

## 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](mailto:ece239as.w18@gmail.com).

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

In [7]:

```
## Import and setups

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

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

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

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

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

```
%reload_ext autoreload
```

In [9]:

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

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

## Building upon your HW #3 implementation

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

- `affine_forward` in `nndl/layers.py`
- `affine_backward` in `nndl/layers.py`
- `relu_forward` in `nndl/layers.py`
- `relu_backward` in `nndl/layers.py`
- `affine_relu_forward` in `nndl/layer_utils.py`
- `affine_relu_backward` in `nndl/layer_utils.py`
- The `FullyConnectedNet` class in `nndl/fc_net.py`

## Test all functions you copy and pasted

In [16]:

```
from nndl.layer_tests import *

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

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

difference: 9.769849468192957e-10

If affine\_backward is working, error should be less than 1e-9::

dx error: 2.8743685843251695e-10

dw error: 6.975022567890211e-10

db error: 3.9255446974462864e-11

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

difference: 4.999999798022158e-08

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

dx error: 3.275631636895981e-12

If affine\_relu\_forward and affine\_relu\_backward are working, error should be less than 1e-9::

dx error: 4.2236749436221695e-11

dw error: 9.589155550192892e-10

db error: 7.826689997629487e-12

Running check with reg = 0

Initial loss: 2.303526895251343

W1 relative error: 1.5841698922860941e-06

W2 relative error: 4.193331243584077e-06

W3 relative error: 6.720220636835107e-08

b1 relative error: 3.333143081548786e-08

b2 relative error: 9.174673340253911e-08

b3 relative error: 1.0003585687814655e-10

Running check with reg = 3.14

Initial loss: 7.639756633953297

W1 relative error: 1.7456418356571706e-08

W2 relative error: 5.273495549776094e-08

W3 relative error: 5.119768890898601e-09

b1 relative error: 7.335370292651373e-08

b2 relative error: 6.86606117889474e-09

b3 relative error: 3.726638992392252e-10

## Training a larger model

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

## SGD + momentum

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

In [4]:

```
from nndl.optim import sgd_momentum

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

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

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

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

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

## SGD + Nesterov momentum

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

In [5]:

```
from nndl.optim import sgd_nesterov_momentum

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

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

expected_next_w = np.asarray([
    [0.08714,      0.15246105,  0.21778211,  0.28310316,  0.34842421],
    [0.41374526,  0.47906632,  0.54438737,  0.60970842,  0.67502947],
    [0.74035053,  0.80567158,  0.87099263,  0.93631368,  1.00163474],
    [1.06695579,  1.13227684,  1.19759789,  1.26291895,  1.32824   ]])
expected_velocity = np.asarray([
    [ 0.5406,      0.55475789,  0.56891579,  0.58307368,  0.59723158],
    [ 0.61138947,  0.62554737,  0.63970526,  0.65386316,  0.66802105],
    [ 0.68217895,  0.69633684,  0.71049474,  0.72465263,  0.73881053],
    [ 0.75296842,  0.76712632,  0.78128421,  0.79544211,  0.8096   ]])

