: 6.76740588610292
W1 relative error: 6.694281324716406e-07
W2 relative error: 1.8351544605456008e-07
W3 relative error: 1.8796349959419661e-07
b1 relative error: 1.2971891978615749e-08
b2 relative error: 8.565763643112226e-08
b3 relative error: 1.420542925364598e-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, which is provided by CS231n, 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)

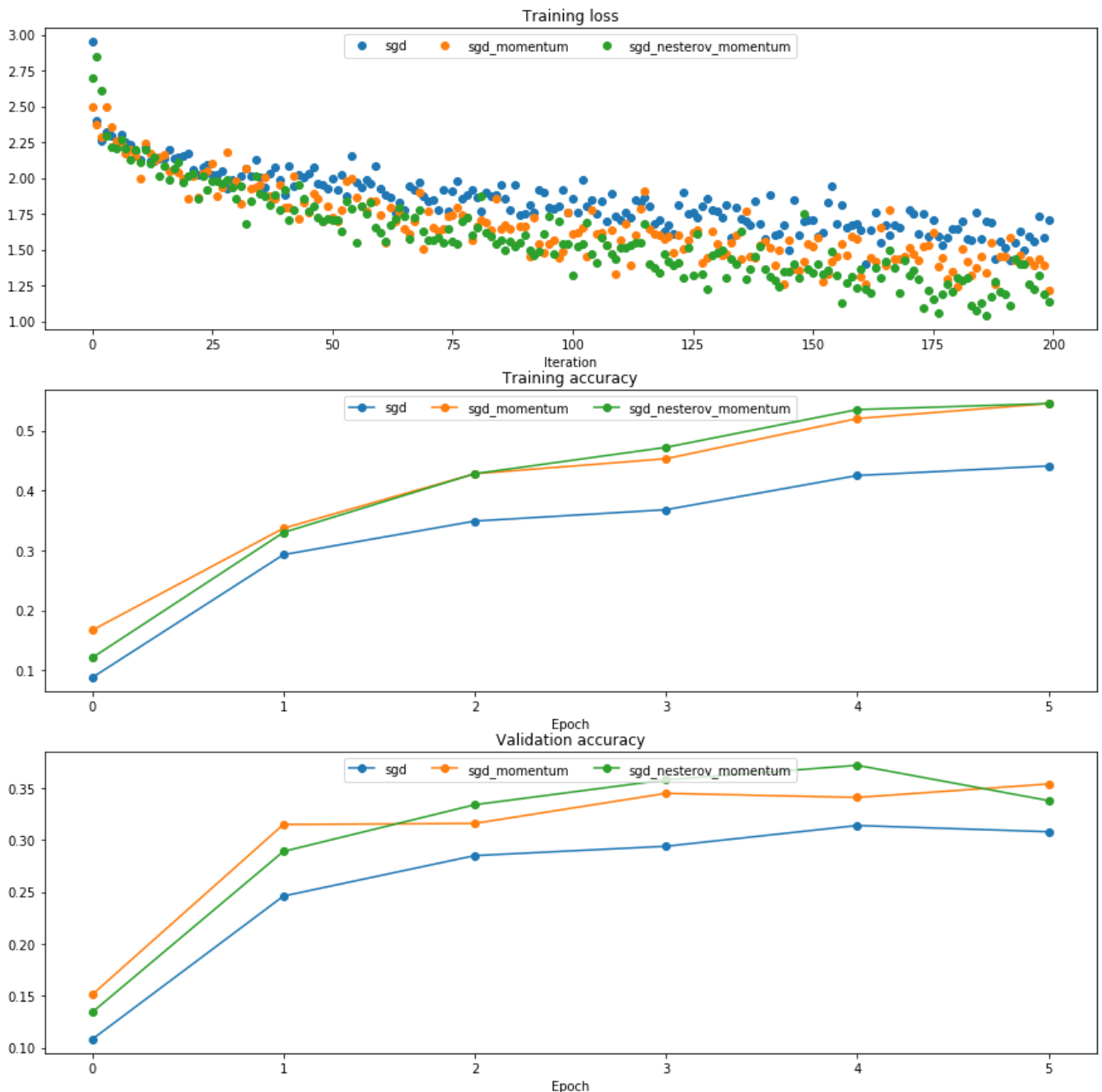
        for i in [1, 2, 3]:
            plt.subplot(3, 1, i)
            plt.legend(loc='upper center', ncol=4)
        plt.gcf().set_size_inches(15, 15)
        plt.show()

```

```
Optimizing with sgd
Optimizing with sgd_momentum
Optimizing with sgd_nesterov_momentum
```

/Users/quentintruong/Desktop/UCLA-ECE239AS-W19/hw4/env/lib/python3.6/site-packages/matplotlib/cbook/deprecation.py:106: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

```
warnings.warn(message, mplDeprecation, stacklevel=1)
```



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)

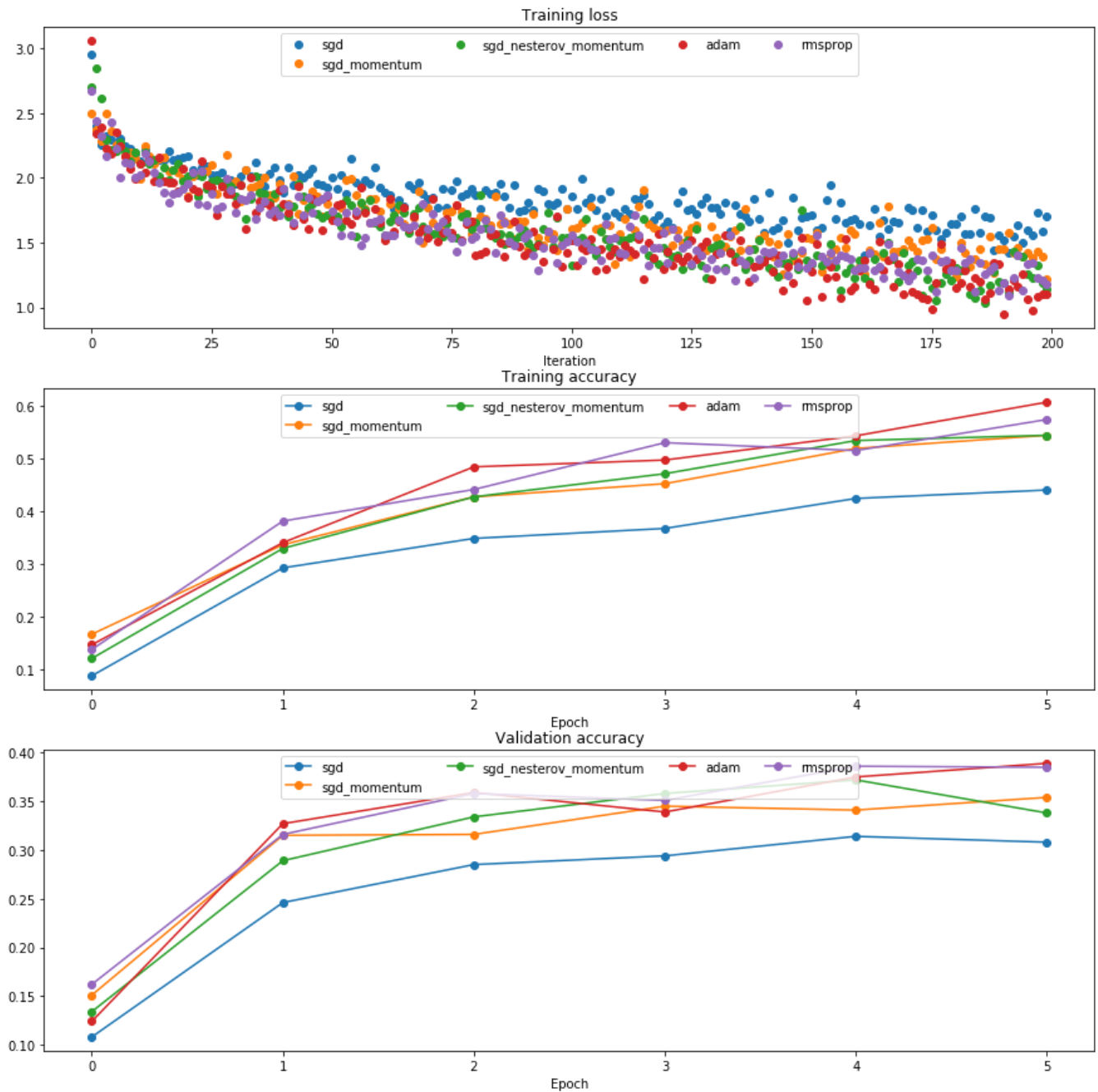
for i in [1, 2, 3]:
    plt.subplot(3, 1, i)
    plt.legend(loc='upper center', ncol=4)
plt.gcf().set_size_inches(15, 15)
plt.show()

```

Optimizing with adam
Optimizing with rmsprop

/Users/quentintruong/Desktop/UCLA-ECE239AS-W19/hw4/env/lib/python3.6/site-packages/matplotlib/cbook/deprecation.py:106: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

```
warnings.warn(message, mplDeprecation, stacklevel=1)
```



