# 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 using 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]:    1  ## Import and setups
           2
           3  import time
           4  import numpy as np
           5  import matplotlib.pyplot as plt
           6  from nndl.fc_net import *
           7  from cs231n.data_utils import get_CIFAR10_data
           8  from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
           9  from cs231n.solver import Solver
          10
          11  %matplotlib inline
          12  plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
          13  plt.rcParams['image.interpolation'] = 'nearest'
          14  plt.rcParams['image.cmap'] = 'gray'
          15
          16  # for auto-reloading external modules
          17  # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
          18  %load_ext autoreload
          19  %autoreload 2
          20
          21  def rel_error(x, y):
          22    """ returns relative error """
          23    return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

```
In [2]:    1  # Load the (preprocessed) CIFAR10 data.
           2
           3  data = get_CIFAR10_data()
           4  for k in data.keys():
           5    print('{}: {} '.format(k, data[k].shape))
```

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

## 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 [5]:   1  from nndl.layer_tests import *
          2
          3  affine_forward_test(); print('\n')
          4  affine_backward_test(); print('\n')
          5  relu_forward_test(); print('\n')
          6  relu_backward_test(); print('\n')
          7  affine_relu_test(); print('\n')
          8  fc_net_test()
```

```
If affine_forward function is working, difference should be less than 1e-9:
difference: 9.769849468192957e-10


If affine_backward is working, error should be less than 1e-9::
dx error: 3.666950908980358e-10
dw error: 8.201783430817534e-10
db error: 1.6263680637581037e-11


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: 3.2755836707832798e-12


If affine_relu_forward and affine_relu_backward are working, error should be less than 1e-9::
dx error: 5.543539459666791e-11
dw error: 2.363803633156294e-10
db error: 7.826635255953652e-12


Running check with reg = 0
Initial loss: 2.308535257408285
W1 relative error: 7.372336841535986e-07
W2 relative error: 3.66674068452486e-07
W3 relative error: 1.343554782509405e-07
b1 relative error: 1.4252530105730905e-08
b2 relative error: 5.616642824420512e-09
b3 relative error: 1.1821438210204362e-10
Running check with reg = 3.14
Initial loss: 7.3016750389748655
W1 relative error: 0.004730420797436833
W2 relative error: 0.33008883827369184
W3 relative error: 5.44330801643248e-08
b1 relative error: 1.742017280608633e-08
b2 relative error: 0.9602337605136528
b3 relative error: 4.064532921805321e-10
```

# Training a larger model

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

## SGD + momentum

In the following section, implement SGD with momentum. Read the `nndl/optim.py` API, which is provided by CS231n, and be sure you understand it. After, implement `sgd_momentum` in `nndl/optim.py`. Test your implementation of `sgd_momentum` by running the cell below.

```
In [6]:   1  from nndl.optim import sgd_momentum
          2
          3  N, D = 4, 5
          4  w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
          5  dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
          6  v = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)
          7
          8  config = {'learning_rate': 1e-3, 'velocity': v}
          9  next_w, _ = sgd_momentum(w, dw, config=config)
         10
         11  expected_next_w = np.asarray([
         12    [ 0.1406,      0.20738947,  0.27417895,  0.34096842,  0.40775789],
         13    [ 0.47454737,  0.54133684,  0.60812632,  0.67491579,  0.74170526],
         14    [ 0.80849474,  0.87528421,  0.94207368,  1.00886316,  1.07565263],
         15    [ 1.14244211,  1.20923158,  1.27602105,  1.34281053,  1.4096     ]])
         16  expected_velocity = np.asarray([
         17    [ 0.5406,      0.55475789, 0.56891579, 0.58307368,  0.59723158],
         18    [ 0.61138947,  0.62554737, 0.63970526, 0.65386316,  0.66802105],
         19    [ 0.68217895,  0.69633684, 0.71049474, 0.72465263,  0.73881053],
         20    [ 0.75296842,  0.76712632, 0.78128421, 0.79544211,  0.8096     ]])
         21
         22  print('next_w error: {}'.format(rel_error(next_w, expected_next_w)))
         23  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 [7]:   1  from nndl.optim import sgd_nesterov_momentum
          2
          3  N, D = 4, 5
          4  w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
          5  dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
          6  v = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)
          7
          8  config = {'learning_rate': 1e-3, 'velocity': v}
          9  next_w, _ = sgd_nesterov_momentum(w, dw, config=config)
         10
         11  expected_next_w = np.asarray([
         12    [0.08714,      0.15246105,  0.21778211,  0.28310316,  0.34842421],
         13    [0.41374526,   0.47906632,  0.54438737,  0.60970842,  0.67502947],
         14    [0.74035053,   0.80567158,  0.87099263,  0.93631368,  1.00163474],
         15    [1.06695579,   1.13227684,  1.19759789,  1.26291895,  1.32824    ]])
         16  expected_velocity = np.asarray([
         17    [ 0.5406,      0.55475789, 0.56891579, 0.58307368,  0.59723158],
         18    [ 0.61138947,  0.62554737, 0.63970526, 0.65386316,  0.66802105],
         19    [ 0.68217895,  0.69633684, 0.71049474, 0.72465263,  0.73881053],
         20    [ 0.75296842,  0.76712632, 0.78128421, 0.79544211,  0.8096     ]])
         21
         22  print('next_w error: {}'.format(rel_error(next_w, expected_next_w)))
         23  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 [8]:    1  num_train = 4000
           2  small_data = {
           3    'X_train': data['X_train'][:num_train],
           4    'y_train': data['y_train'][:num_train],
           5    'X_val': data['X_val'],
           6    'y_val': data['y_val'],
           7  }
           8
           9  solvers = {}
          10
          11  for update_rule in ['sgd', 'sgd_momentum', 'sgd_nesterov_momentum']:
          12    print('Optimizing with {}'.format(update_rule))
          13    model = FullyConnectedNet([100, 100, 100, 100, 100], weight_scale=5e-2)
          14
          15    solver = Solver(model, small_data,
          16                    num_epochs=5, batch_size=100,
          17                    update_rule=update_rule,
          18                    optim_config={
          19                      'learning_rate': 1e-2,
          20                    },
          21                    verbose=False)
          22    solvers[update_rule] = solver
          23    solver.train()
          24    print
          25
          26  plt.subplot(3, 1, 1)
          27  plt.title('Training loss')
          28  plt.xlabel('Iteration')
          29
          30  plt.subplot(3, 1, 2)
          31  plt.title('Training accuracy')
          32  plt.xlabel('Epoch')
          33
          34  plt.subplot(3, 1, 3)
          35  plt.title('Validation accuracy')
          36  plt.xlabel('Epoch')
          37
          38  for update_rule, solver in solvers.items():
          39    plt.subplot(3, 1, 1)
          40    plt.plot(solver.loss_history, 'o', label=update_rule)
          41
          42    plt.subplot(3, 1, 2)
          43    plt.plot(solver.train_acc_history, '-o', label=update_rule)
          44
          45    plt.subplot(3, 1, 3)
          46    plt.plot(solver.val_acc_history, '-o', label=update_rule)
          47
          48  for i in [1, 2, 3]:
          49    plt.subplot(3, 1, i)
          50    plt.legend(loc='upper center', ncol=4)
          51  plt.gcf().set_size_inches(15, 15)
          52  plt.show()
```
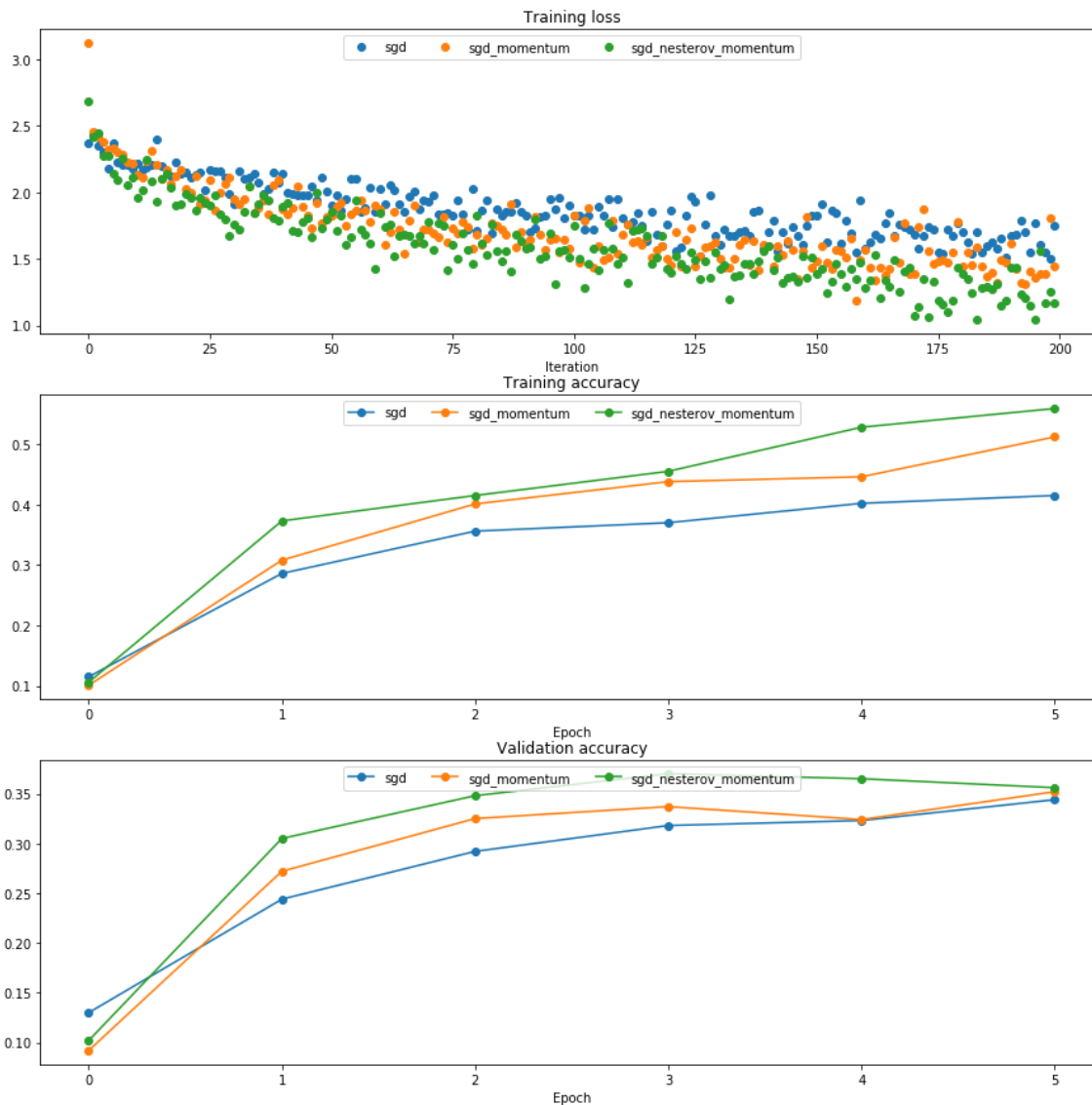
```
Optimizing with sgd
Optimizing with sgd_momentum
Optimizing with sgd_nesterov_momentum

/Users/hannah_wang/anaconda3/lib/python3.6/site-packages/matplotlib/figure.py:98: 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.
  "Adding an axes using the same arguments as a previous axes "
```

## 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 [9]:    1  from nndl.optim import rmsprop
           2
           3  N, D = 4, 5
           4  w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
           5  dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
           6  a = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)
           7
           8  config = {'learning_rate': 1e-2, 'a': a}
           9  next_w, _ = rmsprop(w, dw, config=config)
          10
          11  expected_next_w = np.asarray([
          12    [-0.39223849, -0.34037513, -0.28849239, -0.23659121, -0.18467247],
          13    [-0.132737,   -0.08078555, -0.02881884,  0.02316247,  0.07515774],
          14    [ 0.12716641,  0.17918792,  0.23122175,  0.28326742,  0.33532447],
          15    [ 0.38739248,  0.43947102,  0.49155973,  0.54365823,  0.59576619]])
          16  expected_cache = np.asarray([
          17    [ 0.5976,      0.6126277,   0.6277108,   0.64284931,  0.65804321],
          18    [ 0.67329252,  0.68859723,  0.70395734,  0.71937285,  0.73484377],
          19    [ 0.75037008,  0.7659518,   0.78158892,  0.79728144,  0.81302936],
          20    [ 0.82883269,  0.84469141,  0.86060554,  0.87657507,  0.8926    ]])
          21
          22  print('next_w error: {}'.format(rel_error(expected_next_w, next_w)))
          23  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.