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

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

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

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

In [17]:

```
num_train = 4000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

solvers = {}

for update_rule in ['sgd', 'sgd_momentum', 'sgd_nesterov_momentum']:
    print('Optimizing with {}'.format(update_rule))
    model = FullyConnectedNet([100, 100, 100, 100, 100], weight_scale=5e-2)

    solver = Solver(model, small_data,
                    num_epochs=5, batch_size=100,
                    update_rule=update_rule,
                    optim_config={
                        'learning_rate': 1e-2,
                    },
                    verbose=False)
    solvers[update_rule] = solver
    solver.train()
    print

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

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

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

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

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

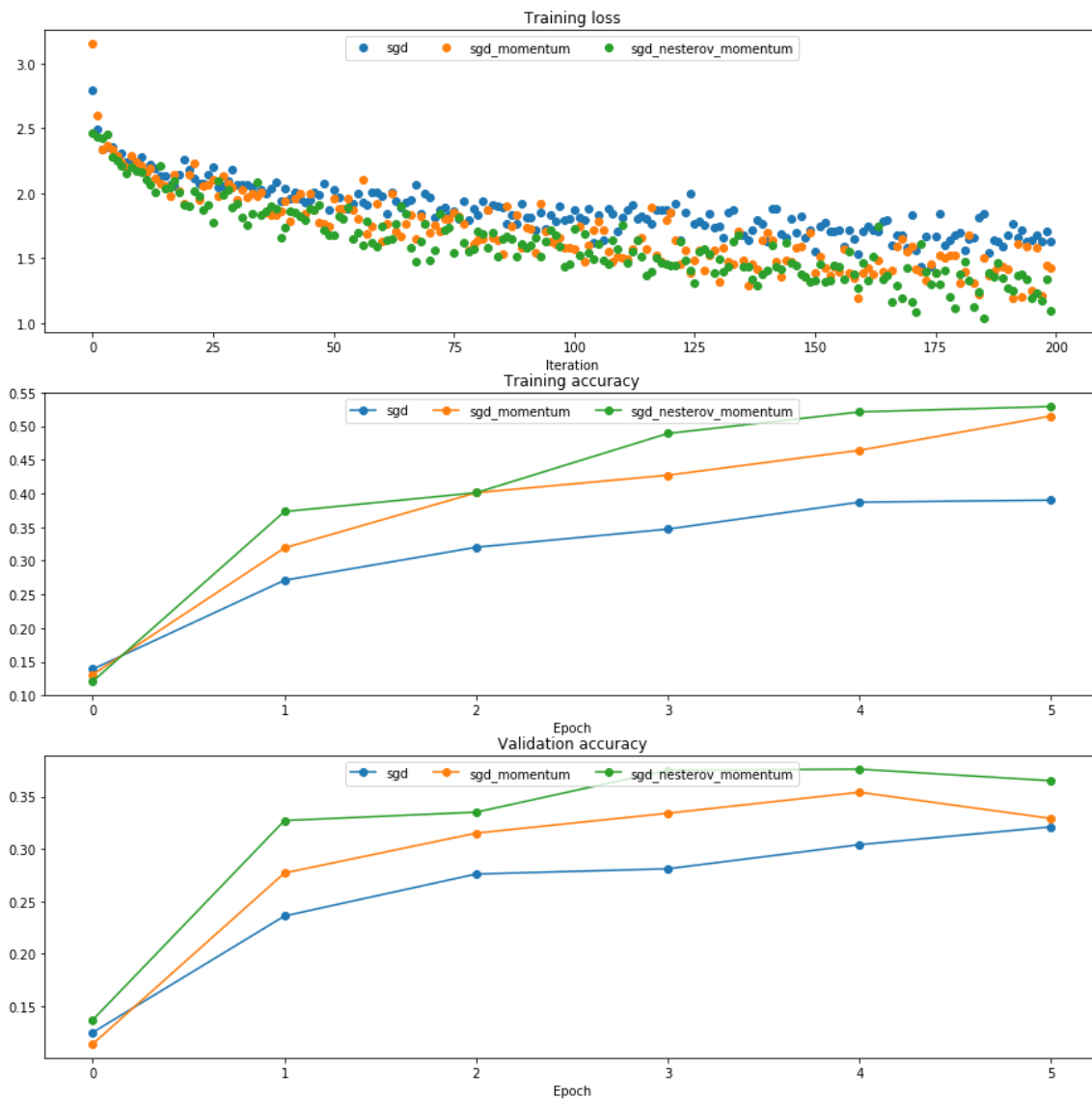
    plt.subplot(3, 1, 3)
    plt.plot(solver.val_acc_history, '-o', label=update_rule)

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

Optimizing with sgd  
Optimizing with sgd\_momentum  
Optimizing with sgd\_nesterov\_momentum

```
c:\users\aksha\appdata\local\programs\python\python36\lib\site-packages\matplotlib\cbook\deprecation.py:106: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.
```

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



## RMSProp

Now we go to techniques that adapt the gradient. Implement `rmsprop` in `nnd1/optim.py`. Test your implementation by running the cell below.



In [23]:

```
from nndl.optim import rmsprop

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

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

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

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

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

## Adaptive moments

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

In [25]:

```
# Test Adam implementation; you should see errors around 1e-7 or less
from nndl.optim import adam

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

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

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

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

next\_w error: 1.1395691798535431e-07

a error: 4.208314038113071e-09

v error: 4.214963193114416e-09

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

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

In [26]:

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

for update_rule in ['adam', 'rmsprop']:
    print('Optimizing with {}'.format(update_rule))
    model = FullyConnectedNet([100, 100, 100, 100, 100], weight_scale=5e-2)

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

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

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

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

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

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

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

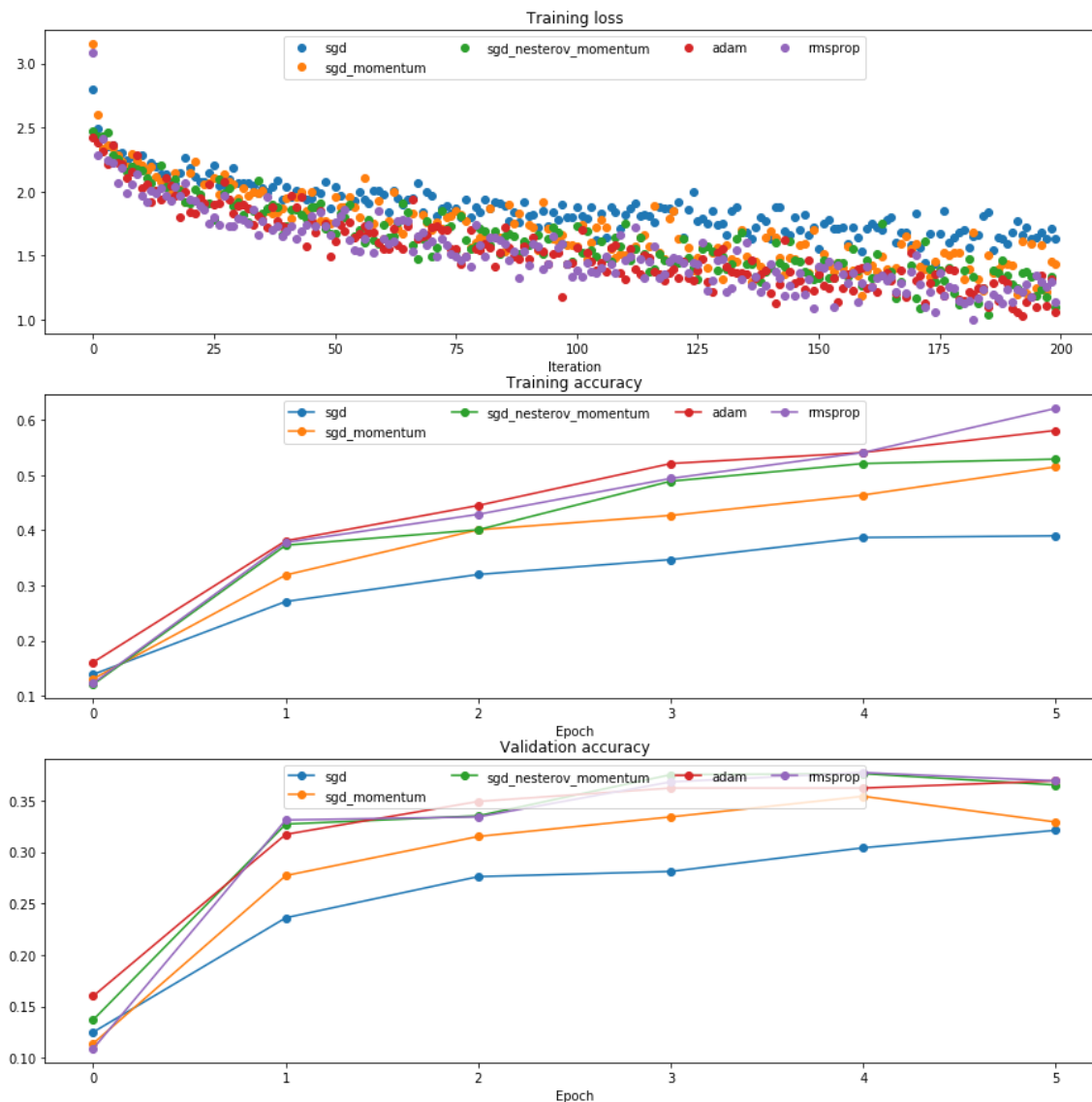
for i in [1, 2, 3]:
    plt.subplot(3, 1, i)
    plt.legend(loc='upper center', ncol=4)
plt.gcf().set_size_inches(15, 15)
plt.show()
```

Optimizing with adam

Optimizing with rmsprop

```
c:\users\aksha\appdata\local\programs\python\python36\lib\site-packages\matplotlib\cbook\deprecation.py:106: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.
```

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



## Easier optimization

In the following cell, we'll train a 4 layer neural network having 500 units in each hidden layer with the different optimizers, and find that it is far easier to get up to 50+% performance on CIFAR-10. After we implement batchnorm and dropout, we'll ask you to get 60+% on CIFAR-10.