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 60+% 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: 2.320879
(Epoch 0 / 10) train acc: 0.257000; val_acc: 0.203000
(Iteration 51 / 4900) loss: 1.855257
(Iteration 101 / 4900) loss: 1.649930
(Iteration 151 / 4900) loss: 1.682393
(Iteration 201 / 4900) loss: 1.579466
(Iteration 251 / 4900) loss: 1.491922
(Iteration 301 / 4900) loss: 1.444749
(Iteration 351 / 4900) loss: 1.718596
(Iteration 401 / 4900) loss: 1.312029
(Iteration 451 / 4900) loss: 1.516429
(Epoch 1 / 10) train acc: 0.505000; val_acc: 0.504000
(Iteration 501 / 4900) loss: 1.490690
(Iteration 551 / 4900) loss: 1.334185
(Iteration 601 / 4900) loss: 1.451640
(Iteration 651 / 4900) loss: 1.406012
(Iteration 701 / 4900) loss: 1.506934
(Iteration 751 / 4900) loss: 1.383689
(Iteration 801 / 4900) loss: 1.341524
(Iteration 851 / 4900) loss: 1.269531
(Iteration 901 / 4900) loss: 1.298273
(Iteration 951 / 4900) loss: 1.108998
(Epoch 2 / 10) train acc: 0.535000; val_acc: 0.514000
(Iteration 1001 / 4900) loss: 1.345785
(Iteration 1051 / 4900) loss: 1.212798
(Iteration 1101 / 4900) loss: 1.149350
(Iteration 1151 / 4900) loss: 1.146366
(Iteration 1201 / 4900) loss: 1.117880
(Iteration 1251 / 4900) loss: 1.033107
(Iteration 1301 / 4900) loss: 1.269095
(Iteration 1351 / 4900) loss: 1.286481
(Iteration 1401 / 4900) loss: 1.056997
(Iteration 1451 / 4900) loss: 1.446774
(Epoch 3 / 10) train acc: 0.598000; val_acc: 0.529000
(Iteration 1501 / 4900) loss: 1.085261
(Iteration 1551 / 4900) loss: 1.038095
(Iteration 1601 / 4900) loss: 1.214963
(Iteration 1651 / 4900) loss: 1.127744
(Iteration 1701 / 4900) loss: 1.211224
(Iteration 1751 / 4900) loss: 1.197089
(Iteration 1801 / 4900) loss: 1.170208
(Iteration 1851 / 4900) loss: 1.170528
(Iteration 1901 / 4900) loss: 0.900901
(Iteration 1951 / 4900) loss: 0.987784
(Epoch 4 / 10) train acc: 0.622000; val_acc: 0.544000
(Iteration 2001 / 4900) loss: 0.864718
(Iteration 2051 / 4900) loss: 0.979695
(Iteration 2101 / 4900) loss: 0.921122
(Iteration 2151 / 4900) loss: 0.882806
(Iteration 2201 / 4900) loss: 0.957249
(Iteration 2251 / 4900) loss: 1.209450
(Iteration 2301 / 4900) loss: 0.889809
(Iteration 2351 / 4900) loss: 1.146533
(Iteration 2401 / 4900) loss: 1.065676
(Epoch 5 / 10) train acc: 0.648000; val_acc: 0.545000
(Iteration 2451 / 4900) loss: 0.933104
(Iteration 2501 / 4900) loss: 0.741917
(Iteration 2551 / 4900) loss: 0.857239
(Iteration 2601 / 4900) loss: 0.925064
(Iteration 2651 / 4900) loss: 1.009138
(Iteration 2701 / 4900) loss: 1.098043

```

(Iteration 2751 / 4900) loss: 0.929387
(Iteration 2801 / 4900) loss: 0.896967
(Iteration 2851 / 4900) loss: 1.177850
(Iteration 2901 / 4900) loss: 0.779549
(Epoch 6 / 10) train acc: 0.701000; val_acc: 0.545000
(Iteration 2951 / 4900) loss: 0.897904
(Iteration 3001 / 4900) loss: 0.943141
(Iteration 3051 / 4900) loss: 0.955763
(Iteration 3101 / 4900) loss: 0.667067
(Iteration 3151 / 4900) loss: 0.916187
(Iteration 3201 / 4900) loss: 0.851208
(Iteration 3251 / 4900) loss: 0.881747
(Iteration 3301 / 4900) loss: 0.716645
(Iteration 3351 / 4900) loss: 0.680211
(Iteration 3401 / 4900) loss: 0.845933
(Epoch 7 / 10) train acc: 0.699000; val_acc: 0.551000
(Iteration 3451 / 4900) loss: 0.682974
(Iteration 3501 / 4900) loss: 0.916752
(Iteration 3551 / 4900) loss: 0.776289
(Iteration 3601 / 4900) loss: 0.692955
(Iteration 3651 / 4900) loss: 0.694823
(Iteration 3701 / 4900) loss: 0.842871
(Iteration 3751 / 4900) loss: 0.589021
(Iteration 3801 / 4900) loss: 0.765869
(Iteration 3851 / 4900) loss: 0.706277
(Iteration 3901 / 4900) loss: 0.808544
(Epoch 8 / 10) train acc: 0.757000; val_acc: 0.543000
(Iteration 3951 / 4900) loss: 0.648728
(Iteration 4001 / 4900) loss: 0.697785
(Iteration 4051 / 4900) loss: 0.628622
(Iteration 4101 / 4900) loss: 0.662943
(Iteration 4151 / 4900) loss: 0.904270
(Iteration 4201 / 4900) loss: 0.408750
(Iteration 4251 / 4900) loss: 0.639366
(Iteration 4301 / 4900) loss: 0.460691
(Iteration 4351 / 4900) loss: 0.681030
(Iteration 4401 / 4900) loss: 0.546384
(Epoch 9 / 10) train acc: 0.799000; val_acc: 0.539000
(Iteration 4451 / 4900) loss: 0.781065
(Iteration 4501 / 4900) loss: 0.587474
(Iteration 4551 / 4900) loss: 0.581917
(Iteration 4601 / 4900) loss: 0.634287
(Iteration 4651 / 4900) loss: 0.470900
(Iteration 4701 / 4900) loss: 0.507806
(Iteration 4751 / 4900) loss: 0.606882
(Iteration 4801 / 4900) loss: 0.498933
(Iteration 4851 / 4900) loss: 0.500319
(Epoch 10 / 10) train acc: 0.830000; val_acc: 0.564000

```

```

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: 0.566
Test set accuracy: 0.575

```

Batch Normalization

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

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

```
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 cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
from cs231n.solver import Solver

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

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

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

```
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: [-32.01112059 -9.94546852 -41.84035895]
  stds: [35.3568188 32.19471298 34.01701357]
After batch normalization (gamma=1, beta=0)
  mean: [-4.20774526e-16  8.71264866e-17 -7.19424520e-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:  [-0.00385673  0.04523241 -0.05831553]
  stds:   [1.10601571  1.08595284  1.00428327]
```

Batchnorm backward pass

Implement the backward pass for the batchnorm layer, `batchnorm_backward` in `nnd1/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:  1.3992698334533834e-08
dgamma error:  5.746786090518256e-12
dbeta error:  6.2021112657457505e-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: 2.188787718661933
W1 relative error: 0.00037497800149038587
W2 relative error: 2.9802196480990068e-05
W3 relative error: 4.888275254235284e-10
b1 relative error: 1.3988810110276972e-06
b2 relative error: 6.661338147750939e-08
b3 relative error: 1.4612611236765422e-10
beta1 relative error: 5.951068872273817e-08
beta2 relative error: 8.59207933341679e-09
gamma1 relative error: 6.559388629007613e-08
gamma2 relative error: 4.915022499153862e-09

```

```

Running check with reg = 3.14
Initial loss: 6.387021516091027
W1 relative error: 2.0159561734290893e-05
W2 relative error: 7.613736541079917e-06
W3 relative error: 4.301795856380999e-08
b1 relative error: 0.004440892098500625
b2 relative error: 0.004440819934004025
b3 relative error: 1.887027549448427e-10
beta1 relative error: 6.551191124783562e-09
beta2 relative error: 6.702536344617225e-09
gamma1 relative error: 6.574369250585113e-09
gamma2 relative error: 7.39453183612395e-09

```

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: 2.362126
(Epoch 0 / 10) train acc: 0.168000; val_acc: 0.136000
(Epoch 1 / 10) train acc: 0.364000; val_acc: 0.281000
(Epoch 2 / 10) train acc: 0.428000; val_acc: 0.308000
(Epoch 3 / 10) train acc: 0.505000; val_acc: 0.333000
(Epoch 4 / 10) train acc: 0.540000; val_acc: 0.321000
(Epoch 5 / 10) train acc: 0.600000; val_acc: 0.312000
(Epoch 6 / 10) train acc: 0.656000; val_acc: 0.329000
(Epoch 7 / 10) train acc: 0.732000; val_acc: 0.331000
(Epoch 8 / 10) train acc: 0.753000; val_acc: 0.327000
(Epoch 9 / 10) train acc: 0.775000; val_acc: 0.335000
(Epoch 10 / 10) train acc: 0.822000; val_acc: 0.340000
(Iteration 1 / 200) loss: 2.302878
(Epoch 0 / 10) train acc: 0.135000; val_acc: 0.138000
(Epoch 1 / 10) train acc: 0.244000; val_acc: 0.219000
(Epoch 2 / 10) train acc: 0.334000; val_acc: 0.257000
(Epoch 3 / 10) train acc: 0.362000; val_acc: 0.287000
(Epoch 4 / 10) train acc: 0.376000; val_acc: 0.278000
(Epoch 5 / 10) train acc: 0.445000; val_acc: 0.315000
(Epoch 6 / 10) train acc: 0.464000; val_acc: 0.299000
(Epoch 7 / 10) train acc: 0.530000; val_acc: 0.298000
(Epoch 8 / 10) train acc: 0.573000; val_acc: 0.288000
(Epoch 9 / 10) train acc: 0.581000; val_acc: 0.283000
(Epoch 10 / 10) train acc: 0.626000; val_acc: 0.328000

```

```
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')