```python
In [10]:    1  # Test Adam implementation; you should see errors around 1e-7 or less
            2  from nndl.optim import adam
            3
            4  N, D = 4, 5
            5  w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
            6  dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
            7  v = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)
            8  a = np.linspace(0.7, 0.5, num=N*D).reshape(N, D)
            9
           10  config = {'learning_rate': 1e-2, 'v': v, 'a': a, 't': 5}
           11  next_w, _ = adam(w, dw, config=config)
           12
           13  expected_next_w = np.asarray([
           14    [-0.40094747, -0.34836187, -0.29577703, -0.24319299, -0.19060977],
           15    [-0.1380274,  -0.08544591, -0.03286534,  0.01971428,  0.0722929],
           16    [ 0.1248705,   0.17744702,  0.23002243,  0.28259667,  0.33516969],
           17    [ 0.38774145,  0.44031188,  0.49288093,  0.54544852,  0.59801459]])
           18  expected_a = np.asarray([
           19    [ 0.69966,     0.68908382,  0.67851319,  0.66794809,  0.65738853,],
           20    [ 0.64683452,  0.63628604,  0.6257431,   0.61520571,  0.60467385,],
           21    [ 0.59414753,  0.58362676,  0.57311152,  0.56260183,  0.55209767,],
           22    [ 0.54159906,  0.53110598,  0.52061845,  0.51013645,  0.49966,    ]])
           23  expected_v = np.asarray([
           24    [ 0.48,        0.49947368,  0.51894737,  0.53842105,  0.55789474],
           25    [ 0.57736842,  0.59684211,  0.61631579,  0.63578947,  0.65526316],
           26    [ 0.67473684,  0.69421053,  0.71368421,  0.73315789,  0.75263158],
           27    [ 0.77210526,  0.79157895,  0.81105263,  0.83052632,  0.85      ]])
           28
           29  print('next_w error: {}'.format(rel_error(expected_next_w, next_w)))
           30  print('a error: {}'.format(rel_error(expected_a, config['a'])))
           31  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.
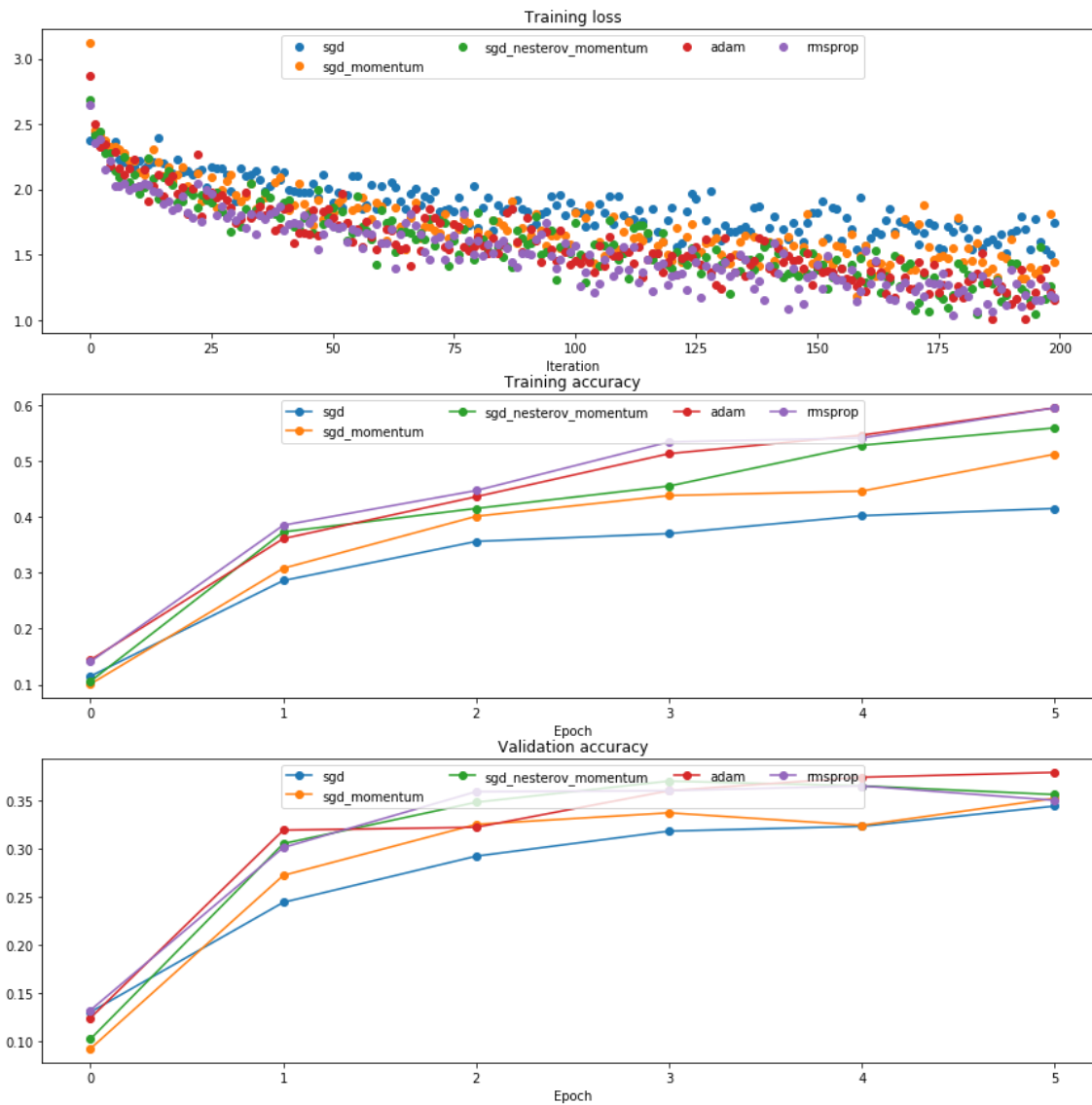
```
In [15]:   1  learning_rates = {'rmsprop': 2e-4, 'adam': 1e-3}
           2
           3  for update_rule in ['adam', 'rmsprop']:
           4    print('Optimizing with {}'.format(update_rule))
           5    model = FullyConnectedNet([100, 100, 100, 100, 100], weight_scale=5e-2)
           6
           7    solver = Solver(model, small_data,
           8                    num_epochs=5, batch_size=100,
           9                    update_rule=update_rule,
          10                    optim_config={
          11                       'learning_rate': learning_rates[update_rule]
          12                    },
          13                    verbose=False)
          14    solvers[update_rule] = solver
          15    solver.train()
          16    print
          17
          18  plt.subplot(3, 1, 1)
          19  plt.title('Training loss')
          20  plt.xlabel('Iteration')
          21
          22  plt.subplot(3, 1, 2)
          23  plt.title('Training accuracy')
          24  plt.xlabel('Epoch')
          25
          26  plt.subplot(3, 1, 3)
          27  plt.title('Validation accuracy')
          28  plt.xlabel('Epoch')
          29
          30  for update_rule, solver in solvers.items():
          31    plt.subplot(3, 1, 1)
          32    plt.plot(solver.loss_history, 'o', label=update_rule)
          33
          34    plt.subplot(3, 1, 2)
          35    plt.plot(solver.train_acc_history, '-o', label=update_rule)
          36
          37    plt.subplot(3, 1, 3)
          38    plt.plot(solver.val_acc_history, '-o', label=update_rule)
          39
          40  for i in [1, 2, 3]:
          41    plt.subplot(3, 1, i)
          42    plt.legend(loc='upper center', ncol=4)
          43  plt.gcf().set_size_inches(15, 15)
          44  plt.show()
```

```
Optimizing with adam
Optimizing with rmsprop

/Users/hannah_wang/anaconda3/lib/python3.6/site-packages/matplotlib/figure.py:98: 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.
  "Adding an axes using the same arguments as a previous axes "
```

## 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 [16]:    1  optimizer = 'adam'
            2  best_model = None
            3
            4  layer_dims = [500, 500, 500]
            5  weight_scale = 0.01
            6  learning_rate = 1e-3
            7  lr_decay = 0.9
            8
            9  model = FullyConnectedNet(layer_dims, weight_scale=weight_scale,
           10                            use_batchnorm=True)
           11
           12  solver = Solver(model, data,
           13                  num_epochs=10, batch_size=100,
           14                  update_rule=optimizer,
           15                  optim_config={
           16                     'learning_rate': learning_rate,
           17                  },
           18                  lr_decay=lr_decay,
           19                  verbose=True, print_every=50)
           20  solver.train()
```

```
(Iteration 1 / 4900) loss: 2.317612
(Epoch 0 / 10) train acc: 0.214000; val_acc: 0.222000
(Iteration 51 / 4900) loss: 1.850729
(Iteration 101 / 4900) loss: 1.529407
(Iteration 151 / 4900) loss: 1.734173
(Iteration 201 / 4900) loss: 1.648241
(Iteration 251 / 4900) loss: 1.607426
(Iteration 301 / 4900) loss: 1.404215
(Iteration 351 / 4900) loss: 1.707623
(Iteration 401 / 4900) loss: 1.347259
(Iteration 451 / 4900) loss: 1.336099
(Epoch 1 / 10) train acc: 0.470000; val_acc: 0.456000
(Iteration 501 / 4900) loss: 1.372245
(Iteration 551 / 4900) loss: 1.223931
(Iteration 601 / 4900) loss: 1.269121
(Iteration 651 / 4900) loss: 1.391613
(Iteration 701 / 4900) loss: 1.319634
(Iteration 751 / 4900) loss: 1.299265
(Iteration 801 / 4900) loss: 1.345726
(Iteration 851 / 4900) loss: 1.238437
(Iteration 901 / 4900) loss: 1.301950
(Iteration 951 / 4900) loss: 1.184255
(Epoch 2 / 10) train acc: 0.553000; val_acc: 0.519000
(Iteration 1001 / 4900) loss: 1.447854
(Iteration 1051 / 4900) loss: 1.135678
(Iteration 1101 / 4900) loss: 1.067297
(Iteration 1151 / 4900) loss: 1.279670
(Iteration 1201 / 4900) loss: 1.213916
(Iteration 1251 / 4900) loss: 1.283601
(Iteration 1301 / 4900) loss: 1.303041
(Iteration 1351 / 4900) loss: 1.313844
(Iteration 1401 / 4900) loss: 1.038823
(Iteration 1451 / 4900) loss: 1.161419
(Epoch 3 / 10) train acc: 0.578000; val_acc: 0.524000
(Iteration 1501 / 4900) loss: 1.080736
(Iteration 1551 / 4900) loss: 1.033311
(Iteration 1601 / 4900) loss: 1.256431
(Iteration 1651 / 4900) loss: 1.165053
(Iteration 1701 / 4900) loss: 1.240109
(Iteration 1751 / 4900) loss: 0.956544
(Iteration 1801 / 4900) loss: 1.123703
(Iteration 1851 / 4900) loss: 1.019171
(Iteration 1901 / 4900) loss: 1.116145
(Iteration 1951 / 4900) loss: 1.156048
(Epoch 4 / 10) train acc: 0.591000; val_acc: 0.544000
(Iteration 2001 / 4900) loss: 1.324829
(Iteration 2051 / 4900) loss: 1.150866
(Iteration 2101 / 4900) loss: 0.935841
(Iteration 2151 / 4900) loss: 0.957888
(Iteration 2201 / 4900) loss: 0.936134
(Iteration 2251 / 4900) loss: 0.978844
(Iteration 2301 / 4900) loss: 0.941161
(Iteration 2351 / 4900) loss: 0.915056
(Iteration 2401 / 4900) loss: 1.104329
(Epoch 5 / 10) train acc: 0.672000; val_acc: 0.532000
(Iteration 2451 / 4900) loss: 0.710098
(Iteration 2501 / 4900) loss: 0.913061
(Iteration 2551 / 4900) loss: 1.135810
(Iteration 2601 / 4900) loss: 0.876673
(Iteration 2651 / 4900) loss: 0.989470
(Iteration 2701 / 4900) loss: 0.990899
(Iteration 2751 / 4900) loss: 0.909822
(Iteration 2801 / 4900) loss: 0.756945
(Iteration 2851 / 4900) loss: 0.910173
(Iteration 2901 / 4900) loss: 0.725418
(Epoch 6 / 10) train acc: 0.692000; val_acc: 0.535000
(Iteration 2951 / 4900) loss: 0.615575
```

```
(Iteration 3001 / 4900) loss: 0.878288
(Iteration 3051 / 4900) loss: 1.034637
(Iteration 3101 / 4900) loss: 0.660970
(Iteration 3151 / 4900) loss: 0.891185
(Iteration 3201 / 4900) loss: 0.765091
(Iteration 3251 / 4900) loss: 0.859637
(Iteration 3301 / 4900) loss: 0.739107
(Iteration 3351 / 4900) loss: 0.663448
(Iteration 3401 / 4900) loss: 0.911083
(Epoch 7 / 10) train acc: 0.711000; val_acc: 0.569000
(Iteration 3451 / 4900) loss: 0.755873
(Iteration 3501 / 4900) loss: 0.817034
(Iteration 3551 / 4900) loss: 0.873263
(Iteration 3601 / 4900) loss: 0.684989
(Iteration 3651 / 4900) loss: 0.784117
(Iteration 3701 / 4900) loss: 0.709356
(Iteration 3751 / 4900) loss: 0.724685
(Iteration 3801 / 4900) loss: 0.639091
(Iteration 3851 / 4900) loss: 0.738882
(Iteration 3901 / 4900) loss: 0.473726
(Epoch 8 / 10) train acc: 0.754000; val_acc: 0.554000
(Iteration 3951 / 4900) loss: 0.587603
(Iteration 4001 / 4900) loss: 0.856407
(Iteration 4051 / 4900) loss: 0.672353
(Iteration 4101 / 4900) loss: 0.781707
(Iteration 4151 / 4900) loss: 0.604926
(Iteration 4201 / 4900) loss: 0.596246
(Iteration 4251 / 4900) loss: 0.630911
(Iteration 4301 / 4900) loss: 0.534743
(Iteration 4351 / 4900) loss: 0.716217
(Iteration 4401 / 4900) loss: 0.587710
(Epoch 9 / 10) train acc: 0.802000; val_acc: 0.569000
(Iteration 4451 / 4900) loss: 0.556148
(Iteration 4501 / 4900) loss: 0.547765
(Iteration 4551 / 4900) loss: 0.704546
(Iteration 4601 / 4900) loss: 0.596978
(Iteration 4651 / 4900) loss: 0.642664
(Iteration 4701 / 4900) loss: 0.584378
(Iteration 4751 / 4900) loss: 0.542869
(Iteration 4801 / 4900) loss: 0.612463
(Iteration 4851 / 4900) loss: 0.613176
(Epoch 10 / 10) train acc: 0.818000; val_acc: 0.557000
```

In [17]:
```python
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: 0.578
Test set accuracy: 0.569
```

In [ ]:
```
1
```

# Batch Normalization

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

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

```
In [1]:   1  ## Import and setups
          2
          3  import time
          4  import numpy as np
          5  import matplotlib.pyplot as plt
          6  from nndl.fc_net import *
          7  from nndl.layers import *
          8  from cs231n.data_utils import get_CIFAR10_data
          9  from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
         10  from cs231n.solver import Solver
         11
         12  %matplotlib inline
         13  plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
         14  plt.rcParams['image.interpolation'] = 'nearest'
         15  plt.rcParams['image.cmap'] = 'gray'
         16
         17  # for auto-reloading external modules
         18  # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
         19  %load_ext autoreload
         20  %autoreload 2
         21
         22  def rel_error(x, y):
         23    """ returns relative error """
         24    return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

```
In [2]:   1  # Load the (preprocessed) CIFAR10 data.
          2
          3  data = get_CIFAR10_data()
          4  for k in data.keys():
          5    print('{}: {} '.format(k, data[k].shape))
```