In [30]:

```
optimizer = 'adam'
best_model = None

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

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

solver = Solver(model, data,
                 num_epochs=10, batch_size=100,
                 update_rule=optimizer,
                 optim_config={
                     'learning_rate': learning_rate,
                 },
                 lr_decay=lr_decay,
                 verbose=True, print_every=50)

solver.train()
```

```
(Iteration 1 / 4900) loss: 2.301169
(Epoch 0 / 10) train acc: 0.137000; val_acc: 0.140000
(Iteration 51 / 4900) loss: 1.685701
(Iteration 101 / 4900) loss: 1.775679
(Iteration 151 / 4900) loss: 1.709662
(Iteration 201 / 4900) loss: 1.745853
(Iteration 251 / 4900) loss: 1.670498
(Iteration 301 / 4900) loss: 1.611355
(Iteration 351 / 4900) loss: 1.633561
(Iteration 401 / 4900) loss: 1.572479
(Iteration 451 / 4900) loss: 1.508511
(Epoch 1 / 10) train acc: 0.446000; val_acc: 0.436000
(Iteration 501 / 4900) loss: 1.738661
(Iteration 551 / 4900) loss: 1.524508
(Iteration 601 / 4900) loss: 1.520372
(Iteration 651 / 4900) loss: 1.547033
(Iteration 701 / 4900) loss: 1.622758
(Iteration 751 / 4900) loss: 1.415125
(Iteration 801 / 4900) loss: 1.601579
(Iteration 851 / 4900) loss: 1.520090
(Iteration 901 / 4900) loss: 1.396371
(Iteration 951 / 4900) loss: 1.702087
(Epoch 2 / 10) train acc: 0.446000; val_acc: 0.451000
(Iteration 1001 / 4900) loss: 1.450721
(Iteration 1051 / 4900) loss: 1.377479
(Iteration 1101 / 4900) loss: 1.309599
(Iteration 1151 / 4900) loss: 1.395568
(Iteration 1201 / 4900) loss: 1.493719
(Iteration 1251 / 4900) loss: 1.330925
(Iteration 1301 / 4900) loss: 1.565845
(Iteration 1351 / 4900) loss: 1.363037
(Iteration 1401 / 4900) loss: 1.293119
(Iteration 1451 / 4900) loss: 1.281568
(Epoch 3 / 10) train acc: 0.515000; val_acc: 0.487000
(Iteration 1501 / 4900) loss: 1.555297
(Iteration 1551 / 4900) loss: 1.471259
(Iteration 1601 / 4900) loss: 1.388649
(Iteration 1651 / 4900) loss: 1.241068
(Iteration 1701 / 4900) loss: 1.468328
(Iteration 1751 / 4900) loss: 1.201056
(Iteration 1801 / 4900) loss: 1.132421
(Iteration 1851 / 4900) loss: 1.171905
(Iteration 1901 / 4900) loss: 1.280849
(Iteration 1951 / 4900) loss: 1.197774
(Epoch 4 / 10) train acc: 0.507000; val_acc: 0.492000
(Iteration 2001 / 4900) loss: 1.191567
(Iteration 2051 / 4900) loss: 1.261319
(Iteration 2101 / 4900) loss: 1.187541
(Iteration 2151 / 4900) loss: 1.024772
(Iteration 2201 / 4900) loss: 1.081363
(Iteration 2251 / 4900) loss: 1.185666
(Iteration 2301 / 4900) loss: 1.273112
(Iteration 2351 / 4900) loss: 1.211875
(Iteration 2401 / 4900) loss: 1.267167
(Epoch 5 / 10) train acc: 0.562000; val_acc: 0.515000
(Iteration 2451 / 4900) loss: 1.107632
(Iteration 2501 / 4900) loss: 1.178251
(Iteration 2551 / 4900) loss: 1.097444
(Iteration 2601 / 4900) loss: 1.382508
(Iteration 2651 / 4900) loss: 1.245508
(Iteration 2701 / 4900) loss: 1.105382
```

```

(Iteration 2751 / 4900) loss: 1.053378
(Iteration 2801 / 4900) loss: 1.170962
(Iteration 2851 / 4900) loss: 1.153845
(Iteration 2901 / 4900) loss: 1.095130
(Epoch 6 / 10) train acc: 0.585000; val_acc: 0.534000
(Iteration 2951 / 4900) loss: 1.360417
(Iteration 3001 / 4900) loss: 1.207257
(Iteration 3051 / 4900) loss: 1.016950
(Iteration 3101 / 4900) loss: 1.222481
(Iteration 3151 / 4900) loss: 1.063637
(Iteration 3201 / 4900) loss: 1.165328
(Iteration 3251 / 4900) loss: 1.210303
(Iteration 3301 / 4900) loss: 1.162433
(Iteration 3351 / 4900) loss: 1.026439
(Iteration 3401 / 4900) loss: 1.012245
(Epoch 7 / 10) train acc: 0.615000; val_acc: 0.539000
(Iteration 3451 / 4900) loss: 1.053071
(Iteration 3501 / 4900) loss: 1.078127
(Iteration 3551 / 4900) loss: 0.909880
(Iteration 3601 / 4900) loss: 1.017019
(Iteration 3651 / 4900) loss: 1.128381
(Iteration 3701 / 4900) loss: 0.946742
(Iteration 3751 / 4900) loss: 0.858189
(Iteration 3801 / 4900) loss: 1.047041
(Iteration 3851 / 4900) loss: 0.966694
(Iteration 3901 / 4900) loss: 1.045555
(Epoch 8 / 10) train acc: 0.626000; val_acc: 0.549000
(Iteration 3951 / 4900) loss: 0.877196
(Iteration 4001 / 4900) loss: 0.848048
(Iteration 4051 / 4900) loss: 0.962812
(Iteration 4101 / 4900) loss: 0.895925
(Iteration 4151 / 4900) loss: 1.296155
(Iteration 4201 / 4900) loss: 1.006795
(Iteration 4251 / 4900) loss: 0.863569
(Iteration 4301 / 4900) loss: 1.078027
(Iteration 4351 / 4900) loss: 1.105667
(Iteration 4401 / 4900) loss: 1.061704
(Epoch 9 / 10) train acc: 0.651000; val_acc: 0.532000
(Iteration 4451 / 4900) loss: 0.775457
(Iteration 4501 / 4900) loss: 1.186613
(Iteration 4551 / 4900) loss: 0.613210
(Iteration 4601 / 4900) loss: 1.021387
(Iteration 4651 / 4900) loss: 0.982887
(Iteration 4701 / 4900) loss: 1.120728
(Iteration 4751 / 4900) loss: 0.857032
(Iteration 4801 / 4900) loss: 1.062320
(Iteration 4851 / 4900) loss: 0.889548
(Epoch 10 / 10) train acc: 0.679000; val_acc: 0.530000

