

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_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,)
```

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: $1.1736859492107945e-09$

dw error: $9.903414542360811e-11$

db error: $3.658900127324957e-12$

If relu_forward function is working, difference should be around $1e-8$:

difference: $4.999999798022158e-08$

If relu_forward function is working, error should be less than $1e-9$:

dx error: $3.2756218084232666e-12$

If affine_relu_forward and affine_relu_backward are working, error should be less than $1e-9$:

dx error: $2.1553226229056188e-10$

dw error: $6.07535197721702e-10$

db error: $1.8928964784418166e-11$

Running check with reg = 0

Initial loss: 2.302737553701701

W1 relative error: $5.697158087460033e-06$

W2 relative error: $6.383999016022497e-07$

W3 relative error: $8.680262003512137e-07$

b1 relative error: $5.127074305892729e-07$

b2 relative error: $5.041477971652994e-07$

b3 relative error: $5.040832277334391e-07$

Running check with reg = 3.14

Initial loss: 7.1870429520308585

W1 relative error: $6.680587775733287e-07$

W2 relative error: $1.1663211181671765e-06$

W3 relative error: $1.5416056763702263e-06$

b1 relative error: $5.039402779395593e-07$

b2 relative error: $5.033985011124235e-07$

b3 relative error: $5.005245125813661e-07$

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)

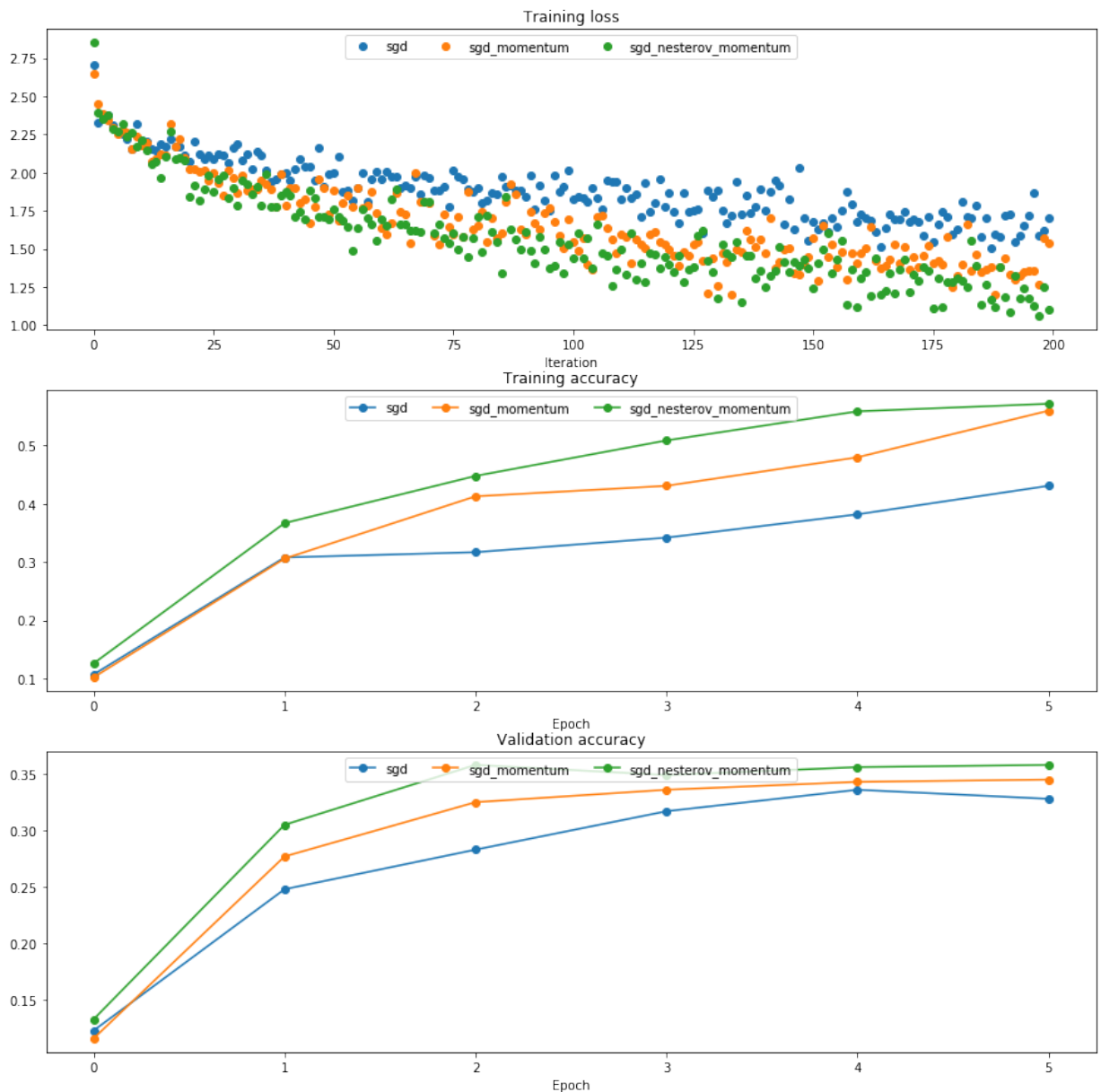
for i in [1, 2, 3]:
    plt.subplot(3, 1, i)
    plt.legend(loc='upper center', ncol=4)
plt.gcf().set_size_inches(15, 15)
plt.show()

```

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

/Users/Jonny/anaconda3/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()