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

## 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 [4]:    1  # Check the training-time forward pass by checking means and variances
           2  # of features both before and after batch normalization
           3
           4  # Simulate the forward pass for a two-layer network
           5  N, D1, D2, D3 = 200, 50, 60, 3
           6  X = np.random.randn(N, D1)
           7  W1 = np.random.randn(D1, D2)
           8  W2 = np.random.randn(D2, D3)
           9  a = np.maximum(0, X.dot(W1)).dot(W2)
          10
          11  print('Before batch normalization:')
          12  print('  means: ', a.mean(axis=0))
          13  print('  stds: ', a.std(axis=0))
          14
          15  # Means should be close to zero and stds close to one
          16  print('After batch normalization (gamma=1, beta=0)')
          17  a_norm, _ = batchnorm_forward(a, np.ones(D3), np.zeros(D3), {'mode': 'train'})
          18  print('  mean: ', a_norm.mean(axis=0))
          19  print('  std: ', a_norm.std(axis=0))
          20
          21  # Now means should be close to beta and stds close to gamma
          22  gamma = np.asarray([1.0, 2.0, 3.0])
          23  beta = np.asarray([11.0, 12.0, 13.0])
          24  a_norm, _ = batchnorm_forward(a, gamma, beta, {'mode': 'train'})
          25  print('After batch normalization (nontrivial gamma, beta)')
          26  print('  means: ', a_norm.mean(axis=0))
          27  print('  stds: ', a_norm.std(axis=0))
```

```
Before batch normalization:
  means:  [-9.27589484 21.75456393 17.54660778]
  stds:  [29.33259377 40.84649232 29.16922165]
After batch normalization (gamma=1, beta=0)
  mean:  [-5.10702591e-17 -5.41927614e-17 -1.74860126e-16]
  std:  [0.99999999 1.          0.99999999]
After batch normalization (nontrivial gamma, beta)
  means:  [11. 12. 13.]
  stds:  [0.99999999 1.99999999 2.99999998]
```

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

```
In [6]:    1  # Check the test-time forward pass by running the training-time
           2  # forward pass many times to warm up the running averages, and then
           3  # checking the means and variances of activations after a test-time
           4  # forward pass.
           5
           6  N, D1, D2, D3 = 200, 50, 60, 3
           7  W1 = np.random.randn(D1, D2)
           8  W2 = np.random.randn(D2, D3)
           9
          10  bn_param = {'mode': 'train'}
          11  gamma = np.ones(D3)
          12  beta = np.zeros(D3)
          13  for t in np.arange(50):
          14    X = np.random.randn(N, D1)
          15    a = np.maximum(0, X.dot(W1)).dot(W2)
          16    batchnorm_forward(a, gamma, beta, bn_param)
          17  bn_param['mode'] = 'test'
          18  X = np.random.randn(N, D1)
          19  a = np.maximum(0, X.dot(W1)).dot(W2)
          20  a_norm, _ = batchnorm_forward(a, gamma, beta, bn_param)
          21
          22  # Means should be close to zero and stds close to one, but will be
          23  # noisier than training-time forward passes.
          24  print('After batch normalization (test-time):')
          25  print('  means: ', a_norm.mean(axis=0))
          26  print('  stds: ', a_norm.std(axis=0))
```

```
After batch normalization (test-time):
  means:  [-0.04310833  0.03218559 -0.01450178]
  stds:  [0.98114234 1.08661348 1.01088546]
```

## 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 [7]:   1  # Gradient check batchnorm backward pass
          2
          3  N, D = 4, 5
          4  x = 5 * np.random.randn(N, D) + 12
          5  gamma = np.random.randn(D)
          6  beta = np.random.randn(D)
          7  dout = np.random.randn(N, D)
          8
          9  bn_param = {'mode': 'train'}
         10  fx = lambda x: batchnorm_forward(x, gamma, beta, bn_param)[0]
         11  fg = lambda a: batchnorm_forward(x, gamma, beta, bn_param)[0]
         12  fb = lambda b: batchnorm_forward(x, gamma, beta, bn_param)[0]
         13
         14  dx_num = eval_numerical_gradient_array(fx, x, dout)
         15  da_num = eval_numerical_gradient_array(fg, gamma, dout)
         16  db_num = eval_numerical_gradient_array(fb, beta, dout)
         17
         18  _, cache = batchnorm_forward(x, gamma, beta, bn_param)
         19  dx, dgamma, dbeta = batchnorm_backward(dout, cache)
         20  print('dx error: ', rel_error(dx_num, dx))
         21  print('dgamma error: ', rel_error(da_num, dgamma))
         22  print('dbeta error: ', rel_error(db_num, dbeta))
```

```
dx error:  1.3407317993977258e-09
dgamma error:  3.2369133023643904e-12
dbeta error:  3.275420131545755e-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 [8]:    1  N, D, H1, H2, C = 2, 15, 20, 30, 10
           2  X = np.random.randn(N, D)
           3  y = np.random.randint(C, size=(N,))
           4
           5  for reg in [0, 3.14]:
           6    print('Running check with reg = ', reg)
           7    model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
           8                              reg=reg, weight_scale=5e-2, dtype=np.float64,
           9                              use_batchnorm=True)
          10
          11    loss, grads = model.loss(X, y)
          12    print('Initial loss: ', loss)
          13
          14    for name in sorted(grads):
          15      f = lambda _: model.loss(X, y)[0]
          16      grad_num = eval_numerical_gradient(f, model.params[name], verbose=False, h=1e-5)
          17      print('{} relative error: {}'.format(name, rel_error(grad_num, grads[name])))
          18    if reg == 0: print('\n')
```

```
Running check with reg =  0
Initial loss:  2.116495647317618
W1 relative error: 0.003059418529722701
W2 relative error: 3.848515927934831e-06
W3 relative error: 4.093183258233747e-10
b1 relative error: 4.336808689942018e-11
b2 relative error: 1.1102230246251565e-08
b3 relative error: 1.237819602511532e-10
beta1 relative error: 3.265429352574524e-07
beta2 relative error: 7.424791031728775e-09
gamma1 relative error: 3.984791350136109e-07
gamma2 relative error: 3.9423807659745554e-09


Running check with reg =  3.14
Initial loss:  6.945176771933103
W1 relative error: 1.602456195923461e-05
W2 relative error: 2.03407132663365e-06
W3 relative error: 1.4177185950422472e-08
b1 relative error: 1.1102230246251565e-08
b2 relative error: 2.7755575615628914e-09
b3 relative error: 1.47406506027144e-10
beta1 relative error: 4.766663710604228e-09
beta2 relative error: 2.5601940267967956e-08
gamma1 relative error: 4.697072145384006e-09
gamma2 relative error: 5.502280624431884e-09
```

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

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

```
In [9]:    1  # Try training a very deep net with batchnorm
           2  hidden_dims = [100, 100, 100, 100, 100]
           3
           4  num_train = 1000
           5  small_data = {
           6    'X_train': data['X_train'][:num_train],
           7    'y_train': data['y_train'][:num_train],
           8    'X_val': data['X_val'],
           9    'y_val': data['y_val'],
          10  }
          11
          12  weight_scale = 2e-2
          13  bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=True)
          14  model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=False)
          15
          16  bn_solver = Solver(bn_model, small_data,
          17                  num_epochs=10, batch_size=50,
          18                  update_rule='adam',
          19                  optim_config={
          20                     'learning_rate': 1e-3,
          21                  },
          22                  verbose=True, print_every=200)
          23  bn_solver.train()
          24
          25  solver = Solver(model, small_data,
          26                  num_epochs=10, batch_size=50,
          27                  update_rule='adam',
          28                  optim_config={
          29                     'learning_rate': 1e-3,
          30                  },
          31                  verbose=True, print_every=200)
          32  solver.train()
```

```
(Iteration 1 / 200) loss: 2.306698
(Epoch 0 / 10) train acc: 0.139000; val_acc: 0.118000
(Epoch 1 / 10) train acc: 0.340000; val_acc: 0.279000
(Epoch 2 / 10) train acc: 0.424000; val_acc: 0.303000
(Epoch 3 / 10) train acc: 0.509000; val_acc: 0.311000
(Epoch 4 / 10) train acc: 0.572000; val_acc: 0.323000
(Epoch 5 / 10) train acc: 0.619000; val_acc: 0.318000
(Epoch 6 / 10) train acc: 0.653000; val_acc: 0.333000
(Epoch 7 / 10) train acc: 0.732000; val_acc: 0.345000
(Epoch 8 / 10) train acc: 0.754000; val_acc: 0.329000
(Epoch 9 / 10) train acc: 0.800000; val_acc: 0.329000
(Epoch 10 / 10) train acc: 0.807000; val_acc: 0.313000
(Iteration 1 / 200) loss: 2.302677
(Epoch 0 / 10) train acc: 0.165000; val_acc: 0.155000
(Epoch 1 / 10) train acc: 0.241000; val_acc: 0.229000
(Epoch 2 / 10) train acc: 0.290000; val_acc: 0.255000
(Epoch 3 / 10) train acc: 0.343000; val_acc: 0.303000
(Epoch 4 / 10) train acc: 0.361000; val_acc: 0.275000
(Epoch 5 / 10) train acc: 0.415000; val_acc: 0.280000
(Epoch 6 / 10) train acc: 0.488000; val_acc: 0.317000
(Epoch 7 / 10) train acc: 0.447000; val_acc: 0.286000
(Epoch 8 / 10) train acc: 0.524000; val_acc: 0.298000
(Epoch 9 / 10) train acc: 0.550000; val_acc: 0.334000
(Epoch 10 / 10) train acc: 0.658000; val_acc: 0.338000
```
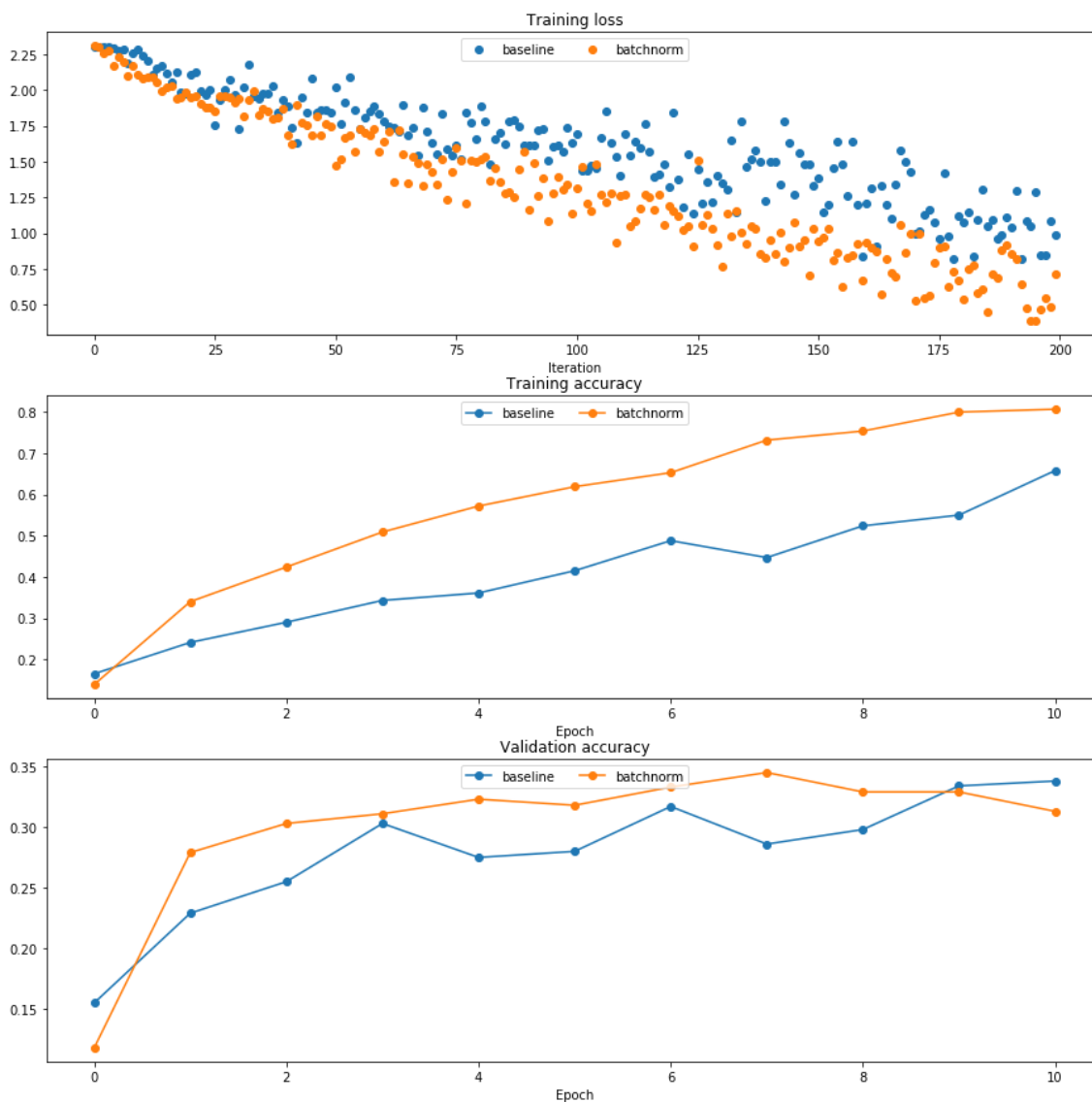
```
In [10]:   1  plt.subplot(3, 1, 1)
           2  plt.title('Training loss')
           3  plt.xlabel('Iteration')
           4
           5  plt.subplot(3, 1, 2)
           6  plt.title('Training accuracy')
           7  plt.xlabel('Epoch')
           8
           9  plt.subplot(3, 1, 3)
          10  plt.title('Validation accuracy')
          11  plt.xlabel('Epoch')
          12
          13  plt.subplot(3, 1, 1)
          14  plt.plot(solver.loss_history, 'o', label='baseline')
          15  plt.plot(bn_solver.loss_history, 'o', label='batchnorm')
          16
          17  plt.subplot(3, 1, 2)
          18  plt.plot(solver.train_acc_history, '-o', label='baseline')
          19  plt.plot(bn_solver.train_acc_history, '-o', label='batchnorm')
          20
          21  plt.subplot(3, 1, 3)
          22  plt.plot(solver.val_acc_history, '-o', label='baseline')
          23  plt.plot(bn_solver.val_acc_history, '-o', label='batchnorm')
          24
          25  for i in [1, 2, 3]:
          26      plt.subplot(3, 1, i)
          27      plt.legend(loc='upper center', ncol=4)
          28  plt.gcf().set_size_inches(15, 15)
          29  plt.show()
```