```

In [31]:

```

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: 0.549

Test set accuracy: 0.538



# 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](http://cs231n.stanford.edu)).

In [1]:

```
## Import and setups

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

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

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

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

In [2]:

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

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

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

## Batchnorm forward pass

Implement the training time batchnorm forward pass, `batchnorm_forward`, in `nnd1/layers.py`. After that, test your implementation by running the following cell.

In [3]:

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

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

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

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

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

```
Before batch normalization:
  means: [16.63073881  8.91888657  4.24324471]
  stds:  [26.43922081 36.9924524 27.1979339 ]
After batch normalization (gamma=1, beta=0)
  mean: [-1.47104551e-16 -1.21014310e-16 -3.16413562e-17]
  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 `nnd1/layers.py`. After that, test your implementation by running the following cell.

In [4]:

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

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

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

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

```
After batch normalization (test-time):
  means: [-0.08507578  0.06929814  0.11251958]
  stds:  [0.96879105  0.95061937  1.05113349]
```

## Batchnorm backward pass

Implement the backward pass for the batchnorm layer, `batchnorm_backward` in `nnd1/layers.py`. Check your implementation by running the following cell.

In [5]:

```
# Gradient check batchnorm backward pass

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

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

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

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

```
dx error:  7.556355825421427e-09
dgamma error:  2.7525915236159743e-11
dbeta error:  3.9609388546541704e-11
```

## 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]:

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

for reg in [0, 3.14]:
    print('Running check with reg = ', reg)
    model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
                              reg=reg, weight_scale=5e-2, dtype=np.float64,
                              use_batchnorm=True)