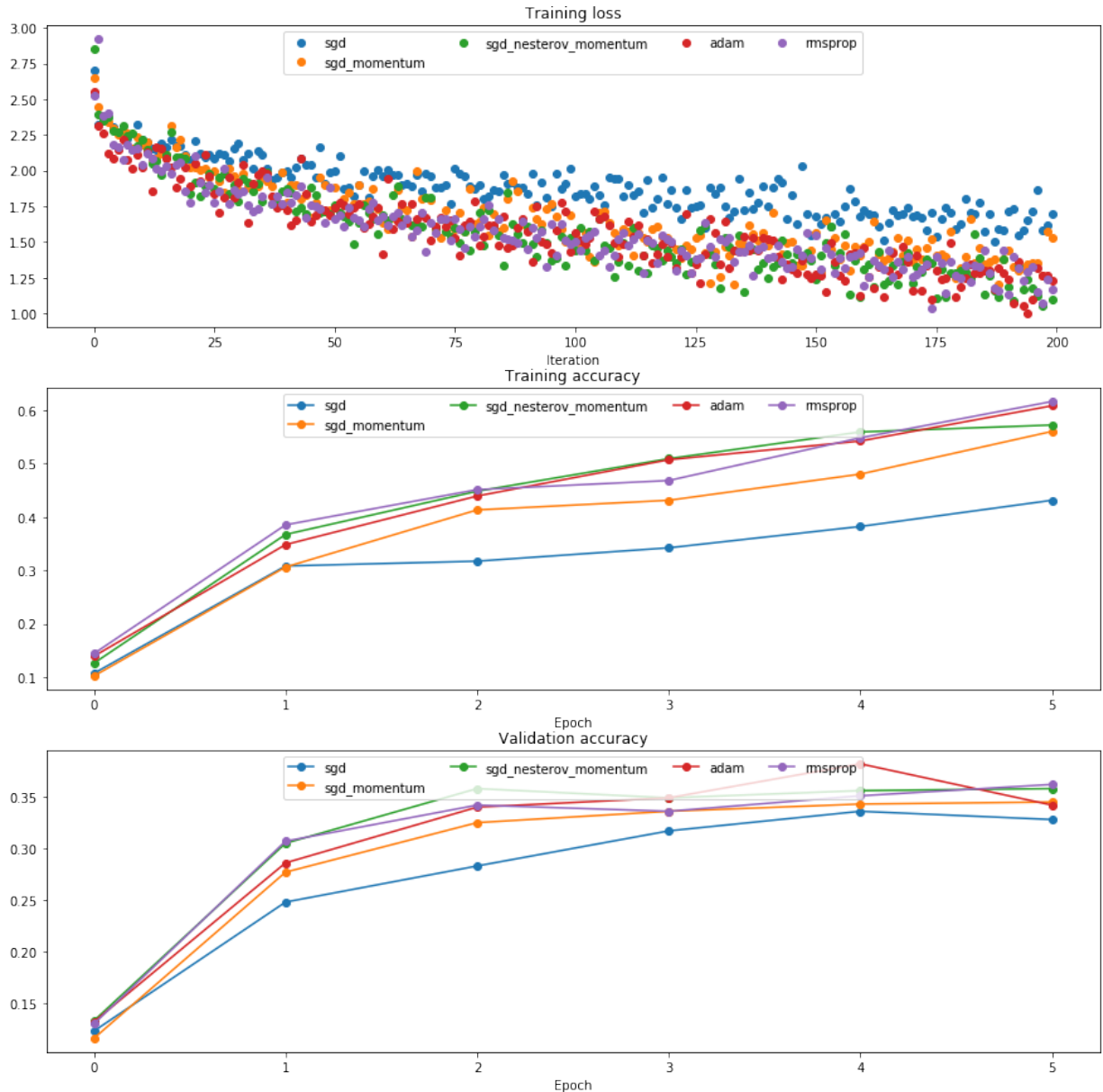
```

Optimizing with adam

Optimizing with rmsprop

/Users/Jonny/anaconda3/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: 2.318704
(Epoch 0 / 10) train acc: 0.215000; val_acc: 0.229000
(Iteration 51 / 4900) loss: 1.619281
(Iteration 101 / 4900) loss: 1.582920
(Iteration 151 / 4900) loss: 1.744086
(Iteration 201 / 4900) loss: 1.448380
(Iteration 251 / 4900) loss: 1.603279
(Iteration 301 / 4900) loss: 1.523390
(Iteration 351 / 4900) loss: 1.654481
(Iteration 401 / 4900) loss: 1.427959
(Iteration 451 / 4900) loss: 1.456330
(Epoch 1 / 10) train acc: 0.505000; val_acc: 0.482000
(Iteration 501 / 4900) loss: 1.540869
(Iteration 551 / 4900) loss: 1.412860
(Iteration 601 / 4900) loss: 1.452870
(Iteration 651 / 4900) loss: 1.350827
(Iteration 701 / 4900) loss: 1.319707
(Iteration 751 / 4900) loss: 1.182061
(Iteration 801 / 4900) loss: 1.307882
(Iteration 851 / 4900) loss: 1.343506
```

```
(Iteration 901 / 4900) loss: 1.225088
(Iteration 951 / 4900) loss: 1.201702
(Epoch 2 / 10) train acc: 0.560000; val_acc: 0.504000
(Iteration 1001 / 4900) loss: 1.350344
(Iteration 1051 / 4900) loss: 1.334488
(Iteration 1101 / 4900) loss: 1.451310
(Iteration 1151 / 4900) loss: 1.106676
(Iteration 1201 / 4900) loss: 1.085034
(Iteration 1251 / 4900) loss: 1.121413
(Iteration 1301 / 4900) loss: 1.126679
(Iteration 1351 / 4900) loss: 1.062616
(Iteration 1401 / 4900) loss: 1.268062
(Iteration 1451 / 4900) loss: 1.159597
(Epoch 3 / 10) train acc: 0.601000; val_acc: 0.544000
(Iteration 1501 / 4900) loss: 1.128401
(Iteration 1551 / 4900) loss: 1.159889
(Iteration 1601 / 4900) loss: 1.182442
(Iteration 1651 / 4900) loss: 1.091373
(Iteration 1701 / 4900) loss: 1.015887
(Iteration 1751 / 4900) loss: 1.175737
(Iteration 1801 / 4900) loss: 1.185392
(Iteration 1851 / 4900) loss: 1.157714
(Iteration 1901 / 4900) loss: 1.067200
(Iteration 1951 / 4900) loss: 0.981427
(Epoch 4 / 10) train acc: 0.642000; val_acc: 0.545000
(Iteration 2001 / 4900) loss: 1.035615
(Iteration 2051 / 4900) loss: 0.883135
(Iteration 2101 / 4900) loss: 1.094307
(Iteration 2151 / 4900) loss: 1.174929
(Iteration 2201 / 4900) loss: 0.947830
(Iteration 2251 / 4900) loss: 0.797088
(Iteration 2301 / 4900) loss: 0.953508
(Iteration 2351 / 4900) loss: 0.975237
(Iteration 2401 / 4900) loss: 1.168018
(Epoch 5 / 10) train acc: 0.644000; val_acc: 0.537000
(Iteration 2451 / 4900) loss: 1.053042
(Iteration 2501 / 4900) loss: 0.895468
(Iteration 2551 / 4900) loss: 0.920849
(Iteration 2601 / 4900) loss: 0.881266
(Iteration 2651 / 4900) loss: 0.836560
(Iteration 2701 / 4900) loss: 0.850926
(Iteration 2751 / 4900) loss: 0.882054
(Iteration 2801 / 4900) loss: 1.052449
(Iteration 2851 / 4900) loss: 0.899877
(Iteration 2901 / 4900) loss: 0.902509
(Epoch 6 / 10) train acc: 0.694000; val_acc: 0.542000
(Iteration 2951 / 4900) loss: 0.936894
(Iteration 3001 / 4900) loss: 0.759375
(Iteration 3051 / 4900) loss: 0.959044
(Iteration 3101 / 4900) loss: 0.875181
```

```

(Iteration 3151 / 4900) loss: 0.886862
(Iteration 3201 / 4900) loss: 0.643409
(Iteration 3251 / 4900) loss: 0.626558
(Iteration 3301 / 4900) loss: 0.849284
(Iteration 3351 / 4900) loss: 0.553705
(Iteration 3401 / 4900) loss: 0.820701
(Epoch 7 / 10) train acc: 0.736000; val_acc: 0.541000
(Iteration 3451 / 4900) loss: 0.691496
(Iteration 3501 / 4900) loss: 0.726798
(Iteration 3551 / 4900) loss: 0.809308
(Iteration 3601 / 4900) loss: 0.612413
(Iteration 3651 / 4900) loss: 0.711555
(Iteration 3701 / 4900) loss: 0.679913
(Iteration 3751 / 4900) loss: 0.579952
(Iteration 3801 / 4900) loss: 0.767239
(Iteration 3851 / 4900) loss: 0.668577
(Iteration 3901 / 4900) loss: 0.680020
(Epoch 8 / 10) train acc: 0.773000; val_acc: 0.556000
(Iteration 3951 / 4900) loss: 0.749292
(Iteration 4001 / 4900) loss: 0.790182
(Iteration 4051 / 4900) loss: 0.639617
(Iteration 4101 / 4900) loss: 0.576958
(Iteration 4151 / 4900) loss: 0.858437
(Iteration 4201 / 4900) loss: 0.583717
(Iteration 4251 / 4900) loss: 0.588226
(Iteration 4301 / 4900) loss: 0.760246
(Iteration 4351 / 4900) loss: 0.633393
(Iteration 4401 / 4900) loss: 0.481371
(Epoch 9 / 10) train acc: 0.777000; val_acc: 0.532000
(Iteration 4451 / 4900) loss: 0.637874
(Iteration 4501 / 4900) loss: 0.660831
(Iteration 4551 / 4900) loss: 0.717537
(Iteration 4601 / 4900) loss: 0.515468
(Iteration 4651 / 4900) loss: 0.533865
(Iteration 4701 / 4900) loss: 0.613792
(Iteration 4751 / 4900) loss: 0.483881
(Iteration 4801 / 4900) loss: 0.619108
(Iteration 4851 / 4900) loss: 0.400438
(Epoch 10 / 10) train acc: 0.832000; val_acc: 0.551000

```

```

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: 0.56
Test set accuracy: 0.556

```

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 [288]: ## 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))))
```

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