/Users/hannah_wang/anaconda3/lib/python3.6/site-packages/matplotlib/figure.py:98: 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.
  "Adding an axes using the same arguments as a previous axes "

## 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 [11]:
 1  # Try training a very deep net with batchnorm
 2  hidden_dims = [50, 50, 50, 50, 50, 50, 50]
 3
 4  num_train = 1000
 5  small_data = {
 6    'X_train': data['X_train'][:num_train],
 7    'y_train': data['y_train'][:num_train],
 8    'X_val': data['X_val'],
 9    'y_val': data['y_val'],
10  }
11
12  bn_solvers = {}
13  solvers = {}
14  weight_scales = np.logspace(-4, 0, num=20)
15  for i, weight_scale in enumerate(weight_scales):
16    print('Running weight scale {} / {}'.format(i + 1, len(weight_scales)))
17    bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=True)
18    model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=False)
19
20    bn_solver = Solver(bn_model, small_data,
21                    num_epochs=10, batch_size=50,
22                    update_rule='adam',
23                    optim_config={
24                        'learning_rate': 1e-3,
25                    },
26                    verbose=False, print_every=200)
27    bn_solver.train()
28    bn_solvers[weight_scale] = bn_solver
29
30    solver = Solver(model, small_data,
31                    num_epochs=10, batch_size=50,
32                    update_rule='adam',
33                    optim_config={
34                        'learning_rate': 1e-3,
35                    },
36                    verbose=False, print_every=200)
37    solver.train()
38    solvers[weight_scale] = solver
```

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

/Users/hannah_wang/Desktop/hw4/code/nndl/layers.py:438: RuntimeWarning: divide by zero encountered in log
  loss = -np.sum(np.log(probs[np.arange(N), y])) / N

Running weight scale 17 / 20
Running weight scale 18 / 20
Running weight scale 19 / 20
Running weight scale 20 / 20
```
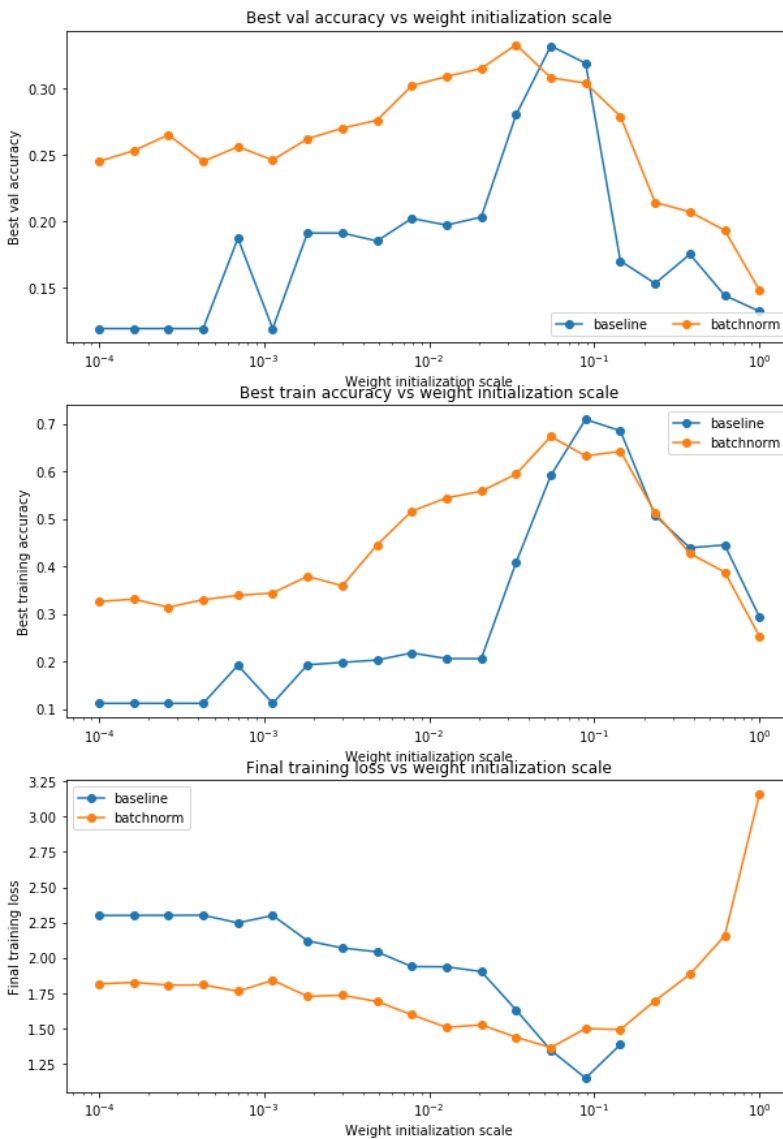
```
In [12]:   1  # Plot results of weight scale experiment
           2  best_train_accs, bn_best_train_accs = [], []
           3  best_val_accs, bn_best_val_accs = [], []
           4  final_train_loss, bn_final_train_loss = [], []
           5
           6  for ws in weight_scales:
           7    best_train_accs.append(max(solvers[ws].train_acc_history))
           8    bn_best_train_accs.append(max(bn_solvers[ws].train_acc_history))
           9
          10    best_val_accs.append(max(solvers[ws].val_acc_history))
          11    bn_best_val_accs.append(max(bn_solvers[ws].val_acc_history))
          12
          13    final_train_loss.append(np.mean(solvers[ws].loss_history[-100:]))
          14    bn_final_train_loss.append(np.mean(bn_solvers[ws].loss_history[-100:]))
          15
          16  plt.subplot(3, 1, 1)
          17  plt.title('Best val accuracy vs weight initialization scale')
          18  plt.xlabel('Weight initialization scale')
          19  plt.ylabel('Best val accuracy')
          20  plt.semilogx(weight_scales, best_val_accs, '-o', label='baseline')
          21  plt.semilogx(weight_scales, bn_best_val_accs, '-o', label='batchnorm')
          22  plt.legend(ncol=2, loc='lower right')
          23
          24  plt.subplot(3, 1, 2)
          25  plt.title('Best train accuracy vs weight initialization scale')
          26  plt.xlabel('Weight initialization scale')
          27  plt.ylabel('Best training accuracy')
          28  plt.semilogx(weight_scales, best_train_accs, '-o', label='baseline')
          29  plt.semilogx(weight_scales, bn_best_train_accs, '-o', label='batchnorm')
          30  plt.legend()
          31
          32  plt.subplot(3, 1, 3)
          33  plt.title('Final training loss vs weight initialization scale')
          34  plt.xlabel('Weight initialization scale')
          35  plt.ylabel('Final training loss')
          36  plt.semilogx(weight_scales, final_train_loss, '-o', label='baseline')
          37  plt.semilogx(weight_scales, bn_final_train_loss, '-o', label='batchnorm')
          38  plt.legend()
          39
          40  plt.gcf().set_size_inches(10, 15)
          41  plt.show()
```

Best val accuracy vs weight initialization scale

Best train accuracy vs weight initialization scale

Final training loss vs weight initialization scale

## Question:

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

## Answer:

From the figures above, on the one hand, batchnorm can achieve more accurate results. We can find that when weight initialization scale is smaller than 10^(-1), the accuracies of model with batchnorm are almost larger than those of model without batchnorm. And the loss of model with batchnorm are almost smaller than those of model without batchnorm. While when weight initialization scale is larger than 10^(-1), the model with batchnorm has poorer performance than model without batchnorm. On the other hand, we find that the model with batchnorm is much less sensitive to the weight initialization scale. The orange curves have wider stable ranges than blue curves.

As we know, batchnorm can be used as a regularizer to make results more stable. Thus, the model with batchnorm should be less affected by weight initialization scale. And the results we observed meet this theory and make sense.

In [ ]: | 1 |

# Dropout

In this notebook, you will implement dropout. Then we will ask you to train a network with batchnorm and dropout, and acheive 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 [2]:    1  ## Import and setups
           2
           3  import time
           4  import numpy as np
           5  import matplotlib.pyplot as plt
           6  from nndl.fc_net import *
           7  from nndl.layers import *
           8  from cs231n.data_utils import get_CIFAR10_data
           9  from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
          10  from cs231n.solver import Solver
          11
          12  %matplotlib inline
          13  plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
          14  plt.rcParams['image.interpolation'] = 'nearest'
          15  plt.rcParams['image.cmap'] = 'gray'
          16
          17  # for auto-reloading external modules
          18  # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
          19  %load_ext autoreload
          20  %autoreload 2
          21
          22  def rel_error(x, y):
          23    """ returns relative error """
          24    return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

```
In [3]:    1  # Load the (preprocessed) CIFAR10 data.
           2
           3  data = get_CIFAR10_data()
           4  for k in data.keys():
           5    print('{}: {} '.format(k, data[k].shape))
```

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

## 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 [4]:    1  x = np.random.randn(500, 500) + 10
           2
           3  for p in [0.3, 0.6, 0.75]:
           4      out, _ = dropout_forward(x, {'mode': 'train', 'p': p})
           5      out_test, _ = dropout_forward(x, {'mode': 'test', 'p': p})
           6
           7      print('Running tests with p = ', p)
           8      print('Mean of input: ', x.mean())
           9      print('Mean of train-time output: ', out.mean())
          10      print('Mean of test-time output: ', out_test.mean())
          11      print('Fraction of train-time output set to zero: ', (out == 0).mean())
          12      print('Fraction of test-time output set to zero: ', (out_test == 0).mean())
```

```
Running tests with p =  0.3
Mean of input:  9.997675409022476
Mean of train-time output:  9.982471546063854
Mean of test-time output:  9.997675409022476
Fraction of train-time output set to zero:  0.30116
Fraction of test-time output set to zero:  0.0
Running tests with p =  0.6
Mean of input:  9.997675409022476
Mean of train-time output:  9.953386552374035
Mean of test-time output:  9.997675409022476
Fraction of train-time output set to zero:  0.601796
Fraction of test-time output set to zero:  0.0
Running tests with p =  0.75
Mean of input:  9.997675409022476
Mean of train-time output:  9.96677999564878
Mean of test-time output:  9.997675409022476
Fraction of train-time output set to zero:  0.750764
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 [5]:    1  x = np.random.randn(10, 10) + 10
           2  dout = np.random.randn(*x.shape)
           3
           4  dropout_param = {'mode': 'train', 'p': 0.8, 'seed': 123}
           5  out, cache = dropout_forward(x, dropout_param)
           6  dx = dropout_backward(dout, cache)
           7  dx_num = eval_numerical_gradient_array(lambda xx: dropout_forward(xx, dropout_param)[0], x, dout)
           8
           9  print('dx relative error: ', rel_error(dx, dx_num))
```

```
dx relative error:  1.892905802538152e-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).

```
 1  N, D, H1, H2, C = 2, 15, 20, 30, 10
 2  X = np.random.randn(N, D)
 3  y = np.random.randint(C, size=(N,))
 4
 5  for dropout in [0, 0.25, 0.5]:
 6    print('Running check with dropout = ', dropout)
 7    model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
 8                               weight_scale=5e-2, dtype=np.float64,
 9                               dropout=dropout, seed=123)
10
11    loss, grads = model.loss(X, y)
12    print('Initial loss: ', loss)
13
14    for name in sorted(grads):
15      f = lambda _: model.loss(X, y)[0]
16      grad_num = eval_numerical_gradient(f, model.params[name], verbose=False, h=1e-5)
17      print('{} relative error: {}'.format(name, rel_error(grad_num, grads[name])))
18    print('\n')
```

```
Running check with dropout =  0
Initial loss:  2.3051948273987857
W1 relative error: 2.5272575344376073e-07
W2 relative error: 1.5034484929313676e-05
W3 relative error: 2.753446833630168e-07
b1 relative error: 2.936957476400148e-06
b2 relative error: 5.051339805546953e-08
b3 relative error: 1.1740467838205477e-10


Running check with dropout =  0.25
Initial loss:  2.3052077546540826
W1 relative error: 2.613846944812385e-07
W2 relative error: 5.022056536108928e-07
W3 relative error: 4.456316077044505e-08
b1 relative error: 7.39711723790801e-08
b2 relative error: 7.151678402730031e-10
b3 relative error: 1.003974732116764e-10


Running check with dropout =  0.5
Initial loss:  2.3035667586595423
W1 relative error: 1.1401257458777745e-06
W2 relative error: 1.847669681023635e-07
W3 relative error: 6.5966195253431734e-09
b1 relative error: 7.71639621892128e-08
b2 relative error: 1.1975910493629166e-09
b3 relative error: 1.4558471033827801e-10
```