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

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

```
Running check with reg = 0
Initial loss: 2.3104140314139707
W1 relative error: 0.000248632509608278
W2 relative error: 0.0004157425164053145
W3 relative error: 3.665257379665477e-10
b1 relative error: 0.0
b2 relative error: 0.0
b3 relative error: 1.4141099806365967e-10
beta1 relative error: 2.686492586822503e-07
beta2 relative error: 2.2319749143096893e-07
gamma1 relative error: 2.269896599054521e-08
gamma2 relative error: 3.1768982937497974e-07
```

```
Running check with reg = 3.14
Initial loss: 7.491569681442019
W1 relative error: 0.0001541161954690109
W2 relative error: 1.9072431172817962e-06
W3 relative error: 7.429189804977281e-08
b1 relative error: 0.0
b2 relative error: 0.0
b3 relative error: 3.0534301145307134e-10
beta1 relative error: 1.8277810921742764e-08
beta2 relative error: 2.2571512254845303e-08
gamma1 relative error: 3.0257070278694645e-08
gamma2 relative error: 3.7687030175145544e-08
```

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

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

In [9]:

```
# Try training a very deep net with batchnorm
hidden_dims = [100, 100, 100, 100, 100]

num_train = 1000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

weight_scale = 2e-2
bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=True)
model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=False)

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

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

```
(Iteration 1 / 200) loss: 2.278569
(Epoch 0 / 10) train acc: 0.161000; val_acc: 0.121000