```
%reload_ext autoreload
```

```
In [289]: # 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 [290]: # 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: [ -4.77273782  29.43839171  13.98384247]
  stds: [ 26.35667822  29.20543313  32.15623214]
After batch normalization (gamma=1, beta=0)
  mean: [ -3.38618023e-17   6.66133815e-18  -7.54951657e-17]
  std: [ 0.99999999  0.99999999  1.          ]
After batch normalization (nontrivial gamma, beta)
  means: [ 11.  12.  13.]
  stds: [ 0.99999999  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 [291]: # 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:  [ 0.08593135 -0.03563948 -0.05269799]
stds:   [ 1.01983728  0.98978012  1.0656028 ]

```

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 [292]: # 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, dgamma, dbeta)

print('dx error: ', rel_error(dx_num, dx))
print('dgamma error: ', rel_error(da_num, dgamma))
print('dbeta error: ', rel_error(db_num, dbeta))

[[ -1.08393230e-01  -7.25486547e-03   6.54049767e-02   8.13297213e-0
 2
    -6.77128463e-02]
 [  1.49957469e-02  -2.06453845e-02  -6.25766188e-02   4.64155138e-0
 2
    1.35075912e-02]
 [ -2.16937699e-04   3.23538320e-02   6.35193145e-02   1.40816420e-0
 1
    1.75741579e-01]
 [  9.36144210e-02  -4.45358205e-03  -6.63476725e-02  -2.68561655e-0
 1
    -1.21536324e-01]] [ 2.46074274  1.17074281 -1.50157428 -0.5387052
 5  0.58101496] [ 1.47707902 -0.19294516  3.33409901  0.8237985  -0.6
 5006074]
dx error:  5.9114822945e-09
dgamma error:  3.47530270348e-12
dbeta error:  3.27572911404e-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 [293]: N, D, H1, H2, C = 2, 15, 20, 30, 10
X = np.random.randn(N, D)
y = np.random.randint(C, size=(N,))

print(X.shape)

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.float
64,
                              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, gra
ds[name])))
        if reg == 0: print('\n')
```

```
(2, 15)
Running check with reg = 0
Initial loss: 2.18570087163
W1 relative error: 0.0017620961330047724
W2 relative error: 1.5155245951274783e-05
W3 relative error: 4.789600547411837e-07
b1 relative error: 6.938893903907228e-09
b2 relative error: 6.661338147750939e-08
b3 relative error: 4.537302024147189e-07
beta1 relative error: 6.149690539945989e-07
beta2 relative error: 4.795865813281534e-07
gamma1 relative error: 6.210851395420151e-07
gamma2 relative error: 4.785530895915207e-07
```

```
Running check with reg = 3.14
Initial loss: 6.93152340975
W1 relative error: 2.124938522100045e-06
W2 relative error: 3.340676923521693e-05
W3 relative error: 6.237179420256357e-05
b1 relative error: 1.6653345369377348e-08
b2 relative error: 1.2212453270876722e-07
b3 relative error: 4.1214720252785485e-07
beta1 relative error: 5.051135620034601e-07
beta2 relative error: 4.374539805719088e-07
gamma1 relative error: 5.0293038836873e-07
gamma2 relative error: 4.37858551157119e-07
```

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 [294]: # 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)

print(small_data['X_train'].shape)
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()
```

```
(1000, 3, 32, 32)
(Iteration 1 / 200) loss: 2.309796
(Epoch 0 / 10) train acc: 0.154000; val_acc: 0.124000
(Epoch 1 / 10) train acc: 0.326000; val_acc: 0.269000
(Epoch 2 / 10) train acc: 0.431000; val_acc: 0.309000
(Epoch 3 / 10) train acc: 0.514000; val_acc: 0.321000
(Epoch 4 / 10) train acc: 0.572000; val_acc: 0.324000
(Epoch 5 / 10) train acc: 0.606000; val_acc: 0.323000
(Epoch 6 / 10) train acc: 0.665000; val_acc: 0.319000
(Epoch 7 / 10) train acc: 0.674000; val_acc: 0.332000
(Epoch 8 / 10) train acc: 0.729000; val_acc: 0.325000
(Epoch 9 / 10) train acc: 0.779000; val_acc: 0.332000
(Epoch 10 / 10) train acc: 0.779000; val_acc: 0.309000
(Iteration 1 / 200) loss: 2.303343
(Epoch 0 / 10) train acc: 0.116000; val_acc: 0.113000
(Epoch 1 / 10) train acc: 0.264000; val_acc: 0.222000
(Epoch 2 / 10) train acc: 0.269000; val_acc: 0.253000
(Epoch 3 / 10) train acc: 0.337000; val_acc: 0.262000
(Epoch 4 / 10) train acc: 0.351000; val_acc: 0.268000
(Epoch 5 / 10) train acc: 0.388000; val_acc: 0.303000
(Epoch 6 / 10) train acc: 0.462000; val_acc: 0.303000
(Epoch 7 / 10) train acc: 0.517000; val_acc: 0.326000
(Epoch 8 / 10) train acc: 0.575000; val_acc: 0.322000
(Epoch 9 / 10) train acc: 0.540000; val_acc: 0.297000
(Epoch 10 / 10) train acc: 0.618000; val_acc: 0.317000
```

```
In [295]: 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')

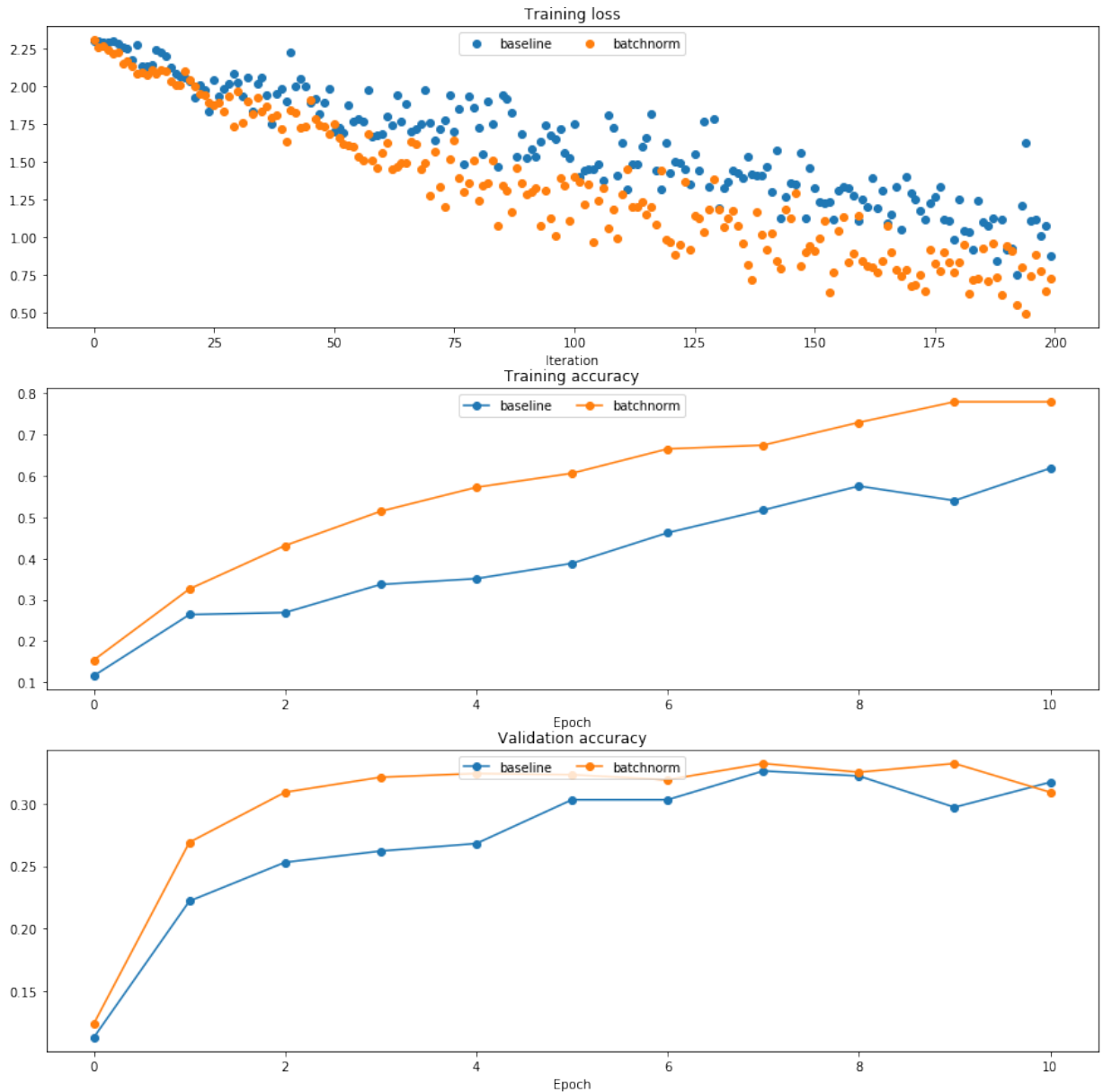
plt.subplot(3, 1, 2)
plt.plot(solver.train_acc_history, '-o', label='baseline')
plt.plot(bn_solver.train_acc_history, '-o', label='batchnorm')

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/Jonny/anaconda3/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 [296]: # 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, us
    e_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 [297]: # 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:]))