## Dropout as a regularizer

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

```python
# 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: 2.298716
(Epoch 0 / 25) train acc: 0.132000; val_acc: 0.146000
(Epoch 1 / 25) train acc: 0.118000; val_acc: 0.131000
(Epoch 2 / 25) train acc: 0.220000; val_acc: 0.214000
(Epoch 3 / 25) train acc: 0.206000; val_acc: 0.180000
(Epoch 4 / 25) train acc: 0.220000; val_acc: 0.193000
(Epoch 5 / 25) train acc: 0.264000; val_acc: 0.229000
(Epoch 6 / 25) train acc: 0.268000; val_acc: 0.203000
(Epoch 7 / 25) train acc: 0.266000; val_acc: 0.212000
(Epoch 8 / 25) train acc: 0.282000; val_acc: 0.236000
(Epoch 9 / 25) train acc: 0.310000; val_acc: 0.255000
(Epoch 10 / 25) train acc: 0.320000; val_acc: 0.267000
(Epoch 11 / 25) train acc: 0.338000; val_acc: 0.273000
(Epoch 12 / 25) train acc: 0.346000; val_acc: 0.278000
(Epoch 13 / 25) train acc: 0.332000; val_acc: 0.279000
(Epoch 14 / 25) train acc: 0.328000; val_acc: 0.284000
(Epoch 15 / 25) train acc: 0.354000; val_acc: 0.271000
(Epoch 16 / 25) train acc: 0.386000; val_acc: 0.277000
(Epoch 17 / 25) train acc: 0.388000; val_acc: 0.297000
(Epoch 18 / 25) train acc: 0.402000; val_acc: 0.280000
(Epoch 19 / 25) train acc: 0.388000; val_acc: 0.274000
(Epoch 20 / 25) train acc: 0.386000; val_acc: 0.274000
(Iteration 101 / 125) loss: 1.919649
(Epoch 21 / 25) train acc: 0.402000; val_acc: 0.272000
(Epoch 22 / 25) train acc: 0.440000; val_acc: 0.286000
(Epoch 23 / 25) train acc: 0.458000; val_acc: 0.295000
(Epoch 24 / 25) train acc: 0.462000; val_acc: 0.311000
(Epoch 25 / 25) train acc: 0.466000; val_acc: 0.297000
```

```
1  # Plot train and validation accuracies of the two models
2
3  train_accs = []
4  val_accs = []
5  for dropout in dropout_choices:
6    solver = solvers[dropout]
7    train_accs.append(solver.train_acc_history[-1])
8    val_accs.append(solver.val_acc_history[-1])
9
10 plt.subplot(3, 1, 1)
11 for dropout in dropout_choices:
12   plt.plot(solvers[dropout].train_acc_history, 'o', label='%.2f dropout' % dropout)
13 plt.title('Train accuracy')
14 plt.xlabel('Epoch')
15 plt.ylabel('Accuracy')
16 plt.legend(ncol=2, loc='lower right')
17
18 plt.subplot(3, 1, 2)
19 for dropout in dropout_choices:
20   plt.plot(solvers[dropout].val_acc_history, 'o', label='%.2f dropout' % dropout)
21 plt.title('Val accuracy')
22 plt.xlabel('Epoch')
23 plt.ylabel('Accuracy')
24 plt.legend(ncol=2, loc='lower right')
25
26 plt.gcf().set_size_inches(15, 15)
27 plt.show()
```



## Question

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

## Answer:

Yes, dropout is performing regularization. From two figures above, we can see that in the second figure, the validation accuracies of model with and without dropout are similar. While in the first figure, the training accuracy of model without dropout is apparently larger than the model with dropout. It indicates that the dropout regularized the overfitting in training data.

## 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(floor((X - 32%)) / 28%, 1) where if you get 60% or higher validation accuracy, you get full points.

```python
# =============================================================== #
# YOUR CODE HERE:
#   Implement a FC-net that achieves at least 60% validation accuracy
#   on CIFAR-10.
# =============================================================== #

# set parameters
hidden_dims = [600, 600, 600, 600]
learning_rate = 2e-3
weight_scale = 0.01
lr_decay = 0.95
dropout = 0.5
update_rule = 'adam'

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

# solve
solver = Solver(model, data,
                num_epochs=80, batch_size=100,
                update_rule=update_rule,
                optim_config={
                    'learning_rate': learning_rate,
                },
                lr_decay=lr_decay,
                verbose=True, print_every=100)
solver.train()

# print out the validation accuracy
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'])))

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

```
(Iteration 1 / 39200) loss: 2.361412
(Epoch 0 / 80) train acc: 0.185000; val_acc: 0.197000
(Iteration 101 / 39200) loss: 1.853467
(Iteration 201 / 39200) loss: 1.860087
(Iteration 301 / 39200) loss: 1.786663
(Iteration 401 / 39200) loss: 1.698297
(Epoch 1 / 80) train acc: 0.442000; val_acc: 0.442000
(Iteration 501 / 39200) loss: 1.532157
(Iteration 601 / 39200) loss: 1.779091
(Iteration 701 / 39200) loss: 1.590496
(Iteration 801 / 39200) loss: 1.645538
(Iteration 901 / 39200) loss: 1.716053
(Epoch 2 / 80) train acc: 0.457000; val_acc: 0.471000
(Iteration 1001 / 39200) loss: 1.560481
(Iteration 1101 / 39200) loss: 1.402321
(Iteration 1201 / 39200) loss: 1.559676
(Iteration 1301 / 39200) loss: 1.514315
(Iteration 1401 / 39200) loss: 1.365991
(Epoch 3 / 80) train acc: 0.503000; val_acc: 0.500000
(Iteration 1501 / 39200) loss: 1.577194
(Iteration 1601 / 39200) loss: 1.380106
(Iteration 1701 / 39200) loss: 1.449796
(Iteration 1801 / 39200) loss: 1.444384
(Iteration 1901 / 39200) loss: 1.465490
(Epoch 4 / 80) train acc: 0.519000; val_acc: 0.518000
(Iteration 2001 / 39200) loss: 1.448264
(Iteration 2101 / 39200) loss: 1.407761
(Iteration 2201 / 39200) loss: 1.635408
(Iteration 2301 / 39200) loss: 1.491990
(Iteration 2401 / 39200) loss: 1.442680
(Epoch 5 / 80) train acc: 0.536000; val_acc: 0.526000
(Iteration 2501 / 39200) loss: 1.322390
(Iteration 2601 / 39200) loss: 1.230148
(Iteration 2701 / 39200) loss: 1.390268
(Iteration 2801 / 39200) loss: 1.354062
(Iteration 2901 / 39200) loss: 1.368472
(Epoch 6 / 80) train acc: 0.522000; val_acc: 0.520000
(Iteration 3001 / 39200) loss: 1.486219
(Iteration 3101 / 39200) loss: 1.269545
(Iteration 3201 / 39200) loss: 1.411731
(Iteration 3301 / 39200) loss: 1.544854
(Iteration 3401 / 39200) loss: 1.208551
(Epoch 7 / 80) train acc: 0.559000; val_acc: 0.538000
(Iteration 3501 / 39200) loss: 1.357756
(Iteration 3601 / 39200) loss: 1.460392
(Iteration 3701 / 39200) loss: 1.351261
(Iteration 3801 / 39200) loss: 1.317291
(Iteration 3901 / 39200) loss: 1.375946
```

```
(Epoch 8 / 80) train acc: 0.593000; val_acc: 0.523000
(Iteration 4001 / 39200) loss: 1.535556
(Iteration 4101 / 39200) loss: 1.370858
(Iteration 4201 / 39200) loss: 1.463170
(Iteration 4301 / 39200) loss: 1.335995
(Iteration 4401 / 39200) loss: 1.238724
(Epoch 9 / 80) train acc: 0.603000; val_acc: 0.543000
(Iteration 4501 / 39200) loss: 1.143931
(Iteration 4601 / 39200) loss: 1.250761
(Iteration 4701 / 39200) loss: 1.309563
(Iteration 4801 / 39200) loss: 1.253994
(Epoch 10 / 80) train acc: 0.601000; val_acc: 0.554000
(Iteration 4901 / 39200) loss: 1.295024
(Iteration 5001 / 39200) loss: 1.377004
(Iteration 5101 / 39200) loss: 1.171046
(Iteration 5201 / 39200) loss: 1.235158
(Iteration 5301 / 39200) loss: 1.346153
(Epoch 11 / 80) train acc: 0.610000; val_acc: 0.562000
(Iteration 5401 / 39200) loss: 1.265616
(Iteration 5501 / 39200) loss: 1.114638
(Iteration 5601 / 39200) loss: 1.352518
(Iteration 5701 / 39200) loss: 1.313555
(Iteration 5801 / 39200) loss: 1.357678
(Epoch 12 / 80) train acc: 0.609000; val_acc: 0.571000
(Iteration 5901 / 39200) loss: 1.306446
(Iteration 6001 / 39200) loss: 1.259089
(Iteration 6101 / 39200) loss: 1.183187
(Iteration 6201 / 39200) loss: 1.309481
(Iteration 6301 / 39200) loss: 1.226226
(Epoch 13 / 80) train acc: 0.604000; val_acc: 0.572000
(Iteration 6401 / 39200) loss: 1.293294
(Iteration 6501 / 39200) loss: 1.186847
(Iteration 6601 / 39200) loss: 1.171038
(Iteration 6701 / 39200) loss: 1.268189
(Iteration 6801 / 39200) loss: 1.128496
(Epoch 14 / 80) train acc: 0.622000; val_acc: 0.564000
(Iteration 6901 / 39200) loss: 1.195287
(Iteration 7001 / 39200) loss: 1.199769
(Iteration 7101 / 39200) loss: 1.265404
(Iteration 7201 / 39200) loss: 1.083662
(Iteration 7301 / 39200) loss: 1.171192
(Epoch 15 / 80) train acc: 0.655000; val_acc: 0.559000
(Iteration 7401 / 39200) loss: 1.132624
(Iteration 7501 / 39200) loss: 1.482053
(Iteration 7601 / 39200) loss: 1.218198
(Iteration 7701 / 39200) loss: 1.353706
(Iteration 7801 / 39200) loss: 1.156821
(Epoch 16 / 80) train acc: 0.653000; val_acc: 0.571000
(Iteration 7901 / 39200) loss: 1.110690
(Iteration 8001 / 39200) loss: 1.061102
(Iteration 8101 / 39200) loss: 1.280265
(Iteration 8201 / 39200) loss: 1.259501
(Iteration 8301 / 39200) loss: 1.400559
(Epoch 17 / 80) train acc: 0.663000; val_acc: 0.558000
(Iteration 8401 / 39200) loss: 1.351456
(Iteration 8501 / 39200) loss: 1.185197
(Iteration 8601 / 39200) loss: 1.114226
(Iteration 8701 / 39200) loss: 1.235766
(Iteration 8801 / 39200) loss: 1.042444
(Epoch 18 / 80) train acc: 0.645000; val_acc: 0.564000
(Iteration 8901 / 39200) loss: 1.190117
(Iteration 9001 / 39200) loss: 1.151252
(Iteration 9101 / 39200) loss: 1.242444
(Iteration 9201 / 39200) loss: 1.226552
(Iteration 9301 / 39200) loss: 1.112046
(Epoch 19 / 80) train acc: 0.686000; val_acc: 0.573000
(Iteration 9401 / 39200) loss: 1.085735
(Iteration 9501 / 39200) loss: 1.150893
(Iteration 9601 / 39200) loss: 1.142801
(Iteration 9701 / 39200) loss: 1.243212
(Epoch 20 / 80) train acc: 0.688000; val_acc: 0.580000
(Iteration 9801 / 39200) loss: 1.100593
(Iteration 9901 / 39200) loss: 1.195339
(Iteration 10001 / 39200) loss: 1.455515
(Iteration 10101 / 39200) loss: 1.102199
(Iteration 10201 / 39200) loss: 1.211633
(Epoch 21 / 80) train acc: 0.689000; val_acc: 0.567000
(Iteration 10301 / 39200) loss: 1.194243
(Iteration 10401 / 39200) loss: 1.110677
(Iteration 10501 / 39200) loss: 1.151281
(Iteration 10601 / 39200) loss: 1.167924
(Iteration 10701 / 39200) loss: 1.079824
(Epoch 22 / 80) train acc: 0.693000; val_acc: 0.569000
(Iteration 10801 / 39200) loss: 1.076431
(Iteration 10901 / 39200) loss: 1.072421
(Iteration 11001 / 39200) loss: 1.313280
(Iteration 11101 / 39200) loss: 1.213015
(Iteration 11201 / 39200) loss: 1.004086
```

```
(Epoch 23 / 80) train acc: 0.691000; val_acc: 0.576000
(Iteration 11301 / 39200) loss: 0.811971
(Iteration 11401 / 39200) loss: 0.999775
(Iteration 11501 / 39200) loss: 1.087524
(Iteration 11601 / 39200) loss: 0.891453
(Iteration 11701 / 39200) loss: 0.944478
(Epoch 24 / 80) train acc: 0.712000; val_acc: 0.576000
(Iteration 11801 / 39200) loss: 1.135463
(Iteration 11901 / 39200) loss: 1.051916
(Iteration 12001 / 39200) loss: 1.058195
(Iteration 12101 / 39200) loss: 1.094761
(Iteration 12201 / 39200) loss: 1.218822
(Epoch 25 / 80) train acc: 0.700000; val_acc: 0.587000
(Iteration 12301 / 39200) loss: 1.095869
(Iteration 12401 / 39200) loss: 0.981534
(Iteration 12501 / 39200) loss: 1.040588
(Iteration 12601 / 39200) loss: 1.101972
(Iteration 12701 / 39200) loss: 0.979700
(Epoch 26 / 80) train acc: 0.727000; val_acc: 0.580000
(Iteration 12801 / 39200) loss: 0.884404
(Iteration 12901 / 39200) loss: 1.039024
(Iteration 13001 / 39200) loss: 0.881299
(Iteration 13101 / 39200) loss: 1.123838
(Iteration 13201 / 39200) loss: 1.139810
(Epoch 27 / 80) train acc: 0.726000; val_acc: 0.585000
(Iteration 13301 / 39200) loss: 0.981287
(Iteration 13401 / 39200) loss: 0.906821
(Iteration 13501 / 39200) loss: 0.915128
(Iteration 13601 / 39200) loss: 1.141380
(Iteration 13701 / 39200) loss: 0.944810
(Epoch 28 / 80) train acc: 0.730000; val_acc: 0.588000
(Iteration 13801 / 39200) loss: 0.835803
(Iteration 13901 / 39200) loss: 0.983092
(Iteration 14001 / 39200) loss: 1.012665
(Iteration 14101 / 39200) loss: 1.056424
(Iteration 14201 / 39200) loss: 0.922939
(Epoch 29 / 80) train acc: 0.750000; val_acc: 0.573000
(Iteration 14301 / 39200) loss: 0.889565
(Iteration 14401 / 39200) loss: 1.063439
(Iteration 14501 / 39200) loss: 0.716205
(Iteration 14601 / 39200) loss: 0.810242
(Epoch 30 / 80) train acc: 0.722000; val_acc: 0.580000
(Iteration 14701 / 39200) loss: 1.299068
(Iteration 14801 / 39200) loss: 0.808461
(Iteration 14901 / 39200) loss: 0.870385
(Iteration 15001 / 39200) loss: 1.024490
(Iteration 15101 / 39200) loss: 1.034015
(Epoch 31 / 80) train acc: 0.754000; val_acc: 0.586000
(Iteration 15201 / 39200) loss: 1.134013
(Iteration 15301 / 39200) loss: 0.978965
(Iteration 15401 / 39200) loss: 0.951913
(Iteration 15501 / 39200) loss: 1.156696
(Iteration 15601 / 39200) loss: 0.845961
(Epoch 32 / 80) train acc: 0.750000; val_acc: 0.582000
(Iteration 15701 / 39200) loss: 0.966734
(Iteration 15801 / 39200) loss: 0.742249