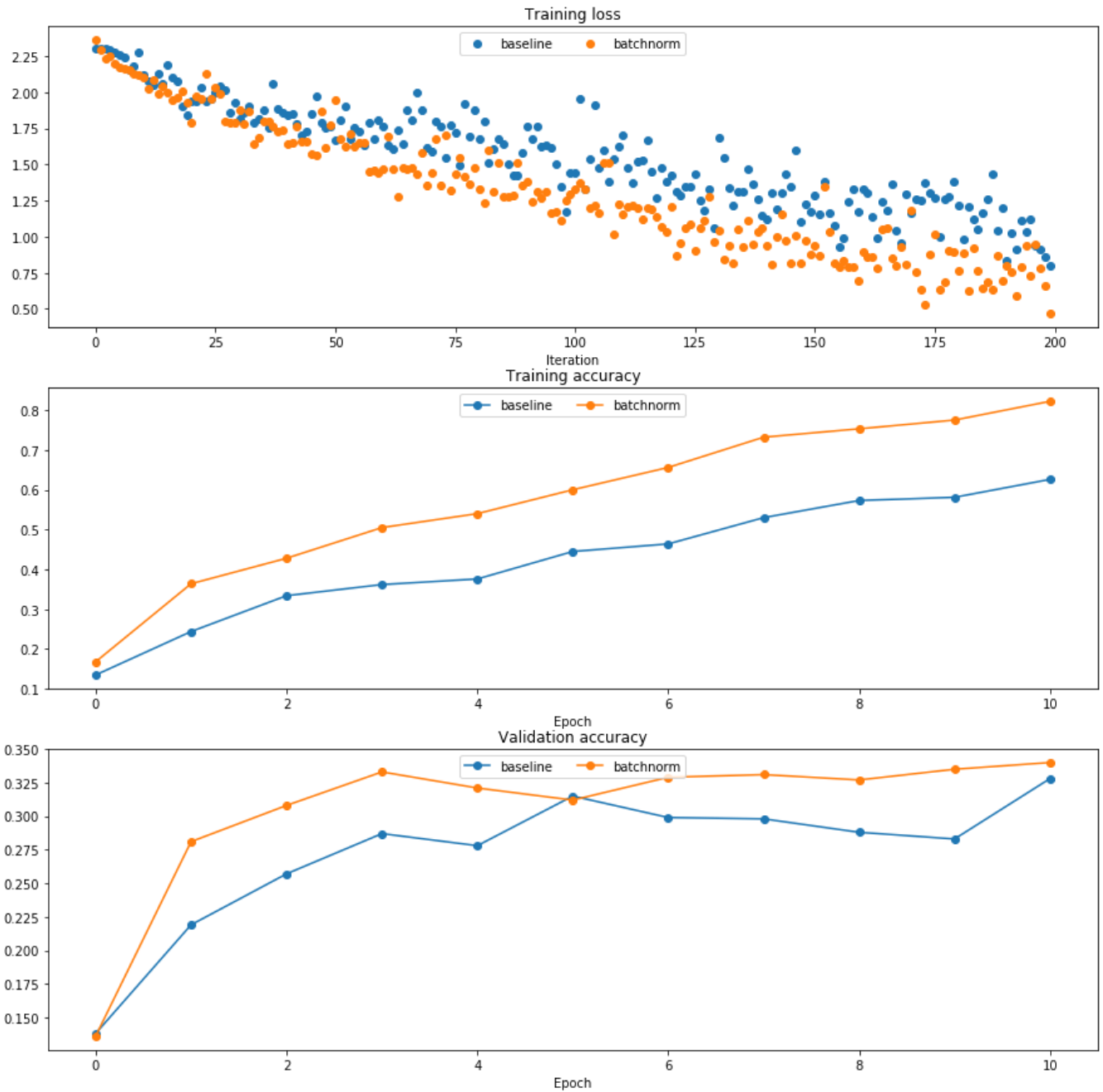
plt.subplot(3, 1, 2)
plt.plot(solver.train_acc_history, '-o', label='baseline')
plt.plot(bn_solver.train_acc_history, '-o', label='batchnorm')

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()
```

/Users/quentintruong/Desktop/UCLA-ECE239AS-W19/hw4/env/lib/python3.6/site-packages/matplotlib/cbook/deprecation.py:106: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

```
warnings.warn(message, mplDeprecation, stacklevel=1)
```



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
```

```
/Users/quentintruong/Desktop/UCLA-ECE239AS-W19/hw4/code/nndl/layers.py:409: RuntimeWarning: divide by zero encountered in log
```

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

```
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:]))

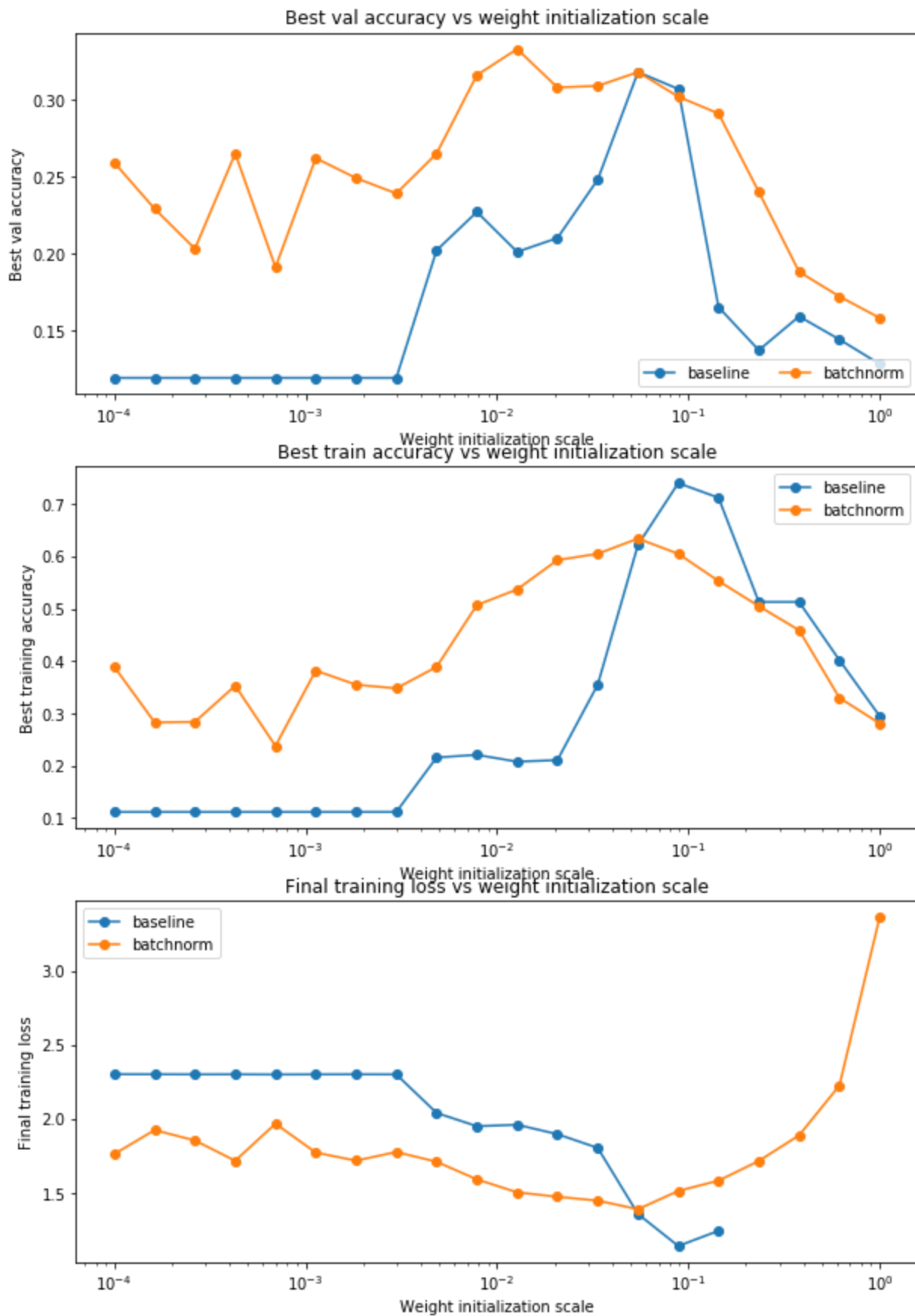
plt.subplot(3, 1, 1)
plt.title('Best val accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Best val accuracy')
plt.semilogx(weight_scales, best_val_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_val_accs, '-o', label='batchnorm')
plt.legend(ncol=2,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

Batchnorm is more capable of achieving higher validation accuracy and training accuracy across a wider range of weight initializations.

First, this makes sense because it prevents the network's weights from going to zero (vanishing); we can see that the network without batchnorm probably died when using weights less than 10^{-3} . Moreover, it prevents the networks weights from going to infinity (exploding); we can see that the network without batchnorm probably exploded when using weights greater than 10^{-1} . Batchnorm helps with weight initializations because it normalizes the data by encouraging zero mean and unit variance, which helps prevent the weights from needing to shrink or grow too much.

Second, batchnorm's higher performance (compared to the network without batchnorm for weights that didn't vanish or explode) is related to the fact that the data is encouraged to be zero mean and unit variance. This unit statistics makes it easier for the network to learn because during gradient descent, the previous layer's data has similar statistics across updates, making the updates more relevant.

Dropout

In this notebook, you will implement dropout. Then we will ask you to train a network with batchnorm and dropout, and achieve over 60% accuracy on CIFAR-10.

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

```
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 cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
from cs231n.solver import Solver

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

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

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

```
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 [15]: 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: 9.999024747460385
Mean of train-time output: 9.973523548843803
Mean of test-time output: 9.999024747460385
Fraction of train-time output set to zero: 0.700644
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.6
Mean of input: 9.999024747460385
Mean of train-time output: 9.980716350624634
Mean of test-time output: 9.999024747460385
Fraction of train-time output set to zero: 0.401164
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.75
Mean of input: 9.999024747460385
Mean of train-time output: 9.99619481359815
Mean of test-time output: 9.999024747460385
Fraction of train-time output set to zero: 0.25018
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 [16]: 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], x, dout)

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

dx relative error: 5.4456103870489836e-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 W_1 gradient relative error is on the order of $1e-6$ (the largest of all the relative errors).