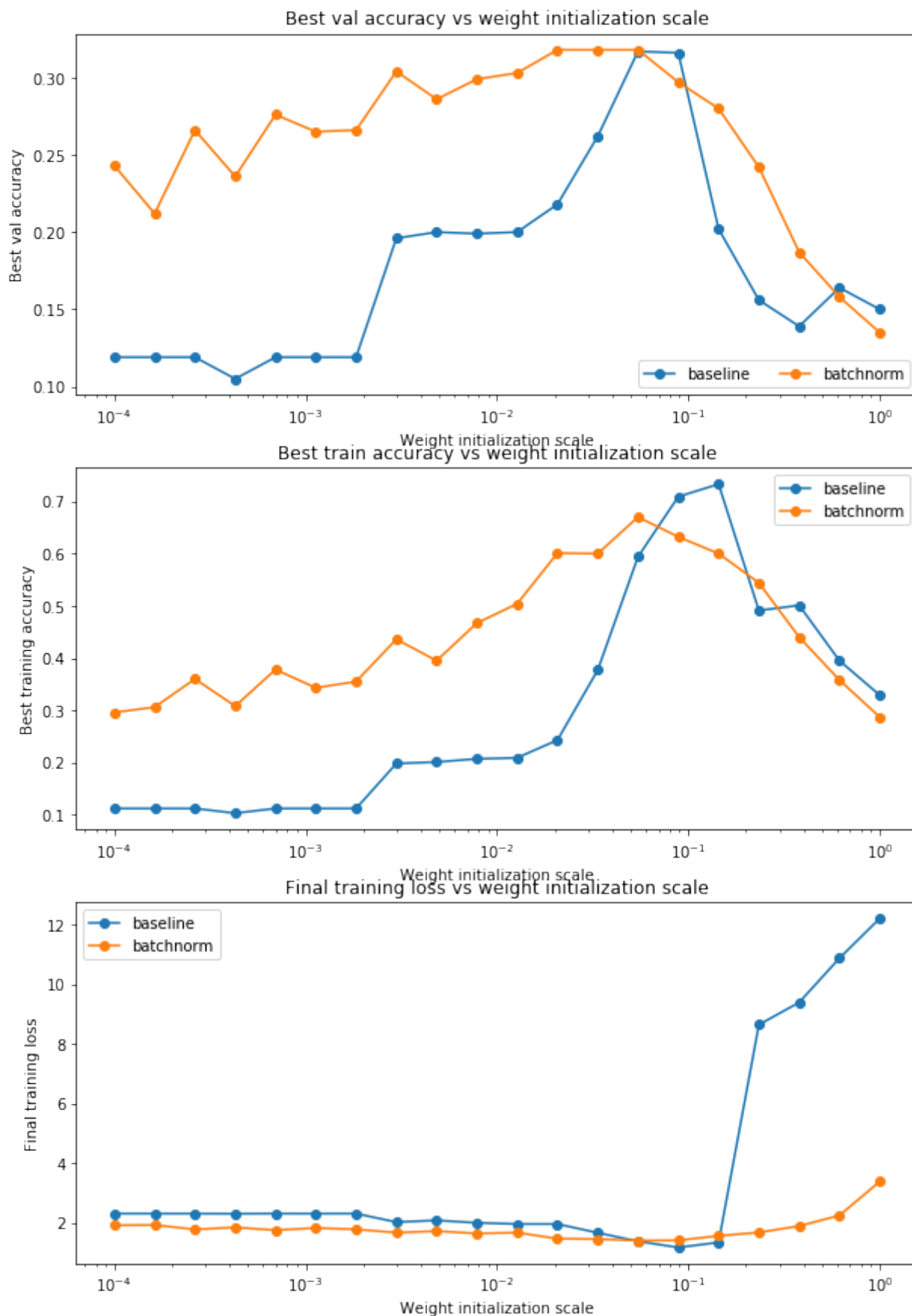
plt.subplot(3, 1, 1)
plt.title('Best val accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Best val accuracy')
plt.semilogx(weight_scales, best_val_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_val_accs, '-o', label='batchnorm')
plt.legend(ncol=2,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

These experiments show that the deep network with batch normalization is less sensitive to weight initializations. In the first experiment, the accuracy for the regular network is extremely low for low weight initialization standard deviations. While an optimal weight initialization is achievable (near 10^{-1}), this would require a hyperparameter sweep. By normalizing the data in each mini-batch, we are essentially making each layer independent in terms of mean and variance. The regular method (no batch normalization) is heavily dependent on weight initializations because these affect the mean and variance of the outputs as data propagates through the network.

The hypothesis is that a batch normalized network will exhibit a much smoother accuracy vs. weight initialization curve vs. the regular network. This is confirmed in all three of the figures. In the last figure, with large stddev for weight initializations, the training loss for the regular network explodes whereas the normalized network's loss increases but does not shoot up. All curves are smoother, confirming our hypothesis.

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 [133]: ## 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))))
```

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

```
%reload_ext autoreload
```

```
In [134]: # 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 [146]: 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: 9.9983572354
Mean of train-time output: 9.99057904364
Mean of test-time output: 9.9983572354
Fraction of train-time output set to zero: 0.300788
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.6
Mean of input: 9.9983572354
Mean of train-time output: 9.96138852574
Mean of test-time output: 9.9983572354
Fraction of train-time output set to zero: 0.601664
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.75
Mean of input: 9.9983572354
Mean of train-time output: 10.0117355182
Mean of test-time output: 9.9983572354
Fraction of train-time output set to zero: 0.749728
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 [147]: 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:  1.89290542075e-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 [148]: 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, gra
ds[name])))
    print('\n')
```

```
Running check with dropout = 0
Initial loss: 2.30519382479
W1 relative error: 7.711419522175973e-07
W2 relative error: 1.5034484932141387e-05
W3 relative error: 6.398873279808276e-07
b1 relative error: 2.9369574464090924e-06
b2 relative error: 6.320762652162454e-07
b3 relative error: 5.016029374797846e-07
```

```
Running check with dropout = 0.25
Initial loss: 2.29898515115
W1 relative error: 6.35222787011374e-06
W2 relative error: 5.687120844147058e-07
W3 relative error: 5.275962590939497e-07
b1 relative error: 5.030395134204798e-07
b2 relative error: 5.077063628933062e-07
b3 relative error: 5.009812500135304e-07
```

```
Running check with dropout = 0.5
Initial loss: 2.30243658786
W1 relative error: 5.836138382264886e-07
W2 relative error: 5.424266743801599e-07
W3 relative error: 8.333513391763223e-07
b1 relative error: 5.000127105947692e-07
b2 relative error: 5.057022318569306e-07
b3 relative error: 5.001185713330327e-07
```

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 [149]: *# 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: 0.350000; val_acc: 0.268000
(Epoch 14 / 25) train acc: 0.362000; val_acc: 0.275000
(Epoch 15 / 25) train acc: 0.394000; val_acc: 0.282000
(Epoch 16 / 25) train acc: 0.436000; val_acc: 0.296000
(Epoch 17 / 25) train acc: 0.438000; val_acc: 0.294000
(Epoch 18 / 25) train acc: 0.410000; val_acc: 0.305000
(Epoch 19 / 25) train acc: 0.388000; val_acc: 0.276000
(Epoch 20 / 25) train acc: 0.386000; val_acc: 0.286000
(Iteration 101 / 125) loss: 1.882976
(Epoch 21 / 25) train acc: 0.410000; val_acc: 0.288000
(Epoch 22 / 25) train acc: 0.448000; val_acc: 0.309000
(Epoch 23 / 25) train acc: 0.500000; val_acc: 0.308000
(Epoch 24 / 25) train acc: 0.486000; val_acc: 0.308000
(Epoch 25 / 25) train acc: 0.482000; val_acc: 0.304000
```

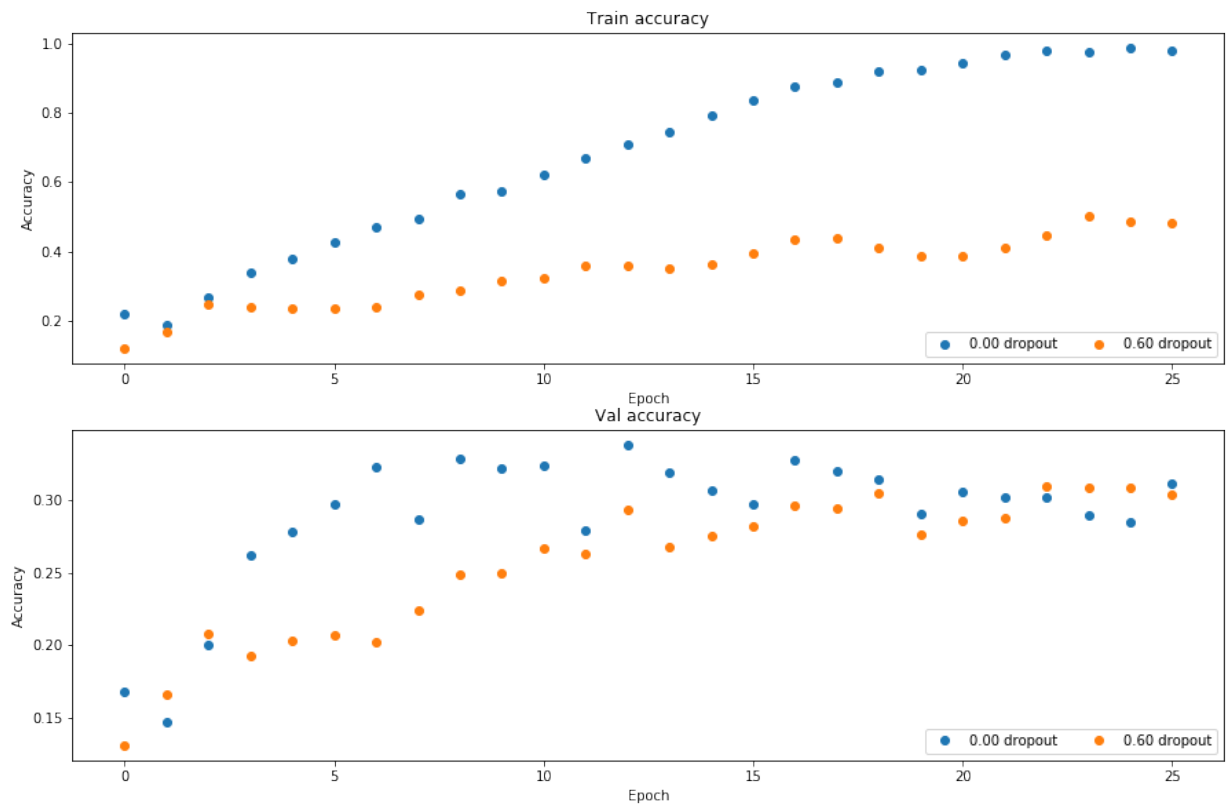
```
In [150]: # Plot train and validation accuracies of the two models

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:

Dropout is performing regularization. The train accuracy is lower and the validation accuracy is marginally higher. Dropout is essentially forcing the trained classifier to generalize well by randomly dropping certain neurons in each hidden layer. The expectation is that for deep neural networks, the validation error will be significantly better.

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.

Test 1000 # ----- #