(Iteration 15901 / 39200) loss: 0.811808
(Iteration 16001 / 39200) loss: 1.023509
(Iteration 16101 / 39200) loss: 0.903440
(Epoch 33 / 80) train acc: 0.727000; val_acc: 0.586000
(Iteration 16201 / 39200) loss: 0.947621
(Iteration 16301 / 39200) loss: 0.775905
(Iteration 16401 / 39200) loss: 1.019085
(Iteration 16501 / 39200) loss: 0.961113
(Iteration 16601 / 39200) loss: 0.977766
(Epoch 34 / 80) train acc: 0.742000; val_acc: 0.583000
(Iteration 16701 / 39200) loss: 0.783518
(Iteration 16801 / 39200) loss: 1.039373
(Iteration 16901 / 39200) loss: 0.858510
(Iteration 17001 / 39200) loss: 0.866176
(Iteration 17101 / 39200) loss: 0.983566
(Epoch 35 / 80) train acc: 0.755000; val_acc: 0.588000
(Iteration 17201 / 39200) loss: 1.048066
(Iteration 17301 / 39200) loss: 0.989369
(Iteration 17401 / 39200) loss: 1.015185
(Iteration 17501 / 39200) loss: 1.068445
(Iteration 17601 / 39200) loss: 0.851720
(Epoch 36 / 80) train acc: 0.755000; val_acc: 0.585000
(Iteration 17701 / 39200) loss: 1.036607
(Iteration 17801 / 39200) loss: 0.897850
(Iteration 17901 / 39200) loss: 0.908984
(Iteration 18001 / 39200) loss: 0.939887
(Iteration 18101 / 39200) loss: 0.835395
(Epoch 37 / 80) train acc: 0.767000; val_acc: 0.578000
(Iteration 18201 / 39200) loss: 1.042282
(Iteration 18301 / 39200) loss: 0.758919
(Iteration 18401 / 39200) loss: 1.002052
(Iteration 18501 / 39200) loss: 0.841024
```

```
(Iteration 18601 / 39200) loss: 0.882245
(Epoch 38 / 80) train acc: 0.765000; val_acc: 0.594000
(Iteration 18701 / 39200) loss: 0.958950
(Iteration 18801 / 39200) loss: 0.851094
(Iteration 18901 / 39200) loss: 0.708583
(Iteration 19001 / 39200) loss: 0.770913
(Iteration 19101 / 39200) loss: 0.733707
(Epoch 39 / 80) train acc: 0.778000; val_acc: 0.590000
(Iteration 19201 / 39200) loss: 0.990309
(Iteration 19301 / 39200) loss: 0.958403
(Iteration 19401 / 39200) loss: 0.992356
(Iteration 19501 / 39200) loss: 1.005794
(Epoch 40 / 80) train acc: 0.752000; val_acc: 0.588000
(Iteration 19601 / 39200) loss: 1.032480
(Iteration 19701 / 39200) loss: 1.035892
(Iteration 19801 / 39200) loss: 0.906792
(Iteration 19901 / 39200) loss: 0.895575
(Iteration 20001 / 39200) loss: 0.825461
(Epoch 41 / 80) train acc: 0.765000; val_acc: 0.591000
(Iteration 20101 / 39200) loss: 0.876272
(Iteration 20201 / 39200) loss: 0.902269
(Iteration 20301 / 39200) loss: 0.834744
(Iteration 20401 / 39200) loss: 0.961422
(Iteration 20501 / 39200) loss: 0.975309
(Epoch 42 / 80) train acc: 0.795000; val_acc: 0.585000
(Iteration 20601 / 39200) loss: 1.038000
(Iteration 20701 / 39200) loss: 0.837005
(Iteration 20801 / 39200) loss: 0.717900
(Iteration 20901 / 39200) loss: 0.850868
(Iteration 21001 / 39200) loss: 1.017165
(Epoch 43 / 80) train acc: 0.752000; val_acc: 0.585000
(Iteration 21101 / 39200) loss: 0.905929
(Iteration 21201 / 39200) loss: 1.025768
(Iteration 21301 / 39200) loss: 0.923585
(Iteration 21401 / 39200) loss: 0.892869
(Iteration 21501 / 39200) loss: 0.887772
(Epoch 44 / 80) train acc: 0.777000; val_acc: 0.590000
(Iteration 21601 / 39200) loss: 0.925660
(Iteration 21701 / 39200) loss: 0.861675
(Iteration 21801 / 39200) loss: 0.809872
(Iteration 21901 / 39200) loss: 0.920505
(Iteration 22001 / 39200) loss: 1.009144
(Epoch 45 / 80) train acc: 0.786000; val_acc: 0.591000
(Iteration 22101 / 39200) loss: 0.856028
(Iteration 22201 / 39200) loss: 0.866785
(Iteration 22301 / 39200) loss: 0.830510
(Iteration 22401 / 39200) loss: 1.279860
(Iteration 22501 / 39200) loss: 0.971482
(Epoch 46 / 80) train acc: 0.806000; val_acc: 0.593000
(Iteration 22601 / 39200) loss: 0.905765
(Iteration 22701 / 39200) loss: 1.066834
(Iteration 22801 / 39200) loss: 0.794572
(Iteration 22901 / 39200) loss: 1.013506
(Iteration 23001 / 39200) loss: 1.105925
(Epoch 47 / 80) train acc: 0.797000; val_acc: 0.591000
(Iteration 23101 / 39200) loss: 0.908778
(Iteration 23201 / 39200) loss: 0.864585
(Iteration 23301 / 39200) loss: 0.902531
(Iteration 23401 / 39200) loss: 0.800342
(Iteration 23501 / 39200) loss: 0.721073
(Epoch 48 / 80) train acc: 0.781000; val_acc: 0.595000
(Iteration 23601 / 39200) loss: 0.840992
(Iteration 23701 / 39200) loss: 1.122623
(Iteration 23801 / 39200) loss: 1.032104
(Iteration 23901 / 39200) loss: 0.909329
(Iteration 24001 / 39200) loss: 0.642137
(Epoch 49 / 80) train acc: 0.799000; val_acc: 0.593000
(Iteration 24101 / 39200) loss: 1.074038
(Iteration 24201 / 39200) loss: 0.888440
(Iteration 24301 / 39200) loss: 0.732228
(Iteration 24401 / 39200) loss: 0.830554
(Epoch 50 / 80) train acc: 0.803000; val_acc: 0.588000
(Iteration 24501 / 39200) loss: 0.690772
(Iteration 24601 / 39200) loss: 1.048687
(Iteration 24701 / 39200) loss: 1.012022
(Iteration 24801 / 39200) loss: 0.850425
(Iteration 24901 / 39200) loss: 0.838619
(Epoch 51 / 80) train acc: 0.796000; val_acc: 0.588000
(Iteration 25001 / 39200) loss: 0.654319
(Iteration 25101 / 39200) loss: 0.797890
(Iteration 25201 / 39200) loss: 0.852705
(Iteration 25301 / 39200) loss: 0.917384
(Iteration 25401 / 39200) loss: 0.845211
(Epoch 52 / 80) train acc: 0.813000; val_acc: 0.595000
(Iteration 25501 / 39200) loss: 0.753940
(Iteration 25601 / 39200) loss: 0.997589
(Iteration 25701 / 39200) loss: 0.775196
(Iteration 25801 / 39200) loss: 0.917092
```

```
(Iteration 25901 / 39200) loss: 0.954966
(Epoch 53 / 80) train acc: 0.785000; val_acc: 0.598000
(Iteration 26001 / 39200) loss: 0.828941
(Iteration 26101 / 39200) loss: 0.922264
(Iteration 26201 / 39200) loss: 0.748944
(Iteration 26301 / 39200) loss: 0.830192
(Iteration 26401 / 39200) loss: 0.816747
(Epoch 54 / 80) train acc: 0.783000; val_acc: 0.591000
(Iteration 26501 / 39200) loss: 0.682134
(Iteration 26601 / 39200) loss: 0.801736
(Iteration 26701 / 39200) loss: 0.749627
(Iteration 26801 / 39200) loss: 1.052627
(Iteration 26901 / 39200) loss: 0.990648
(Epoch 55 / 80) train acc: 0.816000; val_acc: 0.597000
(Iteration 27001 / 39200) loss: 0.950973
(Iteration 27101 / 39200) loss: 0.734878
(Iteration 27201 / 39200) loss: 0.747308
(Iteration 27301 / 39200) loss: 0.884391
(Iteration 27401 / 39200) loss: 0.918879
(Epoch 56 / 80) train acc: 0.797000; val_acc: 0.594000
(Iteration 27501 / 39200) loss: 0.853080
(Iteration 27601 / 39200) loss: 0.990341
(Iteration 27701 / 39200) loss: 0.884415
(Iteration 27801 / 39200) loss: 0.795938
(Iteration 27901 / 39200) loss: 0.919736
(Epoch 57 / 80) train acc: 0.829000; val_acc: 0.592000
(Iteration 28001 / 39200) loss: 0.795244
(Iteration 28101 / 39200) loss: 0.804721
(Iteration 28201 / 39200) loss: 0.958718
(Iteration 28301 / 39200) loss: 1.037316
(Iteration 28401 / 39200) loss: 0.912411
(Epoch 58 / 80) train acc: 0.795000; val_acc: 0.598000
(Iteration 28501 / 39200) loss: 0.786430
(Iteration 28601 / 39200) loss: 0.824360
(Iteration 28701 / 39200) loss: 0.857696
(Iteration 28801 / 39200) loss: 0.716881
(Iteration 28901 / 39200) loss: 0.745188
(Epoch 59 / 80) train acc: 0.827000; val_acc: 0.599000
(Iteration 29001 / 39200) loss: 0.836852
(Iteration 29101 / 39200) loss: 0.952089
(Iteration 29201 / 39200) loss: 0.784935
(Iteration 29301 / 39200) loss: 0.845400
(Epoch 60 / 80) train acc: 0.826000; val_acc: 0.586000
(Iteration 29401 / 39200) loss: 0.868991
(Iteration 29501 / 39200) loss: 0.695192
(Iteration 29601 / 39200) loss: 0.853898
(Iteration 29701 / 39200) loss: 0.978458
(Iteration 29801 / 39200) loss: 0.727721
(Epoch 61 / 80) train acc: 0.802000; val_acc: 0.587000
(Iteration 29901 / 39200) loss: 0.705634
(Iteration 30001 / 39200) loss: 0.829131
(Iteration 30101 / 39200) loss: 0.955335
(Iteration 30201 / 39200) loss: 0.861911
(Iteration 30301 / 39200) loss: 0.960021
(Epoch 62 / 80) train acc: 0.823000; val_acc: 0.583000
(Iteration 30401 / 39200) loss: 1.059026
(Iteration 30501 / 39200) loss: 0.820306
(Iteration 30601 / 39200) loss: 0.889187
(Iteration 30701 / 39200) loss: 0.888176
(Iteration 30801 / 39200) loss: 0.701103
(Epoch 63 / 80) train acc: 0.816000; val_acc: 0.587000
(Iteration 30901 / 39200) loss: 0.750650
(Iteration 31001 / 39200) loss: 0.976211
(Iteration 31101 / 39200) loss: 0.729728
(Iteration 31201 / 39200) loss: 0.829218
(Iteration 31301 / 39200) loss: 0.711588
(Epoch 64 / 80) train acc: 0.818000; val_acc: 0.588000
(Iteration 31401 / 39200) loss: 0.781439
(Iteration 31501 / 39200) loss: 0.746261

(Iteration 31601 / 39200) loss: 0.707213
(Iteration 31701 / 39200) loss: 0.833365
(Iteration 31801 / 39200) loss: 0.814129
(Epoch 65 / 80) train acc: 0.831000; val_acc: 0.578000
(Iteration 31901 / 39200) loss: 0.823512
(Iteration 32001 / 39200) loss: 0.797918
(Iteration 32101 / 39200) loss: 0.741748
(Iteration 32201 / 39200) loss: 0.908855
(Iteration 32301 / 39200) loss: 1.023562
(Epoch 66 / 80) train acc: 0.820000; val_acc: 0.587000
(Iteration 32401 / 39200) loss: 1.022008
(Iteration 32501 / 39200) loss: 0.963455
(Iteration 32601 / 39200) loss: 0.807658
(Iteration 32701 / 39200) loss: 0.779625
(Iteration 32801 / 39200) loss: 1.084706
(Epoch 67 / 80) train acc: 0.814000; val_acc: 0.590000
(Iteration 32901 / 39200) loss: 0.747726
(Iteration 33001 / 39200) loss: 0.897594
(Iteration 33101 / 39200) loss: 0.958651
```