```
In [18]: 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: 2.303043161170242
W1 relative error: 4.795196837428236e-07
W2 relative error: 1.9717710574314515e-07
W3 relative error: 1.5587099943922404e-07
b1 relative error: 2.033615668300116e-08
b2 relative error: 1.6863157228196606e-09
b3 relative error: 1.1144421861081857e-10
```

```
Running check with dropout = 0.25
Initial loss: 2.302354247831908
W1 relative error: 1.0017417771677944e-07
W2 relative error: 2.2591355266316633e-09
W3 relative error: 2.5553916705214888e-05
b1 relative error: 9.368619593647787e-10
b2 relative error: 0.2134290559158092
b3 relative error: 1.2466815697171467e-10
```

```
Running check with dropout = 0.5
Initial loss: 2.304242617164796
W1 relative error: 1.2078835464919935e-07
W2 relative error: 2.454128122918086e-08
W3 relative error: 8.057427967417468e-07
b1 relative error: 2.276255983500775e-08
b2 relative error: 6.836023099402868e-10
b3 relative error: 1.2833211144563337e-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 [19]: *# 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: 2.301328
(Epoch 0 / 25) train acc: 0.154000; val_acc: 0.143000
(Epoch 1 / 25) train acc: 0.214000; val_acc: 0.195000
(Epoch 2 / 25) train acc: 0.252000; val_acc: 0.216000
(Epoch 3 / 25) train acc: 0.276000; val_acc: 0.200000
(Epoch 4 / 25) train acc: 0.308000; val_acc: 0.254000
(Epoch 5 / 25) train acc: 0.316000; val_acc: 0.241000
(Epoch 6 / 25) train acc: 0.322000; val_acc: 0.282000
(Epoch 7 / 25) train acc: 0.354000; val_acc: 0.273000
(Epoch 8 / 25) train acc: 0.364000; val_acc: 0.276000
(Epoch 9 / 25) train acc: 0.408000; val_acc: 0.282000
(Epoch 10 / 25) train acc: 0.454000; val_acc: 0.302000
(Epoch 11 / 25) train acc: 0.472000; val_acc: 0.296000
(Epoch 12 / 25) train acc: 0.496000; val_acc: 0.318000
(Epoch 13 / 25) train acc: 0.512000; val_acc: 0.310000
(Epoch 14 / 25) train acc: 0.532000; val_acc: 0.318000
(Epoch 15 / 25) train acc: 0.558000; val_acc: 0.331000
(Epoch 16 / 25) train acc: 0.574000; val_acc: 0.300000
(Epoch 17 / 25) train acc: 0.624000; val_acc: 0.324000
(Epoch 18 / 25) train acc: 0.608000; val_acc: 0.326000
(Epoch 19 / 25) train acc: 0.622000; val_acc: 0.326000
(Epoch 20 / 25) train acc: 0.664000; val_acc: 0.342000
(Iteration 101 / 125) loss: 1.293077
(Epoch 21 / 25) train acc: 0.688000; val_acc: 0.320000
(Epoch 22 / 25) train acc: 0.708000; val_acc: 0.326000
(Epoch 23 / 25) train acc: 0.734000; val_acc: 0.343000
(Epoch 24 / 25) train acc: 0.756000; val_acc: 0.322000
(Epoch 25 / 25) train acc: 0.788000; val_acc: 0.345000
```

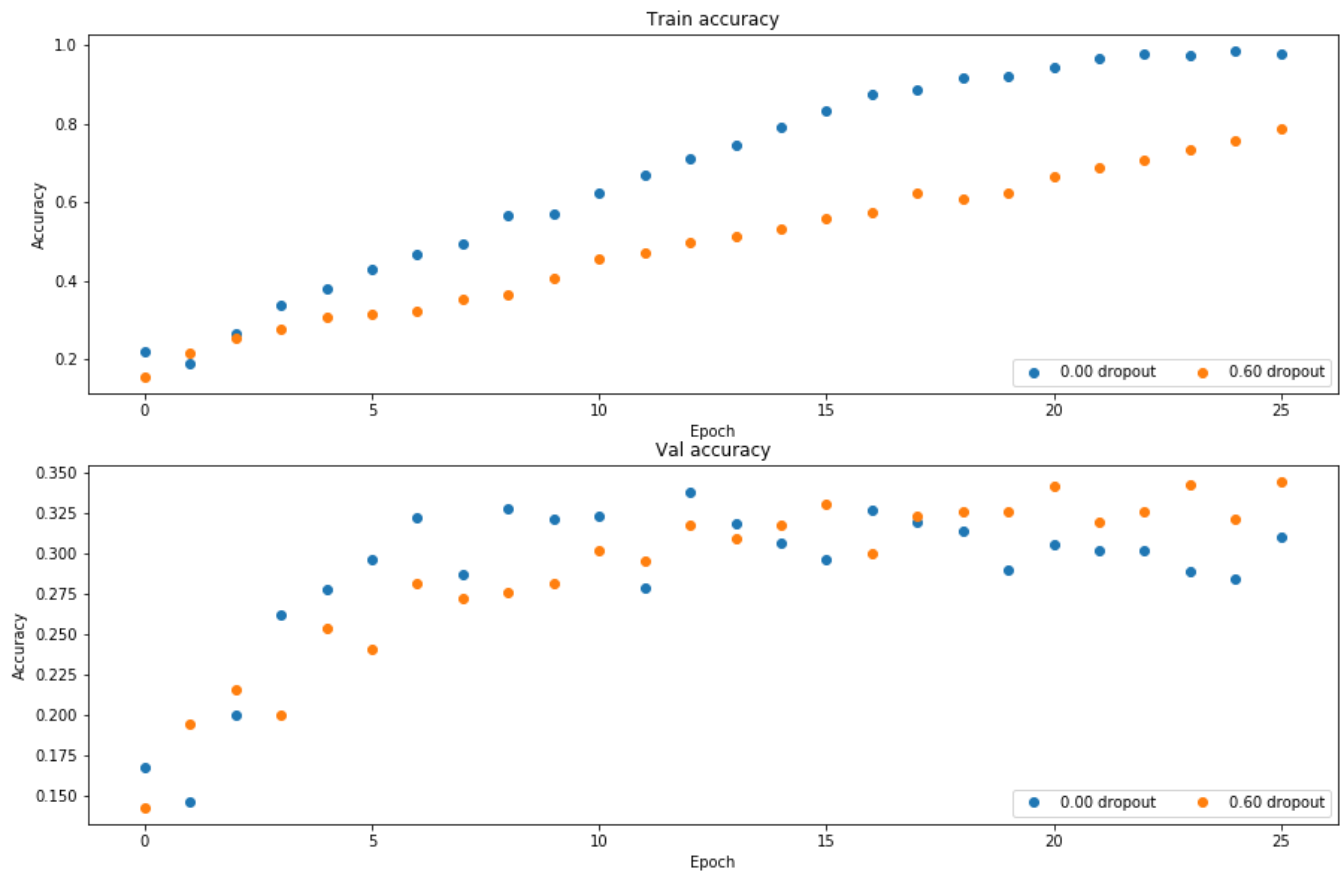
```
In [20]: # Plot train and validation accuracies of the two models
```

```
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, because the difference between the training and validation error is less when using dropout.

Final part of the assignment

Get over 60% 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 [49]: # ===== #
# YOUR CODE HERE:
#   Implement a FC-net that achieves at least 60% validation accuracy
#   on CIFAR-10.
# ===== #
layer_dims = [600, 600, 600, 600]
weight_scale = 0.01
learning_rate = 1e-3
lr_decay = 0.90

model = FullyConnectedNet(layer_dims, weight_scale=weight_scale,
                           dropout=0.65, reg=0.0, use_batchnorm=True)

solver = Solver(model, data,
                 num_epochs=50, batch_size=100,
                 update_rule='adam',
                 optim_config={
                     'learning_rate': learning_rate,
                 },
                 lr_decay=lr_decay,
                 verbose=True, print_every=50)