```

111 [105]: """ ----- """
# YOUR CODE HERE:
# Implement a FC-net that achieves at least 60% validation accuracy
# on CIFAR-10.
# ===== #
import itertools

full_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'],
}

print(full_data['X_train'].shape)

dropout_choices = [0.1, 0.2]#[, 0.1, 0.2, 0.3] #0.3,0.4,0.5, 0.6]
#learning_rates = [5e-3]#[1e-4, 2e-4, 5e-4, 6e-4]
#batch_sizes = [500]#[100, 200, 300, 400, 500]
num_per_layer = [ 500]#[50, 100, 200, 300]
num_layers = [4]#[, 5, 6]

combos = list( itertools.product(dropout_choices, num_per_layer, num_l
ayers))
# Plot train and validation accuracies of the two models

train_accs = []
val_accs = []

i = 0
for combo in combos:
    print("combo: ", i)
    print(combo)
    dropout = combo[0]
    num_per_layer = int(combo[1])
    num_layers = int(combo[2])

    net_config = [num_per_layer for nl in range(num_layers)]
    #print(net_config)
    #break

    optimizer = 'adam'

    weight_scale = 0.01
    learning_rate = 1e-3
    lr_decay = 0.9

    model = FullyConnectedNet(net_config, weight_scale=weight_scale,
                              dropout=dropout, 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()

# solvers[dropout] = solver

    print("Train acc: ", solver.train_acc_history[-1], " || Val acc: "
, solver.val_acc_history[-1])
    train_accs.append(solver.train_acc_history[-1])
    val_accs.append(solver.val_acc_history[-1])

    if(solver.val_acc_history[-1] >= 0.6):
        break

    i += 1
# ===== #
# END YOUR CODE HERE
# ===== #
#
#4 layer
#(Epoch 10 / 10) train acc: 0.740000; val_acc: 0.573000
#Train acc: 0.74 || Val acc: 0.573

(49000, 3, 32, 32)
combo: 0
(0.1, 500, 4)
(Iteration 1 / 4900) loss: 2.296913
(Epoch 0 / 10) train acc: 0.180000; val_acc: 0.190000
(Iteration 51 / 4900) loss: 1.741179
(Iteration 101 / 4900) loss: 1.616187
(Iteration 151 / 4900) loss: 1.627492
(Iteration 201 / 4900) loss: 2.017857
(Iteration 251 / 4900) loss: 1.714556
(Iteration 301 / 4900) loss: 1.477805
(Iteration 351 / 4900) loss: 1.526273
(Iteration 401 / 4900) loss: 1.432306
(Iteration 451 / 4900) loss: 1.496049
(Epoch 1 / 10) train acc: 0.466000; val_acc: 0.465000
(Iteration 501 / 4900) loss: 1.467826
(Iteration 551 / 4900) loss: 1.368999
(Iteration 601 / 4900) loss: 1.402684
(Iteration 651 / 4900) loss: 1.366997
(Iteration 701 / 4900) loss: 1.394733

```

```
(Iteration 751 / 4900) loss: 1.304383
(Iteration 801 / 4900) loss: 1.481736
(Iteration 851 / 4900) loss: 1.434304
(Iteration 901 / 4900) loss: 1.329715
(Iteration 951 / 4900) loss: 1.201515
(Epoch 2 / 10) train acc: 0.549000; val_acc: 0.511000
(Iteration 1001 / 4900) loss: 1.420237
(Iteration 1051 / 4900) loss: 1.309935
(Iteration 1101 / 4900) loss: 1.158514
(Iteration 1151 / 4900) loss: 1.395335
(Iteration 1201 / 4900) loss: 1.244487
(Iteration 1251 / 4900) loss: 1.288770
(Iteration 1301 / 4900) loss: 1.198382
(Iteration 1351 / 4900) loss: 1.246748
(Iteration 1401 / 4900) loss: 1.409010
(Iteration 1451 / 4900) loss: 1.356864
(Epoch 3 / 10) train acc: 0.582000; val_acc: 0.531000
(Iteration 1501 / 4900) loss: 1.193911
(Iteration 1551 / 4900) loss: 1.187355
(Iteration 1601 / 4900) loss: 1.199924
(Iteration 1651 / 4900) loss: 1.310613
(Iteration 1701 / 4900) loss: 1.304499
(Iteration 1751 / 4900) loss: 1.364506
(Iteration 1801 / 4900) loss: 1.320047
(Iteration 1851 / 4900) loss: 1.139209
(Iteration 1901 / 4900) loss: 0.981234
(Iteration 1951 / 4900) loss: 1.194737
(Epoch 4 / 10) train acc: 0.598000; val_acc: 0.540000
(Iteration 2001 / 4900) loss: 1.152996
(Iteration 2051 / 4900) loss: 1.214479
(Iteration 2101 / 4900) loss: 0.941465
(Iteration 2151 / 4900) loss: 1.090359
(Iteration 2201 / 4900) loss: 1.105759
(Iteration 2251 / 4900) loss: 1.242445
(Iteration 2301 / 4900) loss: 1.048798
(Iteration 2351 / 4900) loss: 1.138683
(Iteration 2401 / 4900) loss: 1.180412
(Epoch 5 / 10) train acc: 0.622000; val_acc: 0.552000
(Iteration 2451 / 4900) loss: 1.202777
(Iteration 2501 / 4900) loss: 1.103409
(Iteration 2551 / 4900) loss: 1.079216
(Iteration 2601 / 4900) loss: 0.983463
(Iteration 2651 / 4900) loss: 1.095702
(Iteration 2701 / 4900) loss: 1.096826
(Iteration 2751 / 4900) loss: 1.029970
(Iteration 2801 / 4900) loss: 1.015585
(Iteration 2851 / 4900) loss: 1.034984
(Iteration 2901 / 4900) loss: 0.999851
(Epoch 6 / 10) train acc: 0.660000; val_acc: 0.548000
(Iteration 2951 / 4900) loss: 1.064506
```