```
(Iteration 33201 / 39200) loss: 0.731197
(Iteration 33301 / 39200) loss: 1.067751
(Epoch 68 / 80) train acc: 0.810000; val_acc: 0.593000
(Iteration 33401 / 39200) loss: 0.932779
(Iteration 33501 / 39200) loss: 0.934577
(Iteration 33601 / 39200) loss: 0.699594
(Iteration 33701 / 39200) loss: 0.728579
(Iteration 33801 / 39200) loss: 0.832864
(Epoch 69 / 80) train acc: 0.801000; val_acc: 0.591000
(Iteration 33901 / 39200) loss: 0.767292
(Iteration 34001 / 39200) loss: 0.882040
(Iteration 34101 / 39200) loss: 0.705112
(Iteration 34201 / 39200) loss: 0.948548
(Epoch 70 / 80) train acc: 0.813000; val_acc: 0.586000
(Iteration 34301 / 39200) loss: 0.759878
(Iteration 34401 / 39200) loss: 0.910183
(Iteration 34501 / 39200) loss: 0.665622
(Iteration 34601 / 39200) loss: 0.957756
(Iteration 34701 / 39200) loss: 0.611733
(Epoch 71 / 80) train acc: 0.810000; val_acc: 0.594000
(Iteration 34801 / 39200) loss: 0.929759
(Iteration 34901 / 39200) loss: 0.701143
(Iteration 35001 / 39200) loss: 0.796047
(Iteration 35101 / 39200) loss: 0.673842
(Iteration 35201 / 39200) loss: 0.706945
(Epoch 72 / 80) train acc: 0.825000; val_acc: 0.587000
(Iteration 35301 / 39200) loss: 0.809401
(Iteration 35401 / 39200) loss: 0.721291
(Iteration 35501 / 39200) loss: 0.876259
(Iteration 35601 / 39200) loss: 0.836721
(Iteration 35701 / 39200) loss: 0.792897
(Epoch 73 / 80) train acc: 0.835000; val_acc: 0.589000
(Iteration 35801 / 39200) loss: 1.015434
(Iteration 35901 / 39200) loss: 0.775106
(Iteration 36001 / 39200) loss: 0.770852
(Iteration 36101 / 39200) loss: 0.818432
(Iteration 36201 / 39200) loss: 0.903966
(Epoch 74 / 80) train acc: 0.815000; val_acc: 0.597000
(Iteration 36301 / 39200) loss: 0.890942
(Iteration 36401 / 39200) loss: 0.761177
(Iteration 36501 / 39200) loss: 0.618575
(Iteration 36601 / 39200) loss: 0.836164
(Iteration 36701 / 39200) loss: 0.791301
(Epoch 75 / 80) train acc: 0.832000; val_acc: 0.592000
(Iteration 36801 / 39200) loss: 0.911602
(Iteration 36901 / 39200) loss: 0.706361
(Iteration 37001 / 39200) loss: 0.727411
(Iteration 37101 / 39200) loss: 0.721145
(Iteration 37201 / 39200) loss: 0.926907
(Epoch 76 / 80) train acc: 0.836000; val_acc: 0.594000
(Iteration 37301 / 39200) loss: 0.798053
(Iteration 37401 / 39200) loss: 0.893760
(Iteration 37501 / 39200) loss: 0.773815
(Iteration 37601 / 39200) loss: 0.768303
(Iteration 37701 / 39200) loss: 0.823549
(Epoch 77 / 80) train acc: 0.820000; val_acc: 0.590000
(Iteration 37801 / 39200) loss: 0.743476
(Iteration 37901 / 39200) loss: 0.719722
(Iteration 38001 / 39200) loss: 0.703971
(Iteration 38101 / 39200) loss: 0.725379
(Iteration 38201 / 39200) loss: 1.030217
(Epoch 78 / 80) train acc: 0.813000; val_acc: 0.591000
(Iteration 38301 / 39200) loss: 0.901584
(Iteration 38401 / 39200) loss: 1.047116
(Iteration 38501 / 39200) loss: 0.754022
(Iteration 38601 / 39200) loss: 0.612967
(Iteration 38701 / 39200) loss: 0.597829
(Epoch 79 / 80) train acc: 0.830000; val_acc: 0.590000
(Iteration 38801 / 39200) loss: 0.771769
(Iteration 38901 / 39200) loss: 0.576011
(Iteration 39001 / 39200) loss: 0.735640
(Iteration 39101 / 39200) loss: 0.777838
(Epoch 80 / 80) train acc: 0.830000; val_acc: 0.590000
Validation set accuracy: 0.604
Test set accuracy: 0.597
```

In [ ]: 1

```python
import numpy as np
import pdb

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

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

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

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

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

  # ================================================================ #
  # YOUR CODE HERE:
  #   Calculate the output of the forward pass.  Notice the dimensions
  #   of w are D x M, which is the transpose of what we did in earlier
  #   assignments.
  # ================================================================ #

  x_reshape = x.reshape((x.shape[0], w.shape[0])) # N * D
  out = x_reshape.dot(w) + b.reshape((1, b.shape[0])) # N * M

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

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


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

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

  Returns a tuple of:
  - dx: Gradient with respect to x, of shape (N, d1, ..., d_k)
  - dw: Gradient with respect to w, of shape (D, M)
  - db: Gradient with respect to b, of shape (M,)
  """
  x, w, b = cache
  dx, dw, db = None, None, None

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

  x_reshape = np.reshape(x, (x.shape[0], w.shape[0]))
  dx_reshape = dout.dot(w.T)
  dx = np.reshape(dx_reshape, x.shape) # N * D
  dw = x_reshape.T.dot(dout) # D * M
  db = dout.T.dot(np.ones(x.shape[0]))  # M * 1

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

```python
 84    return dx, dw, db
 85
 86 def relu_forward(x):
 87    """
 88    Computes the forward pass for a layer of rectified linear units (ReLUs).
 89
 90    Input:
 91    - x: Inputs, of any shape
 92
 93    Returns a tuple of:
 94    - out: Output, of the same shape as x
 95    - cache: x
 96    """
 97    # =============================================================== #
 98    # YOUR CODE HERE:
 99    #   Implement the ReLU forward pass.
100    # =============================================================== #
101
102    out = np.maximum(0, x)
103
104    # =============================================================== #
105    # END YOUR CODE HERE
106    # =============================================================== #
107
108    cache = x
109    return out, cache
110
111
112 def relu_backward(dout, cache):
113    """
114    Computes the backward pass for a layer of rectified linear units (ReLUs).
115
116    Input:
117    - dout: Upstream derivatives, of any shape
118    - cache: Input x, of same shape as dout
119
120    Returns:
121    - dx: Gradient with respect to x
122    """
123    x = cache
124
125    # =============================================================== #
126    # YOUR CODE HERE:
127    #   Implement the ReLU backward pass
128    # =============================================================== #
129
130    dx = (x > 0) * (dout)
131
132    # =============================================================== #
133    # END YOUR CODE HERE
134    # =============================================================== #
135
136    return dx
137
138 def batchnorm_forward(x, gamma, beta, bn_param):
139    """
140    Forward pass for batch normalization.
141
142    During training the sample mean and (uncorrected) sample variance are
143    computed from minibatch statistics and used to normalize the incoming data.
144    During training we also keep an exponentially decaying running mean of the
   mean
145    and variance of each feature, and these averages are used to normalize data
146    at test-time.
147
148    At each timestep we update the running averages for mean and variance using
149    an exponential decay based on the momentum parameter:
150
151    running_mean = momentum * running_mean + (1 - momentum) * sample_mean
152    running_var = momentum * running_var + (1 - momentum) * sample_var
153
154    Note that the batch normalization paper suggests a different test-time
155    behavior: they compute sample mean and variance for each feature using a
156    large number of training images rather than using a running average. For
157    this implementation we have chosen to use running averages instead since
158    they do not require an additional estimation step; the torch7
   implementation
159    of batch normalization also uses running averages.
160
161    Input:
162    - x: Data of shape (N, D)
163    - gamma: Scale parameter of shape (D,)
164    - beta: Shift paremeter of shape (D,)
165    - bn_param: Dictionary with the following keys:
```

```python
166      - mode: 'train' or 'test'; required
167      - eps: Constant for numeric stability
168      - momentum: Constant for running mean / variance.
169      - running_mean: Array of shape (D,) giving running mean of features
170      - running_var Array of shape (D,) giving running variance of features
171
172    Returns a tuple of:
173      - out: of shape (N, D)
174      - cache: A tuple of values needed in the backward pass
175    """
176    mode = bn_param['mode']
177    eps = bn_param.get('eps', 1e-5)
178    momentum = bn_param.get('momentum', 0.9)
179
180    N, D = x.shape
181    running_mean = bn_param.get('running_mean', np.zeros(D, dtype=x.dtype))
182    running_var = bn_param.get('running_var', np.zeros(D, dtype=x.dtype))
183
184    out, cache = None, None
185    if mode == 'train':
186
187        # =============================================================== #
188        # YOUR CODE HERE:
189        #   A few steps here:
190        #     (1) Calculate the running mean and variance of the minibatch.
191        #     (2) Normalize the activations with the batch mean and variance.
192        #     (3) Scale and shift the normalized activations.  Store this
193        #         as the variable 'out'
194        #     (4) Store any variables you may need for the backward pass in
195        #         the 'cache' variable.
196        # =============================================================== #
197
198        minibatch_mean = np.mean(x, axis=0)
199        minibatch_var = np.var(x, axis=0)
200        x_normalize = (x - minibatch_mean) / np.sqrt(minibatch_var + eps)
201        out = gamma * x_normalize + beta
202
203        running_mean = momentum * running_mean + (1 - momentum) * minibatch_mean
204        running_var = momentum * running_var + (1 - momentum) * minibatch_var
205        bn_param['running_mean'] = running_mean
206        bn_param['running_var'] = running_var
207
208        cache = {
209            'minibatch_var': minibatch_var,
210            'x_centralize': (x - minibatch_mean),
211            'x_normalize': x_normalize,
212            'gamma': gamma,
213            'eps': eps
214        }
215
216        # =============================================================== #
217        # END YOUR CODE HERE
218        # =============================================================== #
219
220    elif mode == 'test':
221
222        # =============================================================== #
223        # YOUR CODE HERE:
224        #   Calculate the testing time normalized activations.  Normalize using
225        #   the running mean and variance, and then scale and shift
appropriately.
226        #   Store the output as 'out'.
227        # =============================================================== #
228
229        out = gamma * (x - running_mean) / np.sqrt(running_var + eps) + beta
230
231        # =============================================================== #
232        # END YOUR CODE HERE
233        # =============================================================== #
234
235    else:
236        raise ValueError('Invalid forward batchnorm mode "%s"' % mode)
237
238    # Store the updated running means back into bn_param
239    bn_param['running_mean'] = running_mean
240    bn_param['running_var'] = running_var
241
242    return out, cache
243
244  def batchnorm_backward(dout, cache):
245    """
246    Backward pass for batch normalization.
247
248    For this implementation, you should write out a computation graph for
```

```
249      batch normalization on paper and propagate gradients backward through
250      intermediate nodes.
251
252      Inputs:
253      - dout: Upstream derivatives, of shape (N, D)
254      - cache: Variable of intermediates from batchnorm_forward.
255
256      Returns a tuple of:
257      - dx: Gradient with respect to inputs x, of shape (N, D)
258      - dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
259      - dbeta: Gradient with respect to shift parameter beta, of shape (D,)
260      """
261      dx, dgamma, dbeta = None, None, None
262
263      # =============================================================== #
264      # YOUR CODE HERE:
265      #    Implement the batchnorm backward pass, calculating dx, dgamma, and
     dbeta.
266      # =============================================================== #
267
268      # get parameters from cache
269      N = dout.shape[0]
270      minibatch_var = cache.get('minibatch_var')
271      x_centralize = cache.get('x_centralize')
272      x_normalize = cache.get('x_normalize')
273      gamma = cache.get('gamma')
274      eps = cache.get('eps')
275
276      # calculate dx
277      dxhat = dout * gamma
278      dxmu1 = dxhat / np.sqrt(minibatch_var + eps)
279      sqrt_var = np.sqrt(minibatch_var + eps)
280      dsqrt_var = -np.sum(dxhat * x_centralize, axis=0) / (sqrt_var**2)
281      dvar = dsqrt_var * 0.5 / sqrt_var
282      dxmu2 = 2 * x_centralize * dvar * np.ones_like(dout) / N
283      dx1 = dxmu1 + dxmu2
284      dx2 = -np.sum(dx1, axis=0) * np.ones_like(dout) / N
285      dx = dx1 + dx2
286
287      # calculate dbeta and dgamma
288      dbeta = np.sum(dout, axis=0)
289      dgamma = np.sum(dout * x_normalize, axis=0)
290
291      # =============================================================== #
292      # END YOUR CODE HERE
293      # =============================================================== #
294
295      return dx, dgamma, dbeta
296
297  def dropout_forward(x, dropout_param):
298      """
299      Performs the forward pass for (inverted) dropout.
300
301      Inputs:
302      - x: Input data, of any shape
303      - dropout_param: A dictionary with the following keys:
304        - p: Dropout parameter. We drop each neuron output with probability p.
305        - mode: 'test' or 'train'. If the mode is train, then perform dropout;
306          if the mode is test, then just return the input.
307        - seed: Seed for the random number generator. Passing seed makes this
308          function deterministic, which is needed for gradient checking but not
     in
309          real networks.
310
311      Outputs:
312      - out: Array of the same shape as x.
313      - cache: A tuple (dropout_param, mask). In training mode, mask is the
     dropout
314        mask that was used to multiply the input; in test mode, mask is None.
315      """
316      p, mode = dropout_param['p'], dropout_param['mode']
317      if 'seed' in dropout_param:
318          np.random.seed(dropout_param['seed'])
319
320      mask = None
321      out = None
322
323      if mode == 'train':
324          # =============================================================== #
325          # YOUR CODE HERE:
326          #    Implement the inverted dropout forward pass during training time.
327          #    Store the masked and scaled activations in out, and store the
328          #    dropout mask as the variable mask.
329          # =============================================================== #
```

```python
      mask = (np.random.random_sample(x.shape) >= 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)
```

```python
413      dx[margins > 0] = 1
414      dx[np.arange(N), y] -= num_pos
415      dx /= N
416      return loss, dx
417
418
419  def softmax_loss(x, y):
420      """
421      Computes the loss and gradient for softmax classification.
422
423      Inputs:
424      - x: Input data, of shape (N, C) where x[i, j] is the score for the jth
     class
425          for the ith input.
426      - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
427          0 <= y[i] < C
428
429      Returns a tuple of:
430      - loss: Scalar giving the loss
431      - dx: Gradient of the loss with respect to x
432      """
433
434      probs = np.exp(x - np.max(x, axis=1, keepdims=True))
435      probs /= np.sum(probs, axis=1, keepdims=True)
436      N = x.shape[0]
437      loss = -np.sum(np.log(probs[np.arange(N), y])) / N
438      dx = probs.copy()
439      dx[np.arange(N), y] -= 1
440      dx /= N
441      return loss, dx
442
```

```python
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 deteriminstc so we can gradient check the
        model.
    """
    self.use_batchnorm = use_batchnorm
    self.use_dropout = dropout > 0
    self.reg = reg
    self.num_layers = 1 + len(hidden_dims)
    self.dtype = dtype
    self.params = {}

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