solver.train()
# ===== #
# END YOUR CODE HERE
# ===== #
```

(Iteration 1 / 24500) loss: 2.305405
(Epoch 0 / 50) train acc: 0.163000; val_acc: 0.155000
(Iteration 51 / 24500) loss: 1.895922
(Iteration 101 / 24500) loss: 1.807972
(Iteration 151 / 24500) loss: 1.760907
(Iteration 201 / 24500) loss: 1.654568
(Iteration 251 / 24500) loss: 1.774943
(Iteration 301 / 24500) loss: 1.716539
(Iteration 351 / 24500) loss: 1.561507
(Iteration 401 / 24500) loss: 1.624929
(Iteration 451 / 24500) loss: 1.647275
(Epoch 1 / 50) train acc: 0.468000; val_acc: 0.463000
(Iteration 501 / 24500) loss: 1.367936
(Iteration 551 / 24500) loss: 1.555204
(Iteration 601 / 24500) loss: 1.483475
(Iteration 651 / 24500) loss: 1.356630
(Iteration 701 / 24500) loss: 1.536303
(Iteration 751 / 24500) loss: 1.720222
(Iteration 801 / 24500) loss: 1.723152
(Iteration 851 / 24500) loss: 1.568987
(Iteration 901 / 24500) loss: 1.266345
(Iteration 951 / 24500) loss: 1.577630
(Epoch 2 / 50) train acc: 0.518000; val_acc: 0.503000
(Iteration 1001 / 24500) loss: 1.331066
(Iteration 1051 / 24500) loss: 1.541900
(Iteration 1101 / 24500) loss: 1.410722
(Iteration 1151 / 24500) loss: 1.367694
(Iteration 1201 / 24500) loss: 1.413612
(Iteration 1251 / 24500) loss: 1.578218
(Iteration 1301 / 24500) loss: 1.257078
(Iteration 1351 / 24500) loss: 1.445430
(Iteration 1401 / 24500) loss: 1.489271
(Iteration 1451 / 24500) loss: 1.424285
(Epoch 3 / 50) train acc: 0.516000; val_acc: 0.509000
(Iteration 1501 / 24500) loss: 1.378358
(Iteration 1551 / 24500) loss: 1.361894
(Iteration 1601 / 24500) loss: 1.469139
(Iteration 1651 / 24500) loss: 1.356522
(Iteration 1701 / 24500) loss: 1.451606
(Iteration 1751 / 24500) loss: 1.440173
(Iteration 1801 / 24500) loss: 1.413091
(Iteration 1851 / 24500) loss: 1.312794
(Iteration 1901 / 24500) loss: 1.540010
(Iteration 1951 / 24500) loss: 1.465006
(Epoch 4 / 50) train acc: 0.563000; val_acc: 0.515000
(Iteration 2001 / 24500) loss: 1.361087
(Iteration 2051 / 24500) loss: 1.133883
(Iteration 2101 / 24500) loss: 1.361684
(Iteration 2151 / 24500) loss: 1.377699
(Iteration 2201 / 24500) loss: 1.345573
(Iteration 2251 / 24500) loss: 1.205908
(Iteration 2301 / 24500) loss: 1.329005
(Iteration 2351 / 24500) loss: 1.517518
(Iteration 2401 / 24500) loss: 1.401187
(Epoch 5 / 50) train acc: 0.579000; val_acc: 0.542000
(Iteration 2451 / 24500) loss: 1.349786
(Iteration 2501 / 24500) loss: 1.179071
(Iteration 2551 / 24500) loss: 1.292837
(Iteration 2601 / 24500) loss: 1.353099
(Iteration 2651 / 24500) loss: 1.202831
(Iteration 2701 / 24500) loss: 1.382101

(Iteration 2751 / 24500) loss: 1.543414
(Iteration 2801 / 24500) loss: 1.254617
(Iteration 2851 / 24500) loss: 1.235397
(Iteration 2901 / 24500) loss: 1.343249
(Epoch 6 / 50) train acc: 0.616000; val_acc: 0.552000
(Iteration 2951 / 24500) loss: 1.313147
(Iteration 3001 / 24500) loss: 1.352385
(Iteration 3051 / 24500) loss: 1.040495
(Iteration 3101 / 24500) loss: 1.312170
(Iteration 3151 / 24500) loss: 1.159872
(Iteration 3201 / 24500) loss: 1.151120
(Iteration 3251 / 24500) loss: 1.246808
(Iteration 3301 / 24500) loss: 1.202060
(Iteration 3351 / 24500) loss: 1.120763
(Iteration 3401 / 24500) loss: 1.390110
(Epoch 7 / 50) train acc: 0.613000; val_acc: 0.535000
(Iteration 3451 / 24500) loss: 1.157195
(Iteration 3501 / 24500) loss: 1.174080
(Iteration 3551 / 24500) loss: 1.411040
(Iteration 3601 / 24500) loss: 1.153144
(Iteration 3651 / 24500) loss: 1.175108
(Iteration 3701 / 24500) loss: 1.235538
(Iteration 3751 / 24500) loss: 1.281613
(Iteration 3801 / 24500) loss: 1.209705
(Iteration 3851 / 24500) loss: 0.884508
(Iteration 3901 / 24500) loss: 1.189045
(Epoch 8 / 50) train acc: 0.634000; val_acc: 0.549000
(Iteration 3951 / 24500) loss: 1.321094
(Iteration 4001 / 24500) loss: 1.150332
(Iteration 4051 / 24500) loss: 1.435022
(Iteration 4101 / 24500) loss: 0.994263
(Iteration 4151 / 24500) loss: 1.110363
(Iteration 4201 / 24500) loss: 1.063714
(Iteration 4251 / 24500) loss: 1.161106
(Iteration 4301 / 24500) loss: 1.366952
(Iteration 4351 / 24500) loss: 1.268125
(Iteration 4401 / 24500) loss: 1.280903
(Epoch 9 / 50) train acc: 0.680000; val_acc: 0.557000
(Iteration 4451 / 24500) loss: 1.130212
(Iteration 4501 / 24500) loss: 1.152019
(Iteration 4551 / 24500) loss: 0.991677
(Iteration 4601 / 24500) loss: 1.159109
(Iteration 4651 / 24500) loss: 1.210210
(Iteration 4701 / 24500) loss: 1.154131
(Iteration 4751 / 24500) loss: 0.984839
(Iteration 4801 / 24500) loss: 0.986202
(Iteration 4851 / 24500) loss: 1.406942
(Epoch 10 / 50) train acc: 0.655000; val_acc: 0.563000
(Iteration 4901 / 24500) loss: 1.340959
(Iteration 4951 / 24500) loss: 0.970882
(Iteration 5001 / 24500) loss: 1.190687
(Iteration 5051 / 24500) loss: 1.196977
(Iteration 5101 / 24500) loss: 1.253798
(Iteration 5151 / 24500) loss: 1.079625
(Iteration 5201 / 24500) loss: 1.251764
(Iteration 5251 / 24500) loss: 1.108914
(Iteration 5301 / 24500) loss: 1.070281
(Iteration 5351 / 24500) loss: 1.207383
(Epoch 11 / 50) train acc: 0.652000; val_acc: 0.569000
(Iteration 5401 / 24500) loss: 1.116703
(Iteration 5451 / 24500) loss: 1.018663
(Iteration 5501 / 24500) loss: 0.968404

(Iteration 5551 / 24500) loss: 1.160531
(Iteration 5601 / 24500) loss: 1.149543
(Iteration 5651 / 24500) loss: 1.108142
(Iteration 5701 / 24500) loss: 1.040542
(Iteration 5751 / 24500) loss: 1.241276
(Iteration 5801 / 24500) loss: 1.050371
(Iteration 5851 / 24500) loss: 0.945394
(Epoch 12 / 50) train acc: 0.678000; val_acc: 0.584000
(Iteration 5901 / 24500) loss: 0.985231
(Iteration 5951 / 24500) loss: 1.058509
(Iteration 6001 / 24500) loss: 1.320623
(Iteration 6051 / 24500) loss: 1.152069
(Iteration 6101 / 24500) loss: 0.863892
(Iteration 6151 / 24500) loss: 0.976334
(Iteration 6201 / 24500) loss: 1.091637
(Iteration 6251 / 24500) loss: 1.045518
(Iteration 6301 / 24500) loss: 1.049506
(Iteration 6351 / 24500) loss: 1.165354
(Epoch 13 / 50) train acc: 0.695000; val_acc: 0.580000
(Iteration 6401 / 24500) loss: 1.087945
(Iteration 6451 / 24500) loss: 1.185851
(Iteration 6501 / 24500) loss: 1.073525
(Iteration 6551 / 24500) loss: 0.814498
(Iteration 6601 / 24500) loss: 1.023440
(Iteration 6651 / 24500) loss: 1.239899
(Iteration 6701 / 24500) loss: 1.166954
(Iteration 6751 / 24500) loss: 0.965370
(Iteration 6801 / 24500) loss: 1.286594
(Iteration 6851 / 24500) loss: 1.113082
(Epoch 14 / 50) train acc: 0.665000; val_acc: 0.580000
(Iteration 6901 / 24500) loss: 0.839494
(Iteration 6951 / 24500) loss: 0.838638
(Iteration 7001 / 24500) loss: 0.895560
(Iteration 7051 / 24500) loss: 0.967985
(Iteration 7101 / 24500) loss: 1.033964
(Iteration 7151 / 24500) loss: 1.073965
(Iteration 7201 / 24500) loss: 1.231775
(Iteration 7251 / 24500) loss: 1.005998
(Iteration 7301 / 24500) loss: 0.933446
(Epoch 15 / 50) train acc: 0.691000; val_acc: 0.590000
(Iteration 7351 / 24500) loss: 0.952628
(Iteration 7401 / 24500) loss: 0.952984
(Iteration 7451 / 24500) loss: 0.958054
(Iteration 7501 / 24500) loss: 1.042820
(Iteration 7551 / 24500) loss: 1.212823
(Iteration 7601 / 24500) loss: 0.915462
(Iteration 7651 / 24500) loss: 1.112218
(Iteration 7701 / 24500) loss: 1.054046
(Iteration 7751 / 24500) loss: 0.901583
(Iteration 7801 / 24500) loss: 1.155146
(Epoch 16 / 50) train acc: 0.730000; val_acc: 0.588000
(Iteration 7851 / 24500) loss: 0.926943
(Iteration 7901 / 24500) loss: 0.750936
(Iteration 7951 / 24500) loss: 0.961800
(Iteration 8001 / 24500) loss: 0.943148
(Iteration 8051 / 24500) loss: 0.823362
(Iteration 8101 / 24500) loss: 0.931198
(Iteration 8151 / 24500) loss: 0.986520
(Iteration 8201 / 24500) loss: 0.839939
(Iteration 8251 / 24500) loss: 1.148176
(Iteration 8301 / 24500) loss: 1.103267
(Epoch 17 / 50) train acc: 0.732000; val_acc: 0.586000

(Iteration 8351 / 24500) loss: 1.013542
(Iteration 8401 / 24500) loss: 1.033530
(Iteration 8451 / 24500) loss: 0.924949
(Iteration 8501 / 24500) loss: 0.995163
(Iteration 8551 / 24500) loss: 1.072750
(Iteration 8601 / 24500) loss: 1.038492
(Iteration 8651 / 24500) loss: 1.219584
(Iteration 8701 / 24500) loss: 1.007752
(Iteration 8751 / 24500) loss: 0.925297
(Iteration 8801 / 24500) loss: 0.940192
(Epoch 18 / 50) train acc: 0.729000; val_acc: 0.576000
(Iteration 8851 / 24500) loss: 1.077766
(Iteration 8901 / 24500) loss: 1.003778
(Iteration 8951 / 24500) loss: 0.784537
(Iteration 9001 / 24500) loss: 0.818034
(Iteration 9051 / 24500) loss: 0.951942
(Iteration 9101 / 24500) loss: 0.891193
(Iteration 9151 / 24500) loss: 1.045544
(Iteration 9201 / 24500) loss: 0.914384
(Iteration 9251 / 24500) loss: 0.907780
(Iteration 9301 / 24500) loss: 0.966635
(Epoch 19 / 50) train acc: 0.756000; val_acc: 0.588000
(Iteration 9351 / 24500) loss: 0.883855
(Iteration 9401 / 24500) loss: 1.157534
(Iteration 9451 / 24500) loss: 0.968482
(Iteration 9501 / 24500) loss: 0.881435
(Iteration 9551 / 24500) loss: 1.023181
(Iteration 9601 / 24500) loss: 0.985334
(Iteration 9651 / 24500) loss: 0.977189
(Iteration 9701 / 24500) loss: 0.989856
(Iteration 9751 / 24500) loss: 1.021420
(Epoch 20 / 50) train acc: 0.749000; val_acc: 0.574000
(Iteration 9801 / 24500) loss: 0.908041
(Iteration 9851 / 24500) loss: 0.898181
(Iteration 9901 / 24500) loss: 0.764440
(Iteration 9951 / 24500) loss: 1.000073
(Iteration 10001 / 24500) loss: 0.744144
(Iteration 10051 / 24500) loss: 0.988009
(Iteration 10101 / 24500) loss: 0.953059
(Iteration 10151 / 24500) loss: 0.912375
(Iteration 10201 / 24500) loss: 1.080402
(Iteration 10251 / 24500) loss: 0.734249
(Epoch 21 / 50) train acc: 0.763000; val_acc: 0.580000
(Iteration 10301 / 24500) loss: 0.922739
(Iteration 10351 / 24500) loss: 0.885038
(Iteration 10401 / 24500) loss: 0.747679
(Iteration 10451 / 24500) loss: 1.126683
(Iteration 10501 / 24500) loss: 0.891979
(Iteration 10551 / 24500) loss: 0.954819
(Iteration 10601 / 24500) loss: 0.940721
(Iteration 10651 / 24500) loss: 0.861106
(Iteration 10701 / 24500) loss: 0.874840
(Iteration 10751 / 24500) loss: 0.909941
(Epoch 22 / 50) train acc: 0.746000; val_acc: 0.592000
(Iteration 10801 / 24500) loss: 0.801810
(Iteration 10851 / 24500) loss: 0.837487
(Iteration 10901 / 24500) loss: 1.056188
(Iteration 10951 / 24500) loss: 0.938667
(Iteration 11001 / 24500) loss: 0.775945
(Iteration 11051 / 24500) loss: 0.898042
(Iteration 11101 / 24500) loss: 0.892971
(Iteration 11151 / 24500) loss: 0.746334

(Iteration 11201 / 24500) loss: 0.842717
(Iteration 11251 / 24500) loss: 0.795465
(Epoch 23 / 50) train acc: 0.781000; val_acc: 0.593000
(Iteration 11301 / 24500) loss: 0.986391
(Iteration 11351 / 24500) loss: 1.040719
(Iteration 11401 / 24500) loss: 0.778256
(Iteration 11451 / 24500) loss: 1.139385
(Iteration 11501 / 24500) loss: 0.729800
(Iteration 11551 / 24500) loss: 0.931683
(Iteration 11601 / 24500) loss: 0.890840
(Iteration 11651 / 24500) loss: 0.945107
(Iteration 11701 / 24500) loss: 0.752788
(Iteration 11751 / 24500) loss: 0.824466
(Epoch 24 / 50) train acc: 0.791000; val_acc: 0.587000
(Iteration 11801 / 24500) loss: 0.912047
(Iteration 11851 / 24500) loss: 0.807942
(Iteration 11901 / 24500) loss: 0.750175
(Iteration 11951 / 24500) loss: 0.814072
(Iteration 12001 / 24500) loss: 0.892791
(Iteration 12051 / 24500) loss: 0.993192
(Iteration 12101 / 24500) loss: 0.911538
(Iteration 12151 / 24500) loss: 0.672663
(Iteration 12201 / 24500) loss: 1.237944
(Epoch 25 / 50) train acc: 0.773000; val_acc: 0.582000
(Iteration 12251 / 24500) loss: 0.773281
(Iteration 12301 / 24500) loss: 1.026120
(Iteration 12351 / 24500) loss: 0.782128
(Iteration 12401 / 24500) loss: 0.784366
(Iteration 12451 / 24500) loss: 0.760001
(Iteration 12501 / 24500) loss: 0.904179
(Iteration 12551 / 24500) loss: 0.781266
(Iteration 12601 / 24500) loss: 0.727542
(Iteration 12651 / 24500) loss: 1.005716
(Iteration 12701 / 24500) loss: 0.925121
(Epoch 26 / 50) train acc: 0.780000; val_acc: 0.591000
(Iteration 12751 / 24500) loss: 0.984700
(Iteration 12801 / 24500) loss: 0.786625
(Iteration 12851 / 24500) loss: 1.019540
(Iteration 12901 / 24500) loss: 0.721131
(Iteration 12951 / 24500) loss: 0.923106
(Iteration 13001 / 24500) loss: 0.984795
(Iteration 13051 / 24500) loss: 0.909641
(Iteration 13101 / 24500) loss: 0.755178
(Iteration 13151 / 24500) loss: 0.700706
(Iteration 13201 / 24500) loss: 0.833847
(Epoch 27 / 50) train acc: 0.764000; val_acc: 0.589000
(Iteration 13251 / 24500) loss: 0.858937
(Iteration 13301 / 24500) loss: 0.856274
(Iteration 13351 / 24500) loss: 0.912349
(Iteration 13401 / 24500) loss: 1.048571
(Iteration 13451 / 24500) loss: 0.820849
(Iteration 13501 / 24500) loss: 1.022572
(Iteration 13551 / 24500) loss: 0.795427
(Iteration 13601 / 24500) loss: 0.985009
(Iteration 13651 / 24500) loss: 0.933283
(Iteration 13701 / 24500) loss: 0.844097
(Epoch 28 / 50) train acc: 0.797000; val_acc: 0.590000
(Iteration 13751 / 24500) loss: 0.905174
(Iteration 13801 / 24500) loss: 0.843156
(Iteration 13851 / 24500) loss: 0.794416
(Iteration 13901 / 24500) loss: 0.819769
(Iteration 13951 / 24500) loss: 0.978043