(Epoch 1 / 10) train acc: 0.325000; val_acc: 0.276000
(Epoch 2 / 10) train acc: 0.422000; val_acc: 0.320000
(Epoch 3 / 10) train acc: 0.513000; val_acc: 0.323000
(Epoch 4 / 10) train acc: 0.542000; val_acc: 0.309000
(Epoch 5 / 10) train acc: 0.598000; val_acc: 0.332000
(Epoch 6 / 10) train acc: 0.624000; val_acc: 0.314000
(Epoch 7 / 10) train acc: 0.672000; val_acc: 0.343000
(Epoch 8 / 10) train acc: 0.743000; val_acc: 0.339000
(Epoch 9 / 10) train acc: 0.782000; val_acc: 0.339000
(Epoch 10 / 10) train acc: 0.787000; val_acc: 0.338000
(Iteration 1 / 200) loss: 2.301447
(Epoch 0 / 10) train acc: 0.106000; val_acc: 0.116000
(Epoch 1 / 10) train acc: 0.235000; val_acc: 0.224000
(Epoch 2 / 10) train acc: 0.324000; val_acc: 0.280000
(Epoch 3 / 10) train acc: 0.346000; val_acc: 0.277000
(Epoch 4 / 10) train acc: 0.406000; val_acc: 0.316000
(Epoch 5 / 10) train acc: 0.474000; val_acc: 0.307000
(Epoch 6 / 10) train acc: 0.519000; val_acc: 0.316000
(Epoch 7 / 10) train acc: 0.529000; val_acc: 0.285000
(Epoch 8 / 10) train acc: 0.577000; val_acc: 0.308000
(Epoch 9 / 10) train acc: 0.609000; val_acc: 0.327000
(Epoch 10 / 10) train acc: 0.636000; val_acc: 0.326000
```

In [10]:

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

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

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

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

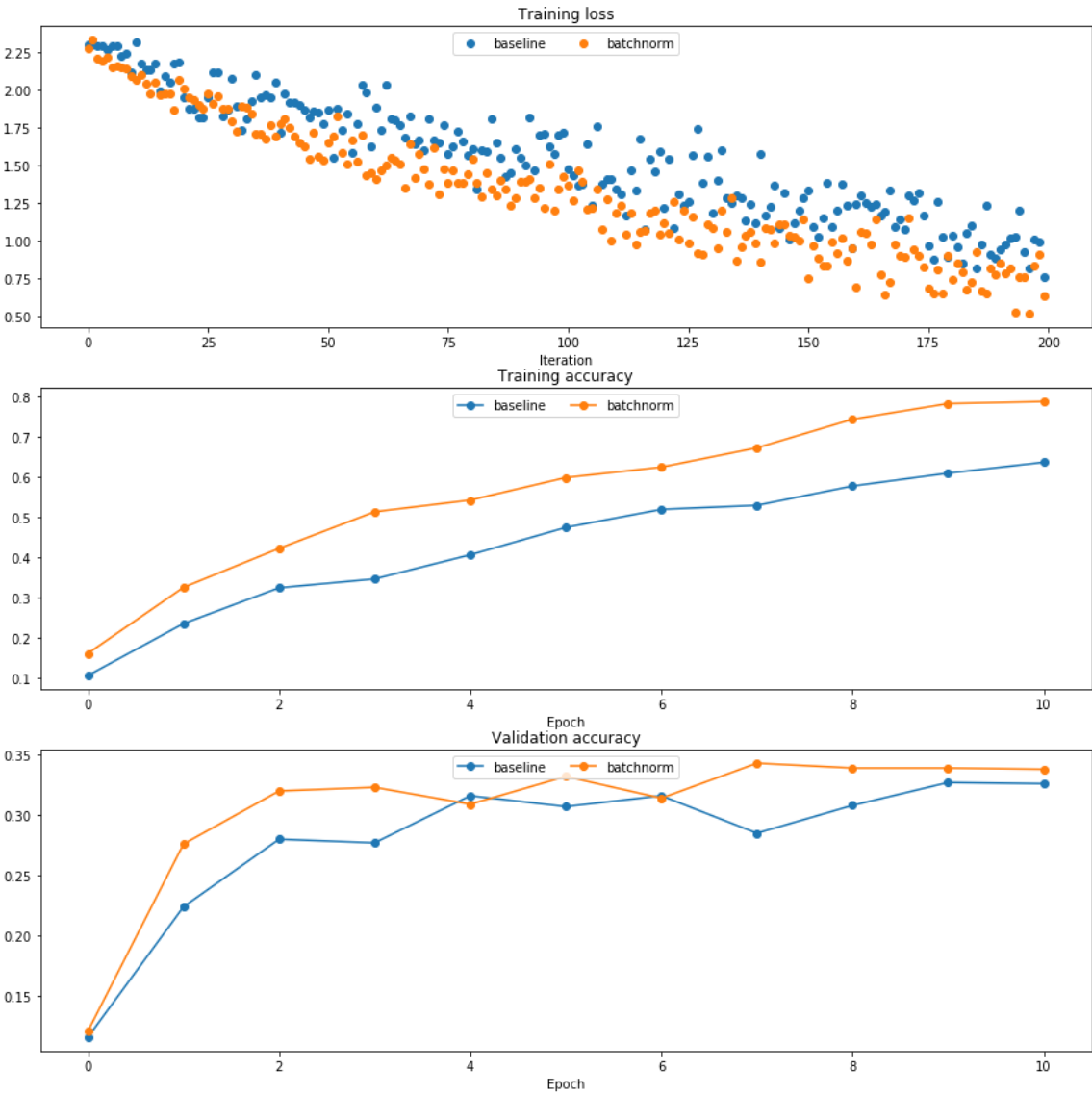
plt.subplot(3, 1, 2)
plt.plot(solver.train_acc_history, '-o', label='baseline')
plt.plot(bn_solver.train_acc_history, '-o', label='batchnorm')

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

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