```
(Iteration 3001 / 4900) loss: 1.141719
(Iteration 3051 / 4900) loss: 0.924702
(Iteration 3101 / 4900) loss: 1.066817
(Iteration 3151 / 4900) loss: 1.068603
(Iteration 3201 / 4900) loss: 0.949488
(Iteration 3251 / 4900) loss: 1.218299
(Iteration 3301 / 4900) loss: 0.914807
(Iteration 3351 / 4900) loss: 0.769179
(Iteration 3401 / 4900) loss: 1.095191
(Epoch 7 / 10) train acc: 0.674000; val_acc: 0.547000
(Iteration 3451 / 4900) loss: 0.982459
(Iteration 3501 / 4900) loss: 1.213621
(Iteration 3551 / 4900) loss: 0.819170
(Iteration 3601 / 4900) loss: 1.089411
(Iteration 3651 / 4900) loss: 0.972153
(Iteration 3701 / 4900) loss: 0.953362
(Iteration 3751 / 4900) loss: 0.939898
(Iteration 3801 / 4900) loss: 0.957945
(Iteration 3851 / 4900) loss: 0.955671
(Iteration 3901 / 4900) loss: 0.960321
(Epoch 8 / 10) train acc: 0.672000; val_acc: 0.553000
(Iteration 3951 / 4900) loss: 1.000361
(Iteration 4001 / 4900) loss: 1.043849
(Iteration 4051 / 4900) loss: 0.832510
(Iteration 4101 / 4900) loss: 0.921644
(Iteration 4151 / 4900) loss: 1.053920
(Iteration 4201 / 4900) loss: 0.816299
(Iteration 4251 / 4900) loss: 0.814988
(Iteration 4301 / 4900) loss: 0.745224
(Iteration 4351 / 4900) loss: 0.936560
(Iteration 4401 / 4900) loss: 0.856821
(Epoch 9 / 10) train acc: 0.718000; val_acc: 0.547000
(Iteration 4451 / 4900) loss: 0.749012
(Iteration 4501 / 4900) loss: 0.530407
(Iteration 4551 / 4900) loss: 0.826853
(Iteration 4601 / 4900) loss: 0.792176
(Iteration 4651 / 4900) loss: 0.802118
(Iteration 4701 / 4900) loss: 0.829370
(Iteration 4751 / 4900) loss: 0.995933
(Iteration 4801 / 4900) loss: 0.690810
(Iteration 4851 / 4900) loss: 0.781138
(Epoch 10 / 10) train acc: 0.762000; val_acc: 0.562000
Train acc: 0.762 || Val acc: 0.562
combo: 1
(0.2, 500, 4)
(Iteration 1 / 4900) loss: 2.315019
(Epoch 0 / 10) train acc: 0.198000; val_acc: 0.194000
(Iteration 51 / 4900) loss: 1.825066
(Iteration 101 / 4900) loss: 1.524264
(Iteration 151 / 4900) loss: 1.743650
```

```
(Iteration 201 / 4900) loss: 1.585431
(Iteration 251 / 4900) loss: 1.786960
(Iteration 301 / 4900) loss: 1.665943
(Iteration 351 / 4900) loss: 1.527417
(Iteration 401 / 4900) loss: 1.400422
(Iteration 451 / 4900) loss: 1.473420
(Epoch 1 / 10) train acc: 0.477000; val_acc: 0.454000
(Iteration 501 / 4900) loss: 1.624516
(Iteration 551 / 4900) loss: 1.615464
(Iteration 601 / 4900) loss: 1.432321
(Iteration 651 / 4900) loss: 1.513342
(Iteration 701 / 4900) loss: 1.328676
(Iteration 751 / 4900) loss: 1.563127
(Iteration 801 / 4900) loss: 1.457114
(Iteration 851 / 4900) loss: 1.265775
(Iteration 901 / 4900) loss: 1.477899
(Iteration 951 / 4900) loss: 1.351059
(Epoch 2 / 10) train acc: 0.517000; val_acc: 0.497000
(Iteration 1001 / 4900) loss: 1.543950
(Iteration 1051 / 4900) loss: 1.354283
(Iteration 1101 / 4900) loss: 1.440769
(Iteration 1151 / 4900) loss: 1.340905
(Iteration 1201 / 4900) loss: 1.339980
(Iteration 1251 / 4900) loss: 1.249155
(Iteration 1301 / 4900) loss: 1.073411
(Iteration 1351 / 4900) loss: 1.254554
(Iteration 1401 / 4900) loss: 1.259344
(Iteration 1451 / 4900) loss: 1.460359
(Epoch 3 / 10) train acc: 0.540000; val_acc: 0.513000
(Iteration 1501 / 4900) loss: 1.204073
(Iteration 1551 / 4900) loss: 1.223408
(Iteration 1601 / 4900) loss: 1.294228
(Iteration 1651 / 4900) loss: 1.224925
(Iteration 1701 / 4900) loss: 1.288375
(Iteration 1751 / 4900) loss: 1.152820
(Iteration 1801 / 4900) loss: 1.149959
(Iteration 1851 / 4900) loss: 1.037948
(Iteration 1901 / 4900) loss: 1.325537
(Iteration 1951 / 4900) loss: 1.293086
(Epoch 4 / 10) train acc: 0.567000; val_acc: 0.522000
(Iteration 2001 / 4900) loss: 1.390587
(Iteration 2051 / 4900) loss: 1.159214
(Iteration 2101 / 4900) loss: 1.299782
(Iteration 2151 / 4900) loss: 1.240472
(Iteration 2201 / 4900) loss: 1.217431
(Iteration 2251 / 4900) loss: 1.019275
(Iteration 2301 / 4900) loss: 1.112467
(Iteration 2351 / 4900) loss: 1.176747
(Iteration 2401 / 4900) loss: 1.128550
(Epoch 5 / 10) train acc: 0.613000; val_acc: 0.558000
```

```
(Iteration 2451 / 4900) loss: 1.173713
(Iteration 2501 / 4900) loss: 1.203973
(Iteration 2551 / 4900) loss: 1.021987
(Iteration 2601 / 4900) loss: 1.366045
(Iteration 2651 / 4900) loss: 1.366005
(Iteration 2701 / 4900) loss: 1.088223
(Iteration 2751 / 4900) loss: 1.144566
(Iteration 2801 / 4900) loss: 1.063326
(Iteration 2851 / 4900) loss: 1.243522
(Iteration 2901 / 4900) loss: 1.194539
(Epoch 6 / 10) train acc: 0.635000; val_acc: 0.557000
(Iteration 2951 / 4900) loss: 1.066216
(Iteration 3001 / 4900) loss: 1.128677
(Iteration 3051 / 4900) loss: 0.954731
(Iteration 3101 / 4900) loss: 1.322600
(Iteration 3151 / 4900) loss: 1.207659
(Iteration 3201 / 4900) loss: 1.011692
(Iteration 3251 / 4900) loss: 0.970870
(Iteration 3301 / 4900) loss: 1.132219
(Iteration 3351 / 4900) loss: 1.038042
(Iteration 3401 / 4900) loss: 1.026004
(Epoch 7 / 10) train acc: 0.659000; val_acc: 0.576000
(Iteration 3451 / 4900) loss: 0.965241
(Iteration 3501 / 4900) loss: 0.998311
(Iteration 3551 / 4900) loss: 1.055191
(Iteration 3601 / 4900) loss: 1.131579
(Iteration 3651 / 4900) loss: 1.052036
(Iteration 3701 / 4900) loss: 1.173005
(Iteration 3751 / 4900) loss: 1.224620
(Iteration 3801 / 4900) loss: 1.020468
(Iteration 3851 / 4900) loss: 0.988655
(Iteration 3901 / 4900) loss: 1.031473
(Epoch 8 / 10) train acc: 0.688000; val_acc: 0.571000
(Iteration 3951 / 4900) loss: 1.031489
(Iteration 4001 / 4900) loss: 0.958201
(Iteration 4051 / 4900) loss: 1.091165
(Iteration 4101 / 4900) loss: 1.182361
(Iteration 4151 / 4900) loss: 1.242961
(Iteration 4201 / 4900) loss: 0.998609
(Iteration 4251 / 4900) loss: 1.025409
(Iteration 4301 / 4900) loss: 0.995536
(Iteration 4351 / 4900) loss: 1.069902
(Iteration 4401 / 4900) loss: 1.035483
(Epoch 9 / 10) train acc: 0.700000; val_acc: 0.564000
(Iteration 4451 / 4900) loss: 0.926242
(Iteration 4501 / 4900) loss: 0.739709
(Iteration 4551 / 4900) loss: 1.189041
(Iteration 4601 / 4900) loss: 1.003469
(Iteration 4651 / 4900) loss: 1.079360
(Iteration 4701 / 4900) loss: 1.342456
```

```
(Iteration 4751 / 4900) loss: 1.079246
(Iteration 4801 / 4900) loss: 1.084932
(Iteration 4851 / 4900) loss: 1.007672
(Epoch 10 / 10) train acc: 0.730000; val_acc: 0.562000
Train acc: 0.73 || Val acc: 0.562
```

```
In [164]: max_train = np.argmax(train_accs)
max_val = np.argmax(val_accs)

print("best train combo: ", combos[max_train])
print("best val combo: ", combos[max_val])

best train combo: (0.1, 500, 4)
best val combo: (0.1, 500, 4)
```

Best Results

From the small sweep above and some testing done in another notebook the best configuration was:

500 neurons per layer

4 layers

1e-3 learning rate

0.01 weight scale

0.9 lr_decay

0.1 dropout