```
(Iteration 14001 / 24500) loss: 0.965652
(Iteration 14051 / 24500) loss: 0.852074
(Iteration 14101 / 24500) loss: 0.845987
(Iteration 14151 / 24500) loss: 0.747855
(Iteration 14201 / 24500) loss: 0.795871
(Epoch 29 / 50) train acc: 0.792000; val_acc: 0.593000
(Iteration 14251 / 24500) loss: 0.801941
(Iteration 14301 / 24500) loss: 0.884767
(Iteration 14351 / 24500) loss: 0.731136
(Iteration 14401 / 24500) loss: 0.794959
(Iteration 14451 / 24500) loss: 1.021526
(Iteration 14501 / 24500) loss: 0.738700
(Iteration 14551 / 24500) loss: 0.828544
(Iteration 14601 / 24500) loss: 0.759870
(Iteration 14651 / 24500) loss: 0.745240
(Epoch 30 / 50) train acc: 0.767000; val_acc: 0.598000
(Iteration 14701 / 24500) loss: 0.971425
(Iteration 14751 / 24500) loss: 0.844364
(Iteration 14801 / 24500) loss: 0.836075
(Iteration 14851 / 24500) loss: 0.773092
(Iteration 14901 / 24500) loss: 1.104318
(Iteration 14951 / 24500) loss: 0.738051
(Iteration 15001 / 24500) loss: 0.981469
(Iteration 15051 / 24500) loss: 0.948608
(Iteration 15101 / 24500) loss: 1.010112
(Iteration 15151 / 24500) loss: 0.828965
(Epoch 31 / 50) train acc: 0.776000; val_acc: 0.585000
(Iteration 15201 / 24500) loss: 0.884839
(Iteration 15251 / 24500) loss: 0.808028
(Iteration 15301 / 24500) loss: 0.814889
(Iteration 15351 / 24500) loss: 0.842192
(Iteration 15401 / 24500) loss: 0.775567
(Iteration 15451 / 24500) loss: 0.868179
(Iteration 15501 / 24500) loss: 0.888033
(Iteration 15551 / 24500) loss: 0.675218
(Iteration 15601 / 24500) loss: 0.809412
(Iteration 15651 / 24500) loss: 0.795681
(Epoch 32 / 50) train acc: 0.802000; val_acc: 0.592000
(Iteration 15701 / 24500) loss: 0.874625
(Iteration 15751 / 24500) loss: 0.776042
(Iteration 15801 / 24500) loss: 0.905111
(Iteration 15851 / 24500) loss: 0.875138
(Iteration 15901 / 24500) loss: 0.782644
(Iteration 15951 / 24500) loss: 1.008065
(Iteration 16001 / 24500) loss: 1.049957
(Iteration 16051 / 24500) loss: 0.593335
(Iteration 16101 / 24500) loss: 0.959024
(Iteration 16151 / 24500) loss: 1.022140
(Epoch 33 / 50) train acc: 0.803000; val_acc: 0.597000
(Iteration 16201 / 24500) loss: 0.581478
(Iteration 16251 / 24500) loss: 0.930686
(Iteration 16301 / 24500) loss: 0.953515
(Iteration 16351 / 24500) loss: 0.870605
(Iteration 16401 / 24500) loss: 0.783915
(Iteration 16451 / 24500) loss: 0.683314
(Iteration 16501 / 24500) loss: 0.817522
(Iteration 16551 / 24500) loss: 1.043455
(Iteration 16601 / 24500) loss: 0.736522
(Iteration 16651 / 24500) loss: 0.765296
(Epoch 34 / 50) train acc: 0.805000; val_acc: 0.591000
(Iteration 16701 / 24500) loss: 0.803287
(Iteration 16751 / 24500) loss: 1.021465
```

```

(Iteration 16801 / 24500) loss: 0.857296
(Iteration 16851 / 24500) loss: 0.862635
(Iteration 16901 / 24500) loss: 0.894754
(Iteration 16951 / 24500) loss: 0.864463
(Iteration 17001 / 24500) loss: 0.984065
(Iteration 17051 / 24500) loss: 0.729083
(Iteration 17101 / 24500) loss: 0.834994
(Epoch 35 / 50) train acc: 0.810000; val_acc: 0.601000
(Iteration 17151 / 24500) loss: 0.862277

```

```

-----
KeyboardInterrupt                                Traceback (most recent call last)
<ipython-input-49-da0f2753ada5> in <module>()
    20                 lr_decay=lr_decay,
    21                 verbose=True, print_every=50)
--> 22 solver.train()
    23 # ===== #
    24 # END YOUR CODE HERE

~/Desktop/UCLA-ECE239AS-W19/hw4/code/cs231n/solver.py in train(self)
    262
    263         for t in range(num_iterations):
--> 264             self._step()
    265
    266             # Maybe print training loss

~/Desktop/UCLA-ECE239AS-W19/hw4/code/cs231n/solver.py in _step(self)
    185             dw = grads[p]
    186             config = self.optim_configs[p]
--> 187             next_w, next_config = self.update_rule(w, dw, config)
    188             self.model.params[p] = next_w
    189             self.optim_configs[p] = next_config

~/Desktop/UCLA-ECE239AS-W19/hw4/code/nndl/optim.py in adam(w, dw, config)
    190     v_tilda = config['v'] / (1-config['beta1'] ** config['t'])
    191     a_tilda = config['a'] / (1-config['beta2'] ** config['t'])
--> 192     next_w = w - config['learning_rate'] / (np.sqrt(a_tilda) + config['eps
ilon']) * v_tilda
    193     # ===== #
    194     # END YOUR CODE HERE

```

KeyboardInterrupt:


```

In [ ]: import numpy as np
import pdb

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

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

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

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

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

    # ===== #
    # YOUR CODE HERE:
    # Calculate the output of the forward pass. Notice the dimensions
    # of w are D x M, which is the transpose of what we did in earlier
    # assignments.
    # ===== #
    out = np.matmul(x.reshape(x.shape[0], -1), w) + b
    # ===== #
    # END YOUR CODE HERE
    # ===== #

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

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

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

    Returns a tuple of:
    - dx: Gradient with respect to x, of shape (N, d_1, ..., d_k)
    - dw: Gradient with respect to w, of shape (D, M)
    - db: Gradient with respect to b, of shape (M,)
    """

```

```

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

# ===== #
# YOUR CODE HERE:
#   Calculate the gradients for the backward pass.
# ===== #
dx = np.matmul(dout, w.T).reshape(x.shape) # gradient of loss wrt x =  $W^T * upstream$ , reshaped back to x's shape
dw = np.matmul(x.reshape(x.shape[0], -1).T, dout) # gradient of loss wrt w =  $upstream * x^T$ 
db = np.sum(dout, axis=0) # gradient of loss wrt b = upstream, summing wrt data points (so it is per class, aka is dimension M)
# ===== #
# END YOUR CODE HERE
# ===== #

return dx, dw, db

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

    Input:
    - x: Inputs, of any shape

    Returns a tuple of:
    - out: Output, of the same shape as x
    - cache: x
    """
    # ===== #
    # YOUR CODE HERE:
    #   Implement the ReLU forward pass.
    # ===== #
    out = np.maximum(0, x)
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    cache = x
    return out, cache

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

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

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

    # ===== #
    # YOUR CODE HERE:
    #   Implement the ReLU backward pass
    # ===== #
    dx = (x > 0) * dout

```

```

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

return dx

def batchnorm_forward(x, gamma, beta, bn_param):
    """
    Forward pass for batch normalization.

    During training the sample mean and (uncorrected) sample variance are
    computed from minibatch statistics and used to normalize the incoming data.
    During training we also keep an exponentially decaying running mean of the mean
    and variance of each feature, and these averages are used to normalize data
    at test-time.

    At each timestep we update the running averages for mean and variance using
    an exponential decay based on the momentum parameter:

    running_mean = momentum * running_mean + (1 - momentum) * sample_mean
    running_var = momentum * running_var + (1 - momentum) * sample_var

    Note that the batch normalization paper suggests a different test-time
    behavior: they compute sample mean and variance for each feature using a
    large number of training images rather than using a running average. For
    this implementation we have chosen to use running averages instead since
    they do not require an additional estimation step; the torch7 implementation
    of batch normalization also uses running averages.

    Input:
    - x: Data of shape (N, D)
    - gamma: Scale parameter of shape (D,)
    - beta: Shift paremeter of shape (D,)
    - bn_param: Dictionary with the following keys:
        - mode: 'train' or 'test'; required
        - eps: Constant for numeric stability
        - momentum: Constant for running mean / variance.
        - running_mean: Array of shape (D,) giving running mean of features
        - running_var Array of shape (D,) giving running variance of features

    Returns a tuple of:
    - out: of shape (N, D)
    - cache: A tuple of values needed in the backward pass
    """
    mode = bn_param['mode']
    eps = bn_param.get('eps', 1e-5)
    momentum = bn_param.get('momentum', 0.9)

    N, D = x.shape
    running_mean = bn_param.get('running_mean', np.zeros(D, dtype=x.dtype))
    running_var = bn_param.get('running_var', np.zeros(D, dtype=x.dtype))

    out, cache = None, None
    if mode == 'train':

        # ===== #
        # YOUR CODE HERE:
        #   A few steps here:
        #   (1) Calculate the running mean and variance of the minibatch.
        #   (2) Normalize the activations with the batch mean and variance.
        #   (3) Scale and shift the normalized activations. Store this
        #       as the variable 'out'

```

```

#         (4) Store any variables you may need for the backward pass in
#         the 'cache' variable.
# ===== #
mean = np.mean(x, axis=0)
variance = np.var(x, axis=0)
x_hat = (x - mean) / np.sqrt(variance + eps)
out = gamma * x_hat + beta

#print(momentum.shape)

#print(running_mean.shape)
#print(mean.shape)
#print()
running_mean = momentum * running_mean + (1 - momentum) * mean
running_var = momentum * running_var + (1 - momentum) * variance
cache = (x_hat, x, mean, variance, eps, gamma, beta)
# ===== #
# END YOUR CODE HERE
# ===== #

elif mode == 'test':

    # ===== #
    # YOUR CODE HERE:
    # Calculate the testing time normalized activations. Normalize using
    # the running mean and variance, and then scale and shift appropriately.
    # Store the output as 'out'.
    # ===== #
    out = gamma * (x - running_mean) / np.sqrt(running_var + eps) + beta
    # ===== #
    # END YOUR CODE HERE
    # ===== #

else:
    raise ValueError('Invalid forward batchnorm mode "%s"' % mode)

# Store the updated running means back into bn_param
bn_param['running_mean'] = running_mean
bn_param['running_var'] = running_var

return out, cache

def batchnorm_backward(dout, cache):
    """
    Backward pass for batch normalization.

    For this implementation, you should write out a computation graph for
    batch normalization on paper and propagate gradients backward through
    intermediate nodes.

    Inputs:
    - dout: Upstream derivatives, of shape (N, D)
    - cache: Variable of intermediates from batchnorm_forward.

    Returns a tuple of:
    - dx: Gradient with respect to inputs x, of shape (N, D)
    - dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
    - dbeta: Gradient with respect to shift parameter beta, of shape (D,)
    """
    dx, dgamma, dbeta = None, None, None

    # ===== #

```



```

# YOUR CODE HERE:
# Implement the batchnorm backward pass, calculating dx, dgamma, and dbeta.
# ===== #
x_hat, x, mean, variance, eps, gamma, beta = cache

dbeta = np.sum(dout, axis=0)

dgamma = np.sum(x_hat * dout, axis=0)

da = dout
db = gamma * dout
dc = 1 / np.sqrt(variance + eps) * db
dx_1 = dc

dd = (x - mean) * db
de = -1 / (variance + eps) * dd
df = 1 / 2 * (variance + eps) ** (-1/2) * de
dg = df
dh = 2 / x.shape[0] * (x - mean) * np.sum(dg, axis=0)
dx_2 = dh

di = - 1 / np.sqrt(variance + eps) * np.sum(db, axis=0) - np.sum((-1 / 2) * 1 *
(variance + eps) ** -1.5 * (x - mean) * db, axis=0) * 2 / x.shape[0] * np.sum(x -
mean)
dx_3 = 1 / x.shape[0] * di

dx = dx_1 + dx_2 + dx_3
# ===== #
# END YOUR CODE HERE
# ===== #

return dx, dgamma, dbeta

def dropout_forward(x, dropout_param):
    """
    Performs the forward pass for (inverted) dropout.

    Inputs:
    - x: Input data, of any shape
    - dropout_param: A dictionary with the following keys:
      - p: Dropout parameter. We drop each neuron output with probability p.
      - mode: 'test' or 'train'. If the mode is train, then perform dropout;
        if the mode is test, then just return the input.
      - seed: Seed for the random number generator. Passing seed makes this
        function deterministic, which is needed for gradient checking but not in
        real networks.

    Outputs:
    - out: Array of the same shape as x.
    - cache: A tuple (dropout_param, mask). In training mode, mask is the dropout
      mask that was used to multiply the input; in test mode, mask is None.
    """
    p, mode = dropout_param['p'], dropout_param['mode']
    if 'seed' in dropout_param:
        np.random.seed(dropout_param['seed'])

    mask = None
    out = None

    if mode == 'train':
        # ===== #
        # YOUR CODE HERE:

```

```

# Implement the inverted dropout forward pass during training time.
# Store the masked and scaled activations in out, and store the
# dropout mask as the variable mask.
# ===== #
mask = (np.random.rand(*x.shape) < p) / p
out = x * mask
# ===== #
# END YOUR CODE HERE
# ===== #

elif mode == 'test':

# ===== #
# YOUR CODE HERE:
# Implement the inverted dropout forward pass during test time.
# ===== #
out = x
# ===== #
# END YOUR CODE HERE
# ===== #

cache = (dropout_param, mask)
out = out.astype(x.dtype, copy=False)

return out, cache

def dropout_backward(dout, cache):
    """
    Perform the backward pass for (inverted) dropout.

    Inputs:
    - dout: Upstream derivatives, of any shape
    - cache: (dropout_param, mask) from dropout_forward.
    """
    dropout_param, mask = cache
    mode = dropout_param['mode']

    dx = None
    if mode == 'train':
# ===== #
# YOUR CODE HERE:
# Implement the inverted dropout backward pass during training time.
# ===== #
dx = dout * mask
# ===== #
# END YOUR CODE HERE
# ===== #
    elif mode == 'test':
# ===== #
# YOUR CODE HERE:
# Implement the inverted dropout backward pass during test time.
# ===== #
dx = dout
# ===== #
# END YOUR CODE HERE
# ===== #
    return dx

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

```

Inputs:

- *x*: Input data, of shape (N, C) where $x[i, j]$ is the score for the j th class for the i th input.
- *y*: Vector of labels, of shape $(N,)$ where $y[i]$ is the label for $x[i]$ and $0 \leq y[i] < C$

Returns a tuple of:

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

"""

```
N = x.shape[0]
correct_class_scores = x[np.arange(N), y]
margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
margins[np.arange(N), y] = 0
loss = np.sum(margins) / N
num_pos = np.sum(margins > 0, axis=1)
dx = np.zeros_like(x)
dx[margins > 0] = 1
dx[np.arange(N), y] -= num_pos
dx /= N
return loss, dx
```

def softmax_loss(x, y):

"""

Computes the loss and gradient for softmax classification.

Inputs:

- *x*: Input data, of shape (N, C) where $x[i, j]$ is the score for the j th class for the i th input.
- *y*: Vector of labels, of shape $(N,)$ where $y[i]$ is the label for $x[i]$ and $0 \leq y[i] < C$

Returns a tuple of:

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

"""

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