```python
     77
     78    cur_dim = input_dim
     79    for idx, hidden_dim in enumerate(hidden_dims):
     80      # initialize weights and bias
     81      self.params['W' + str(idx + 1)] = np.random.randn(cur_dim, hidden_dim)
          * weight_scale
     82      self.params['b' + str(idx + 1)] = np.zeros(hidden_dim)
     83
     84      # initialize gammas and betas
     85      if self.use_batchnorm:
     86        self.params['gamma' + str(idx + 1)] = np.ones(hidden_dim)
     87        self.params['beta' + str(idx + 1)] = np.zeros(hidden_dim)
     88
     89      cur_dim = hidden_dim
     90
     91    self.params['W' + str(self.num_layers)] = np.random.randn(cur_dim,
          num_classes) * weight_scale
     92    self.params['b' + str(self.num_layers)] = np.zeros(num_classes)
     93
     94    # =============================================================== #
     95    # END YOUR CODE HERE
     96    # =============================================================== #
     97
     98    # When using dropout we need to pass a dropout_param dictionary to each
     99    # dropout layer so that the layer knows the dropout probability and the
          mode
    100    # (train / test). You can pass the same dropout_param to each dropout
          layer.
    101    self.dropout_param = {}
    102    if self.use_dropout:
    103      self.dropout_param = {'mode': 'train', 'p': dropout}
    104      if seed is not None:
    105        self.dropout_param['seed'] = seed
    106
    107    # With batch normalization we need to keep track of running means and
    108    # variances, so we need to pass a special bn_param object to each batch
    109    # normalization layer. You should pass self.bn_params[0] to the forward
          pass
    110    # of the first batch normalization layer, self.bn_params[1] to the
          forward
    111    # pass of the second batch normalization layer, etc.
    112    self.bn_params = []
    113    if self.use_batchnorm:
    114      self.bn_params = [{'mode': 'train'} for i in np.arange(self.num_layers
          - 1)]
    115
    116    # Cast all parameters to the correct datatype
    117    for k, v in self.params.items():
    118      self.params[k] = v.astype(dtype)
    119
    120
    121  def loss(self, X, y=None):
    122    """
    123    Compute loss and gradient for the fully-connected net.
    124
    125    Input / output: Same as TwoLayerNet above.
    126    """
    127    X = X.astype(self.dtype)
    128    mode = 'test' if y is None else 'train'
    129
    130    # Set train/test mode for batchnorm params and dropout param since they
    131    # behave differently during training and testing.
    132    if self.dropout_param is not None:
    133      self.dropout_param['mode'] = mode
    134    if self.use_batchnorm:
    135      for bn_param in self.bn_params:
    136        bn_param[mode] = mode
    137
    138    scores = None
    139
    140    # =============================================================== #
    141    # YOUR CODE HERE:
    142    #   Implement the forward pass of the FC net and store the output
    143    #   scores as the variable "scores".
    144    #
    145    #   BATCHNORM: If self.use_batchnorm is true, insert a bathnorm layer
    146    #   between the affine_forward and relu_forward layers.  You may
    147    #   also write an affine_batchnorm_relu() function in layer_utils.py.
    148    #
    149    #   DROPOUT: If dropout is non-zero, insert a dropout layer after
    150    #   every ReLU layer.
    151    # =============================================================== #
    152
    153    # initialize caches
```

```python
154        fc_cache = {}
155        relu_cache = {}
156        batchnorm_cache = {}
157        dropout_cache = {}
158
159        # flatten image
160        X = np.reshape(X, [X.shape[0], -1])
161
162        # go through all layers
163        for i in range(self.num_layers - 1):
164          # fc layer
165          fc_out, fc_cache[str(i + 1)] = affine_forward(X, self.params['W' +
    str(i + 1)], self.params['b' + str(i + 1)])
166
167          # batchnorm layer
168          relu_input = fc_out
169          if self.use_batchnorm:
170            batchnorm_out, batchnorm_cache[str(i + 1)] =
    batchnorm_forward(fc_out, self.params['gamma' + str(i + 1)],
    self.params['beta' + str(i + 1)], self.bn_params[i])
171            relu_input = batchnorm_out
172          relu_out, relu_cache[str(i + 1)] = relu_forward(relu_input)
173
174          # dropout layer
175          if self.use_dropout:
176            relu_out, dropout_cache[str(i + 1)] = dropout_forward(relu_out,
    self.dropout_param)
177
178          # update X
179          X = relu_out.copy()
180
181        # output FC layer with no relu
182        scores, final_cache = affine_forward(X, self.params['W' +
    str(self.num_layers)], self.params['b' + str(self.num_layers)])
183
184        # ================================================================ #
185        # END YOUR CODE HERE
186        # ================================================================ #
187
188        # If test mode return early
189        if mode == 'test':
190          return scores
191
192        loss, grads = 0.0, {}
193        # ================================================================ #
194        # YOUR CODE HERE:
195        #   Implement the backwards pass of the FC net and store the gradients
196        #   in the grads dict, so that grads[k] is the gradient of self.params[k]
197        #   Be sure your L2 regularization includes a 0.5 factor.
198        #
199        #   BATCHNORM: Incorporate the backward pass of the batchnorm.
200        #
201        #   DROPOUT: Incorporate the backward pass of dropout.
202        # ================================================================ #
203
204        # initialize
205        loss, dx = softmax_loss(scores, y)
206        loss += 0.5 * self.reg * (np.sum(np.square(self.params['W' +
    str(self.num_layers)])))
207        dx_back, dw_back, db_back = affine_backward(dx, final_cache)
208        grads['W' + str(self.num_layers)] = dw_back + self.reg * self.params['W'
    + str(self.num_layers)]
209        grads['b' + str(self.num_layers)] = db_back
210
211        # go backward all layers and update weights, bias, gammas and betas
212        for i in range(self.num_layers - 1, 0, -1):
213          # dropout layer
214          if self.use_dropout:
215            dx_back = dropout_backward(dx_back, dropout_cache[str(i)])
216          dx_relu = relu_backward(dx_back, relu_cache[str(i)])
217
218          # batchnorm layer
219          affine_backward_input = dx_relu
220          if self.use_batchnorm:
221            dx_bn, dgamma, dbeta = batchnorm_backward(dx_relu,
    batchnorm_cache[str(i)])
222            grads['gamma' + str(i)] = dgamma
223            grads['beta' + str(i)] = dbeta
224            affine_backward_input = dx_bn
225          dx_back, dw_back, db_back = affine_backward(affine_backward_input,
    fc_cache[str(i)])
226
227          grads['W' + str(i)] = dw_back + self.reg * self.params['W' + str(i)]
228          grads['b' + str(i)] = db_back
```

```python
            loss += 0.5 * self.reg * (np.sum(np.square(self.params['W' + str(i)])))

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

        return loss, grads
```

```
 1 import numpy as np
 2
 3 """
 4 This code was originally written for CS 231n at Stanford University
 5 (cs231n.stanford.edu).  It has been modified in various areas for use in the
 6 ECE 239AS class at UCLA.  This includes the descriptions of what code to
 7 implement as well as some slight potential changes in variable names to be
 8 consistent with class nomenclature.  We thank Justin Johnson & Serena Yeung
   for
 9 permission to use this code.  To see the original version, please visit
10 cs231n.stanford.edu.
11 """
12
13 """
14 This file implements various first-order update rules that are commonly used
   for
15 training neural networks. Each update rule accepts current weights and the
16 gradient of the loss with respect to those weights and produces the next set
   of
17 weights. Each update rule has the same interface:
18
19 def update(w, dw, config=None):
20
21 Inputs:
22   - w: A numpy array giving the current weights.
23   - dw: A numpy array of the same shape as w giving the gradient of the
24     loss with respect to w.
25   - config: A dictionary containing hyperparameter values such as learning
   rate,
26     momentum, etc. If the update rule requires caching values over many
27     iterations, then config will also hold these cached values.
28
29 Returns:
30   - next_w: The next point after the update.
31   - config: The config dictionary to be passed to the next iteration of the
32     update rule.
33
34 NOTE: For most update rules, the default learning rate will probably not
   perform
35 well; however the default values of the other hyperparameters should work
   well
36 for a variety of different problems.
37
38 For efficiency, update rules may perform in-place updates, mutating w and
39 setting next_w equal to w.
40 """
41
42
43 def sgd(w, dw, config=None):
44   """
45   Performs vanilla stochastic gradient descent.
46
47   config format:
48   - learning_rate: Scalar learning rate.
49   """
50   if config is None: config = {}
51   config.setdefault('learning_rate', 1e-2)
52
53   w -= config['learning_rate'] * dw
54   return w, config
55
56
57 def sgd_momentum(w, dw, config=None):
58   """
59   Performs stochastic gradient descent with momentum.
60
61   config format:
62   - learning_rate: Scalar learning rate.
63   - momentum: Scalar between 0 and 1 giving the momentum value.
64     Setting momentum = 0 reduces to sgd.
65   - velocity: A numpy array of the same shape as w and dw used to store a
   moving
66     average of the gradients.
67   """
68   if config is None: config = {}
69   config.setdefault('learning_rate', 1e-2)
70   config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
71   v = config.get('velocity', np.zeros_like(w))   # gets velocity, else sets
   it to zero.
72
73   # ============================================================ #
74   # YOUR CODE HERE:
75   #   Implement the momentum update formula.  Return the updated weights
76   #   as next_w, and store the updated velocity as v.
```

```
77    # ================================================================ #
78
79    v = config['momentum'] * v - config['learning_rate'] * dw
80    next_w = w + v
81
82    # ================================================================ #
83    # END YOUR CODE HERE
84    # ================================================================ #
85
86    config['velocity'] = v
87
88    return next_w, config
89
90  def sgd_nesterov_momentum(w, dw, config=None):
91    """
92    Performs stochastic gradient descent with Nesterov momentum.
93
94    config format:
95    - learning_rate: Scalar learning rate.
96    - momentum: Scalar between 0 and 1 giving the momentum value.
97      Setting momentum = 0 reduces to sgd.
98    - velocity: A numpy array of the same shape as w and dw used to store a
  moving
99      average of the gradients.
100   """
101   if config is None: config = {}
102   config.setdefault('learning_rate', 1e-2)
103   config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
104   v = config.get('velocity', np.zeros_like(w))   # gets velocity, else sets
  it to zero.
105
106   # ================================================================ #
107   # YOUR CODE HERE:
108   #   Implement the momentum update formula.  Return the updated weights
109   #   as next_w, and store the updated velocity as v.
110   # ================================================================ #
111
112   v_old = v
113   v = config['momentum'] * v - config['learning_rate'] * dw
114   w += v + config['momentum'] * (v - v_old)
115   next_w = w
116
117   # ================================================================ #
118   # END YOUR CODE HERE
119   # ================================================================ #
120
121   config['velocity'] = v
122
123   return next_w, config
124
125 def rmsprop(w, dw, config=None):
126   """
127   Uses the RMSProp update rule, which uses a moving average of squared
  gradient
128   values to set adaptive per-parameter learning rates.
129
130   config format:
131   - learning_rate: Scalar learning rate.
132   - decay_rate: Scalar between 0 and 1 giving the decay rate for the squared
133     gradient cache.
134   - epsilon: Small scalar used for smoothing to avoid dividing by zero.
135   - beta: Moving average of second moments of gradients.
136   """
137   if config is None: config = {}
138   config.setdefault('learning_rate', 1e-2)
139   config.setdefault('decay_rate', 0.99)
140   config.setdefault('epsilon', 1e-8)
141   config.setdefault('a', np.zeros_like(w))
142
143   next_w = None
144
145   # ================================================================ #
146   # YOUR CODE HERE:
147   #   Implement RMSProp.  Store the next value of w as next_w.  You need
148   #   to also store in config['a'] the moving average of the second
149   #   moment gradients, so they can be used for future gradients. Concretely,
150   #   config['a'] corresponds to "a" in the lecture notes.
151   # ================================================================ #
152
153   config['a'] = config['decay_rate'] * config['a'] + (1 -
  config['decay_rate']) * (dw**2)
154   next_w = w - config['learning_rate'] * dw / (np.sqrt(config['a']) +
  config['epsilon'])
155
```

```python
    # ================================================================ #
    # 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.
    - beta1: 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('beta1', 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 time.
    # ================================================================ #

    beta1 = config['beta1']
    beta2 = config['beta2']
    t = config['t'] + 1

    v = beta1 * config['v'] + (1 - beta1) * dw
    a = beta2 * config['a'] + (1 - beta2) * (dw**2)
    v_corrected = v / (1 - beta1**t)
    a_corrected = a / (1 - beta2**t)
    next_w = w - config['learning_rate'] * v_corrected / (np.sqrt(a_corrected)
+ config['epsilon'])

    config['v'] = v
    config['a'] = a
    config['t'] = t

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

    return next_w, config
```