```

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.us
e_batchnorm
    # is true and DO NOT batch normalize the output scores.
    # =====e=====
= #
    mu = 0
    stddev = weight_scale
    """
    self.params['W1'] = std * np.random.randn(hidden_size, input_size)
    self.params['b1'] = np.zeros(hidden_size)
    self.params['W2'] = std * np.random.randn(output_size, hidden_size
)
    self.params['b2'] = np.zeros(output_size)

    np.random.normal(mu, stddev, <size>)
    """
    #aggregate all the dims into a single array that we can reference
    #input and output dim (num_classes) will only be used once
    aggregated_dims = [input_dim] + hidden_dims + [num_classes]
    for i in range(self.num_layers):
        #batchnorm on all layers except last one
        #init gammas to 1s and betas to 0
        if self.use_batchnorm and (i != (self.num_layers - 1)):
            self.params['gamma'+str(i+1)] = np.ones(aggregated_dims[i+1])
            self.params['beta'+str(i+1)] = np.zeros(aggregated_dims[i+1])

            self.params['b'+str(i+1)] = np.zeros(aggregated_dims[i+1])
            self.params['W'+str(i+1)] = np.random.normal(mu, stddev, size=(a
ggregated_dims[i], aggregated_dims[i+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 an
d the mode
    # (train / test). You can pass the same dropout_param to each drop
out 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 fo
    rward 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:
        #for i in range(self.num_layers):
        self.bn_params = [{'mode': 'train'} for i in range(self.num_laye
rs - 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 la
    yer
    # 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.
# =====
#

nn_layer = {}
nn_cache = {}
batchnorm_cache = {}
dropout_cache = {}

#initialize the first layer with the inputs
nn_layer[0] = X
#pass through each layer
for i in range(1, self.num_layers):
#    print("iteration", i)
    gamma_idx = 'gamma'+str(i)
    beta_idx = 'beta'+str(i)
    w_idx = 'W'+str(i)
    b_idx = 'b'+str(i)
    #affine relu forward takes (x, w, b)
    if self.use_batchnorm:
        #args: x, gamma, beta, bn_param
        #    nn_layer[i], batchnorm_cache[i] = batchnorm_forward(nn_layer[
i-1], self.params['gamma'+str(i)], self.params['beta'+str(i)], self.bn
_params[i-1])
        #    print(nn_layer[i-1].shape, self.params[w_idx].shape)
        nn_layer[i], nn_cache[i] = affine_batchnorm_relu_forward(nn_la
yer[i-1], self.params[w_idx],
                                                                    self
.params[b_idx], self.params[gamma_idx],
                                                                    self
.params[beta_idx], self.bn_params[i-1])
    else:
        nn_layer[i], nn_cache[i] = affine_relu_forward(nn_layer[i-1],
self.params[w_idx], self.params[b_idx])

    if(self.use_dropout):
        nn_layer[i], dropout_cache[i] = dropout_forward(nn_layer[i], s
elf.dropout_param)
    #all layers will have the affine_relu except for the last layer, w
hich is a passthrough
    #affine_forward takes (x, w, b) and outputs out, cache
    w_idx = 'W'+str(self.num_layers)
    b_idx = 'b'+str(self.num_layers)
    scores, cached_scores = affine_forward(nn_layer[self.num_layers -1
], self.params[w_idx], self.params[b_idx])

    nn_cache[self.num_layers] = cached_scores
# =====

```

```

#
# 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.
# =====
#

#get loss w/ softmax loss
loss, grad_loss = softmax_loss(scores, y)

#add L2 regularization to loss 1/2*np.sum(w**2)
for i in range(1, self.num_layers + 1):
    cur_weight_matrix = self.params['W'+str(i)]
    loss += 0.5 * self.reg * np.sum(cur_weight_matrix**2)

    """
    Backpropping into the (n-1)th layer will be different because we don't have
    the relu. Use affine_backward and then for each previous layer apply affine_relu_backward
    affine_backward takes dout, cache and returns dx, dw, db
    affine_relu_backward takes dout, cache and returns dx, dw, db
    """

    dx={}
    w_idx_nth = 'W'+str(self.num_layers)
    b_idx_nth = 'b'+str(self.num_layers)
    dx[self.num_layers], grads[w_idx_nth], grads[b_idx_nth] = affine_backward(grad_loss, cached_scores)

# print(dx[3])
#regularize
grads[w_idx_nth] += self.reg * self.params[w_idx_nth]

```

```

    #we apply affine_relu_backward now
    for i in range(self.num_layers - 1, 0, -1):
#       print(i, self.num_layers)

        #dx, dw, db
        w_idx = 'W' + str(i)#+1)
        b_idx = 'b' + str(i)#+1)

        gamma_idx = 'gamma' + str(i)#+1)
        beta_idx = 'beta' + str(i)#+1)

        if self.use_dropout:
            dx[i+1] = dropout_backward(dx[i+1], dropout_cache[i])
        if self.use_batchnorm:
            dx[i], grads[w_idx], grads[b_idx], grads[gamma_idx], grads[bet
a_idx] = affine_batchnorm_relu_backward(dx[i+1], nn_cache[i])

        else:
            #dout input to affine_relu_backward is the
            dx[i], grads[w_idx], grads[b_idx] = affine_relu_backward( dx[i
+1], nn_cache[i])

            #regularize
#       print(grads[w_idx].shape, self.params[w_idx].shape )
            grads[w_idx] += self.reg * self.params[w_idx]

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

```

```

In [ ]: from .layers 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.
"""

def affine_relu_forward(x, w, b):
    """
    Convenience layer that performs an affine transform followed by a ReLU

    Inputs:
    - x: Input to the affine layer
    - w, b: Weights for the affine layer

    Returns a tuple of:
    - out: Output from the ReLU
    - cache: Object to give to the backward pass
    """
    a, fc_cache = affine_forward(x, w, b)
    out, relu_cache = relu_forward(a)
    cache = (fc_cache, relu_cache)
    return out, cache

def affine_relu_backward(dout, cache):
    """
    Backward pass for the affine-relu convenience layer
    """
    fc_cache = cache[0]
    relu_cache = cache[1]
    # print("fc cache", fc_cache)
    # print(len(cache))

    da = relu_backward(dout, relu_cache)
    dx, dw, db = affine_backward(da, fc_cache)
    return dx, dw, db

```

```

def affine_batchnorm_relu_forward(x, w, b, gamma, beta, bn_params):
    """
    Performs affine transformation, batchnorm, and ReLU

    Returns all caches

    BN forward takes: def batchnorm_forward(x, gamma, beta, bn_param):

    """
    out, forward_cache = affine_forward(x, w, b)
    # print("beta received: ", beta.shape)
    out, batchnorm_cache = batchnorm_forward(out, gamma, beta, bn_params)
    )
    # print("got dim: ", out.dim)
    out, relu_cache = relu_forward(out)

    total_cache = (forward_cache, relu_cache, batchnorm_cache)
    # print("returning out dim: ", out.shape)
    return out, total_cache

def affine_batchnorm_relu_backward(dout, cache):
    """
    Backward pass
    def batchnorm_backward(dout, cache):
    def relu_backward(dout, cache):

    """
    #unpack the cache tuple
    forward_cache, relu_cache, batchnorm_cache = cache

    dx = relu_backward(dout, relu_cache)
    dx, dgamma, dbeta = batchnorm_backward(dx, batchnorm_cache)
    dx, dw, db = affine_backward(dx, forward_cache)

    gradients = dx, dw, db, dgamma, dbeta
    return gradients

```



```

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
    nd
    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.
    # ===== #

    # N = x.shape[0]

```

```

#D = w.shape[0]
#x_resaped = np.reshape(x, (N,D))
x_shape = x.shape
#Reshaping it as N*D
#x_shape[0] is equal to N
x = x.reshape( [ x_shape[0], np.prod( x_shape[1:]) ] )
out = x.dot(w) + b

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

cache = (x, w, b, x_shape)
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, x_shape = cache
    dx, dw, db = None, None, None

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

    #reshape x matrix to be N, D and multiply upstream for the chain rule
    """
    N = x.shape[0]
    D = w.shape[0]
    reshaped_x = np.reshape(x, (N, D))
    dw = reshaped_x.T.dot(dout)

    #derivative wrt x
    dx_raw = dout.dot(w.T)
    dx = np.reshape(dx_raw, x.shape)

```

```

    #sum derivative for bias
    db = np.sum(dout, axis=0)
    """
# print
dx = np.zeros_like(x)
dw = np.zeros_like(w)
db = np.zeros_like(b)

dx += dout.dot(w.T)
dw += x.T.dot(dout)
db += dout.sum( axis = 0)

# Reshaping dx
dx = dx.reshape(x_shape)

# ===== #
# 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)
out = np.maximum(x, np.zeros_like(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
    # ===== #

    #ReLU backward pass multiplies the dout by the indicator function
    #arr[arr > 255] = x
    dx = dout

    #apply indicator. Uses < and not <= because 0 is undefined for ReLU
    dx[x < 0] = 0
    # ===== #
    # 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

```

```

n
    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 parameter 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)