```

In [ ]: import numpy as np
import pdb

from .layers import *
from .layer_utils import *

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

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

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

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

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

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

        Inputs:
        - hidden_dims: A list of integers giving the size of each hidden layer.
        - input_dim: An integer giving the size of the input.
        - num_classes: An integer giving the number of classes to classify.
        - dropout: Scalar between 0 and 1 giving dropout strength. If dropout=0 then
          the network should not use dropout at all.
        - use_batchnorm: Whether or not the network should use batch normalization.
        - reg: Scalar giving L2 regularization strength.
        - weight_scale: Scalar giving the standard deviation for random
          initialization of the weights.
        - dtype: A numpy datatype object; all computations will be performed using
          this datatype. float32 is faster but less accurate, so you should use
          float64 for numeric gradient checking.
        - seed: If not None, then pass this random seed to the dropout layers. This
          will make the dropout layers deterministic so we can gradient check the
          model.
        """
        self.use_batchnorm = use_batchnorm
        self.use_dropout = dropout > 0
        self.reg = reg
        self.num_layers = 1 + len(hidden_dims)
        self.dtype = dtype
        self.params = {}

```

```

# ===== #
# YOUR CODE HERE:
# Initialize all parameters of the network in the self.params dictionary.
# The weights and biases of layer 1 are W1 and b1; and in general the
# weights and biases of layer i are Wi and bi. The
# biases are initialized to zero and the weights are initialized
# so that each parameter has mean 0 and standard deviation weight_scale.
#
# BATCHNORM: Initialize the gammas of each layer to 1 and the beta
# parameters to zero. The gamma and beta parameters for layer 1 should
# be self.params['gamma1'] and self.params['beta1']. For layer 2, they
# should be gamma2 and beta2, etc. Only use batchnorm if self.use_batchnorm
# is true and DO NOT batch normalize the output scores.
# ===== #
all_dims = [input_dim] + hidden_dims + [num_classes]

#for d in all_dims:
# print(d)

for layer in range(self.num_layers):
    self.params['W{}'.format(layer + 1)] = np.random.normal(0, weight_scale, (a
ll_dims[layer], all_dims[layer + 1]))
    self.params['b{}'.format(layer + 1)] = np.zeros(all_dims[layer + 1])
    #print("{}x{}".format(all_dims[layer], all_dims[layer+1]))

for layer in range(self.num_layers - 1):
    if self.use_batchnorm:
        self.params['gamma{}'.format(layer + 1)] = np.ones(all_dims[layer + 1])
        #print('gamma{}'.format(layer + 1))
        self.params['beta{}'.format(layer + 1)] = np.zeros(all_dims[layer + 1])
        #print('beta{}'.format(layer + 1))
# ===== #
# END YOUR CODE HERE
# ===== #

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

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

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

```

```

def loss(self, X, y=None):
    """

```

Compute loss and gradient for the fully-connected net.

Input / output: Same as TwoLayerNet above.
"""

```
X = X.astype(self.dtype)
mode = 'test' if y is None else 'train'
```

*# Set train/test mode for batchnorm params and dropout param since they
behave differently during training and testing.*

```
if self.dropout_param is not None:
    self.dropout_param['mode'] = mode
if self.use_batchnorm:
    for bn_param in self.bn_params:
        bn_param[mode] = mode
```

```
scores = None
```

*# ===== #
YOUR CODE HERE:*

*# Implement the forward pass of the FC net and store the output
scores as the variable "scores".*

*#
BATCHNORM: If self.use_batchnorm is true, insert a bathnorm layer
between the affine_forward and relu_forward layers. You may
also write an affine_batchnorm_relu() function in layer_utils.py.*

*#
DROPOUT: If dropout is non-zero, insert a dropout layer after
every ReLU layer.*
=====

```
a = {}
```

```
bn = {}
```

```
h = {}
```

```
d = {}
```

```
h[0] = [X]
```

```
for layer in range(self.num_layers - 1):
```

```
    if self.use_batchnorm and not self.use_dropout:
```

```
        a[layer + 1] = affine_forward(h[layer][0], self.params['W{}'.format(layer
+ 1)], self.params['b{}'.format(layer + 1)])
        #print("a{}".format(layer + 1))
        bn[layer + 1] = batchnorm_forward(a[layer + 1][0], self.params['gamma{}'.
format(layer + 1)], self.params['beta{}'.format(layer + 1)], self.bn_params[layer
])
```

```
        #print("bn{}".format(layer + 1))
```

```
        h[layer + 1] = relu_forward(bn[layer + 1][0])
```

```
        #print("h{}".format(layer + 1))
```

```
    elif not self.use_batchnorm and self.use_dropout:
```

```
        if layer == 0:
```

```
            d[layer] = [X]
```

```
            a[layer + 1] = affine_forward(d[layer][0], self.params['W{}'.format(layer
+ 1)], self.params['b{}'.format(layer + 1)])
```

```
            h[layer + 1] = relu_forward(a[layer + 1][0])
```

```
            d[layer + 1] = dropout_forward(h[layer + 1][0], self.dropout_param)
```

```
    elif self.use_batchnorm and self.use_dropout:
```

```
        if layer == 0:
```

```
            d[layer] = [X]
```

```
            a[layer + 1] = affine_forward(d[layer][0], self.params['W{}'.format(layer
+ 1)], self.params['b{}'.format(layer + 1)])
```

```
            bn[layer + 1] = batchnorm_forward(a[layer + 1][0], self.params['gamma{}'.
format(layer + 1)], self.params['beta{}'.format(layer + 1)], self.bn_params[layer
])
```

```
            h[layer + 1] = relu_forward(bn[layer + 1][0])
```

```
            d[layer + 1] = dropout_forward(h[layer + 1][0], self.dropout_param)
```

```

        else:
            a[layer + 1] = affine_forward(h[layer][0], self.params['W{}'.format(layer
+ 1)], self.params['b{}'.format(layer + 1)])
            #print("a{}".format(layer + 1))
            h[layer + 1] = relu_forward(a[layer + 1][0])
            #print("h{}".format(layer + 1))

    if self.use_dropout:
        a[self.num_layers] = affine_forward(d[self.num_layers - 1][0], self.params[
'W{}'.format(self.num_layers)], self.params['b{}'.format(self.num_layers)])
    else:
        a[self.num_layers] = affine_forward(h[self.num_layers - 1][0], self.params[
'W{}'.format(self.num_layers)], self.params['b{}'.format(self.num_layers)])
        #print("a{}".format(self.num_layers))
        scores = a[self.num_layers][0]
        # ===== #
        # END YOUR CODE HERE
        # ===== #

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

    loss, grads = 0.0, {}
    # ===== #
    # YOUR CODE HERE:
    # Implement the backwards pass of the FC net and store the gradients
    # in the grads dict, so that grads[k] is the gradient of self.params[k]
    # Be sure your L2 regularization includes a 0.5 factor.
    #
    # BATCHNORM: Incorporate the backward pass of the batchnorm.
    #
    # DROPOUT: Incorporate the backward pass of dropout.
    # ===== #
    dl_da = {}
    dl_dbn = {}
    dl_dh = {}
    dl_dd = {}
    dl_dw = {}
    dl_db = {}
    loss, dl_da[self.num_layers] = softmax_loss(scores, y)
    #print("dl_da{}".format(self.num_layers))
    loss += 0.5 * self.reg * np.sum([np.sum(self.params['W{}'.format(layer + 1)]
** 2) for layer in range(self.num_layers)])

    dl_dh[self.num_layers - 1], dl_dw[self.num_layers], dl_db[self.num_layers] =
affine_backward(dl_da[self.num_layers], a[self.num_layers][1])
    #print("dl_dh{}".format(self.num_layers - 1))
    for layer in range(self.num_layers)[:0:-1]:
        if self.use_batchnorm and not self.use_dropout:
            dl_dbn[layer] = relu_backward(dl_dh[layer], h[layer][1])
            #print("dl_dbn{}".format(layer))
            dl_da[layer], grads['gamma{}'.format(layer)], grads['beta{}'.format(layer
)] = batchnorm_backward(dl_dbn[layer], bn[layer][1])
            #print("dl_da{}".format(layer))
            dl_dh[layer - 1], dl_dw[layer], dl_db[layer] = affine_backward(dl_da[laye
r], a[layer][1])
            #print("dl_dh{}".format(layer - 1))
        elif not self.use_batchnorm and self.use_dropout:
            if layer == self.num_layers - 1:
                dl_dd[layer] = dl_dh[self.num_layers - 1]
                dl_dh[layer] = dropout_backward(dl_dd[layer], d[layer][1])

```



```

        dl_da[layer] = relu_backward(dl_dh[layer], h[layer][1])
        dl_dd[layer - 1], dl_dw[layer], dl_db[layer] = affine_backward(dl_da[layer], a[layer][1])
    elif self.use_batchnorm and self.use_dropout:
        if layer == self.num_layers - 1:
            dl_dd[layer] = dl_dh[self.num_layers - 1]
            dl_dh[layer] = dropout_backward(dl_dd[layer], d[layer][1])
            dl_dbn[layer] = relu_backward(dl_dh[layer], h[layer][1])
            dl_da[layer], grads['gamma{}'.format(layer)], grads['beta{}'.format(layer)] = batchnorm_backward(dl_dbn[layer], bn[layer][1])
            dl_dd[layer - 1], dl_dw[layer], dl_db[layer] = affine_backward(dl_da[layer], a[layer][1])
        else:
            dl_da[layer] = relu_backward(dl_dh[layer], h[layer][1])
            #print("dl_da{}".format(layer))
            dl_dh[layer - 1], dl_dw[layer], dl_db[layer] = affine_backward(dl_da[layer], a[layer][1])
            #print("dl_dh{}".format(layer - 1))

    for layer in range(self.num_layers):
        grads['W{}'.format(layer + 1)] = dl_dw[layer + 1] + self.reg * self.params['W{}'.format(layer + 1)]
        grads['b{}'.format(layer + 1)] = dl_db[layer + 1]

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

    return loss, grads

```