    # print("received x: ", x.shape)
    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.
ce.
# (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
n
# the 'cache' variable.
# =====
#
sample_mean = np.mean(x, axis=0)
sample_var = np.var(x, axis=0)

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

x_hat = (x - sample_mean) / np.sqrt(sample_var + eps)
out = x_hat*gamma + beta

#store in cache
cache = (mode, x, gamma, sample_mean, sample_var, x_hat, out, eps)

# =====
#
# 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'.
# =====
#

stddev = np.sqrt(running_var + eps)
x_hat = (x - running_mean)/stddev
out = x_hat*gamma + beta

#store in cache
cache = (mode, x, gamma, x_hat, out, eps, stddev)

```

```

# =====
#
# 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

# print(out.shape)
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
    mode = cache[0]
    # ===== #
    # YOUR CODE HERE:
    # Implement the batchnorm backward pass, calculating dx, dgamma, and dbeta.
    # ===== #
    if(mode == 'train'):
        mode, x, gamma, sample_mean, sample_var, x_hat, out, eps = cache
        # print(cache)

        N, D = x.shape

```

```

    dl_dbeta = np.sum(dout, axis=0)
#    print(dout.shape, x_hat.shape)
    dl_dgamma = np.sum(dout*x_hat, axis=0)
    dl_dx = dout*gamma

    dl_da = (1/np.sqrt(sample_var + eps))*dl_dx
    dl_du = -(1/np.sqrt(sample_var+eps))*np.sum(dl_dx, axis=0)

    dl_de = -0.5*(1/(sample_var+eps))*(x_hat)*dl_dx

    dl_dvar = np.sum(dl_de, axis=0)

    dl_da = (1/(np.sqrt(sample_var + eps)))*dl_dx

    dx = dl_da + 2*((x-sample_mean)/N)*dl_dvar + (1/N)*dl_du
    dgamma = dl_dgamma
    dbeta = dl_dbeta

elif(mode == 'test'):
    mode, x, gamma, x_hat, out, eps, stddev = cache
    dl_dbeta = np.sum(dout, axis=0)
    dl_dgamma = np.sum(dout*x_hat, axis=0)
    dx = (gamma*dout)/stddev

# ===== #
# 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.
    # =====
    #
    # print(x.shape)
    mask = (np.random.rand(*x.shape) < (1-p)) / (1-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
jth class
    for the ith input.
- y: Vector of labels, of shape (N,) where y[i] is the label for x[i
] and
    0 <= y[i] < C

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

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

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

    Returns a tuple of:
    - loss: Scalar giving the loss
    - dx: Gradient of the loss with respect to x
    """
    eps = 1e-7
    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] + eps)) / N
    dx = probs.copy()
    dx[np.arange(N), y] -= 1
    dx /= N
    return loss, dx

```


In []: **import numpy as np**

```

"""
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.
"""

"""
This file implements various first-order update rules that are commonl
y used for
training neural networks. Each update rule accepts current weights and
the
gradient of the loss with respect to those weights and produces the ne
xt set of
weights. Each update rule has the same interface:

def update(w, dw, config=None):

Inputs:
- w: A numpy array giving the current weights.
- dw: A numpy array of the same shape as w giving the gradient of th
e
    loss with respect to w.
- config: A dictionary containing hyperparameter values such as lear
ning rate,
    momentum, etc. If the update rule requires caching values over man
y
    iterations, then config will also hold these cached values.

Returns:
- next_w: The next point after the update.
- config: The config dictionary to be passed to the next iteration o
f the
    update rule.

NOTE: For most update rules, the default learning rate will probably n
ot perform
well; however the default values of the other hyperparameters should w

```

ork well
for a variety of different problems.

For efficiency, update rules may perform in-place updates, mutating `w` and setting `next_w` equal to `w`.
"""

```
def sgd(w, dw, config=None):
    """
```

Performs vanilla stochastic gradient descent.

config format:

- learning_rate: Scalar learning rate.
"""

```
if config is None: config = {}
config.setdefault('learning_rate', 1e-2)
```

```
w -= config['learning_rate'] * dw
return w, config
```

```
def sgd_momentum(w, dw, config=None):
    """
```

Performs stochastic gradient descent with momentum.

config format:

- learning_rate: Scalar learning rate.
- momentum: Scalar between 0 and 1 giving the momentum value.
Setting momentum = 0 reduces to sgd.
- velocity: A numpy array of the same shape as w and dw used to store a moving average of the gradients.
"""

```
if config is None: config = {}
config.setdefault('learning_rate', 1e-2)
config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.
```

```
# ===== #
# YOUR CODE HERE:
# Implement the momentum update formula. Return the updated weights
# as next_w, and store the updated velocity as v.
# ===== #
v = config['momentum']*v - config['learning_rate']*dw
next_w = w + v
```

```

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

config['velocity'] = v

return next_w, config

def sgd_nesterov_momentum(w, dw, config=None):
    """
    Performs stochastic gradient descent with Nesterov momentum.

    config format:
    - learning_rate: Scalar learning rate.
    - momentum: Scalar between 0 and 1 giving the momentum value.
      Setting momentum = 0 reduces to sgd.
    - velocity: A numpy array of the same shape as w and dw used to store a moving
      average of the gradients.
    """
    if config is None: config = {}
    config.setdefault('learning_rate', 1e-2)
    config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
    v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.

    # ===== #
    # YOUR CODE HERE:
    # Implement the momentum update formula. Return the updated weights
    # as next_w, and store the updated velocity as v.
    # ===== #
    v_old = v
    v = config['momentum']*v - config['learning_rate']*dw
    next_w = w + v + config['momentum']*(v-v_old)

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

    config['velocity'] = v

    return next_w, config

def rmsprop(w, dw, config=None):
    """
    Uses the RMSProp update rule, which uses a moving average of squared
    gradient

```

```

values to set adaptive per-parameter learning rates.

config format:
- learning_rate: Scalar learning rate.
- decay_rate: Scalar between 0 and 1 giving the decay rate for the squared
  gradient cache.
- epsilon: Small scalar used for smoothing to avoid dividing by zero
.
- beta: Moving average of second moments of gradients.
"""
if config is None: config = {}
config.setdefault('learning_rate', 1e-2)
config.setdefault('decay_rate', 0.99)
config.setdefault('epsilon', 1e-8)
config.setdefault('a', np.zeros_like(w))

next_w = None

# ===== #
# YOUR CODE HERE:
#   Implement RMSProp. Store the next value of w as next_w. You need
#   to also store in config['a'] the moving average of the second
#   moment gradients, so they can be used for future gradients. Conc
retely,
#   config['a'] corresponds to "a" in the lecture notes.
# ===== #

#hadamard product is taken care of by np multiplication
a = config['a']
beta = config['decay_rate']

config['a'] = beta*a + (1-beta)*np.multiply(dw, dw)

#update gradient
next_w = w - np.multiply(config['learning_rate']/(np.sqrt(config['a']
))+config['epsilon']), dw)

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

return next_w, config

def adam(w, dw, config=None):
    """
    Uses the Adam update rule, which incorporates moving averages of both the

```



```

gradient and its square and a bias correction term.

config format:
- learning_rate: Scalar learning rate.
- betal: Decay rate for moving average of first moment of gradient.
- beta2: Decay rate for moving average of second moment of gradient.
- epsilon: Small scalar used for smoothing to avoid dividing by zero
.
- m: Moving average of gradient.
- v: Moving average of squared gradient.
- t: Iteration number.
"""

if config is None: config = {}
config.setdefault('learning_rate', 1e-3)
config.setdefault('betal', 0.9)
config.setdefault('beta2', 0.999)
config.setdefault('epsilon', 1e-8)
config.setdefault('v', np.zeros_like(w))
config.setdefault('a', np.zeros_like(w))
config.setdefault('t', 0)

next_w = None

# ===== #
# YOUR CODE HERE:
# Implement Adam. Store the next value of w as next_w. You need
# to also store in config['a'] the moving average of the second
# moment gradients, and in config['v'] the moving average of the
# first moments. Finally, store in config['t'] the increasing tim
e.
# ===== #

betal = config['betal']
beta2 = config['beta2']
v = config['v']
a = config['a']

#time update
config['t'] = config['t'] + 1
t = config['t']

#first moment update (momentum-like)
config['v'] = betal*v + np.multiply(1-betal, dw)

#second moment update (gradient normalization)
config['a'] = beta2*a + (1-beta2)*np.multiply(dw, dw)

#bias correction in moments
v_bar = (1/(1-betal**t))*config['v']
a_bar = (1/(1-beta2**t))*config['a']

```

```
#gradient
next_w = w - np.multiply(config['learning_rate']/(np.sqrt(a_bar)+config['epsilon']), v_bar)

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

return next_w, config
```