```
c:\users\aksha\appdata\local\programs\python\python36\lib\site-packages\matplotlib\cbook\deprecation.py:106: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.  
warnings.warn(message, mplDeprecation, stacklevel=1)
```





## Batchnorm and initialization

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

In [11]:

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

num_train = 1000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

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

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

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

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

```
C:\Users\aksha\Documents\Winter 2018\EC239\FinalTest\HW4_code\code\nd1\layers.py:418: 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]:

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

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

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

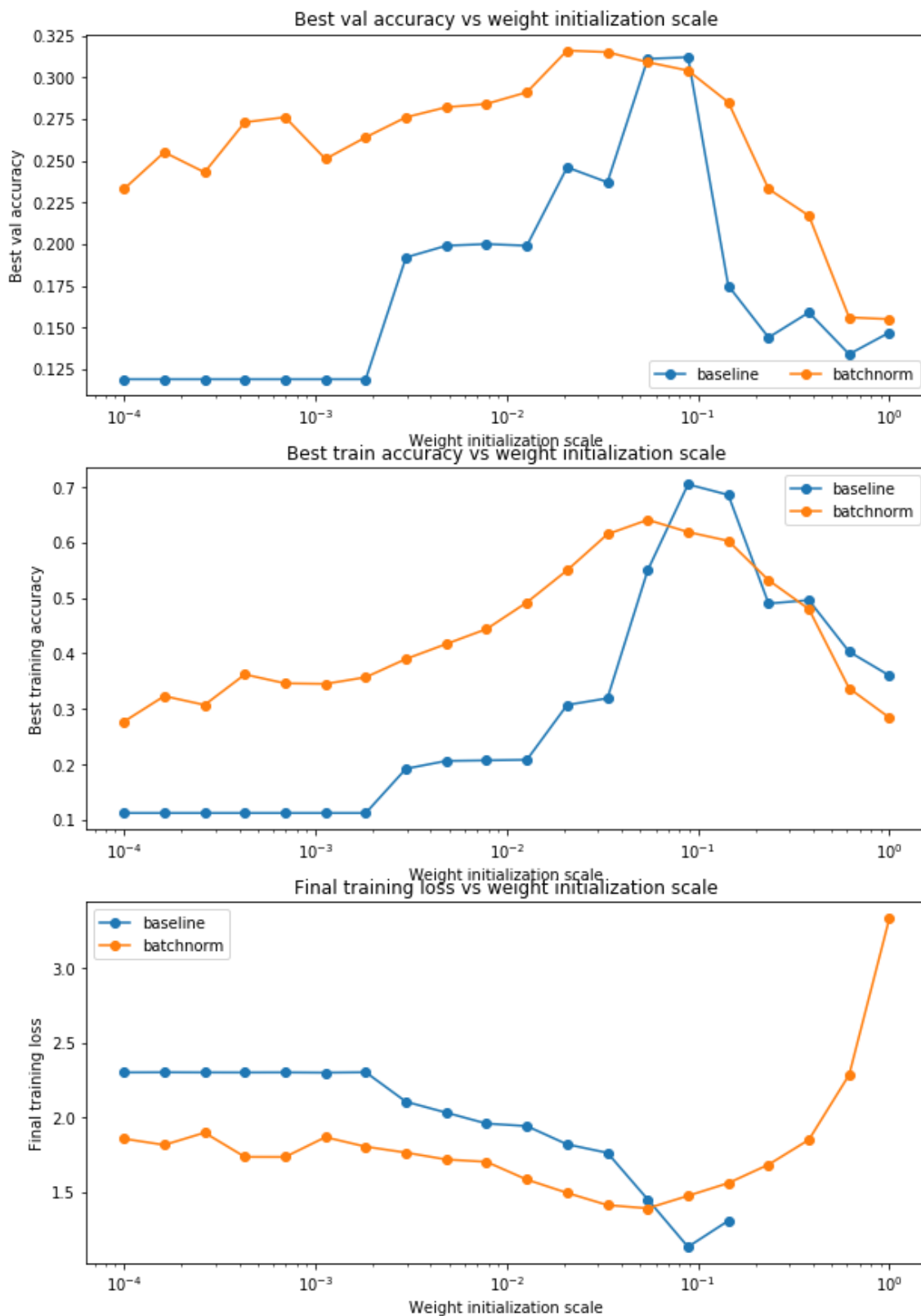
    final_train_loss.append(np.mean(solvers[ws].loss_history[-100:]))
    bn_final_train_loss.append(np.mean(bn_solvers[ws].loss_history[-100:]))

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

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

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

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



## Question:

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

## Answer:

Batchnorm is more robust for initializing weights compared to baseline where we don't normalize activations.

# Dropout ¶

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

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

In [1]:

```
## Import and setups

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

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

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

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

In [2]:

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

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

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

## Dropout forward pass

Implement the training and test time dropout forward pass, `dropout_forward`, in `nnd1/layers.py`. After that, test your implementation by running the following cell.

In [3]:

```
x = np.random.randn(500, 500) + 10

for p in [0.3, 0.6, 0.75]:
    out, _ = dropout_forward(x, {'mode': 'train', 'p': p})
    out_test, _ = dropout_forward(x, {'mode': 'test', 'p': p})

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

```
Running tests with p = 0.3
Mean of input: 10.001829294303864
Mean of train-time output: 9.973812507572998
Mean of test-time output: 10.001829294303864
Fraction of train-time output set to zero: 0.700748
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.6
Mean of input: 10.001829294303864
Mean of train-time output: 10.006583772085632
Mean of test-time output: 10.001829294303864
Fraction of train-time output set to zero: 0.399632
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.75
Mean of input: 10.001829294303864
Mean of train-time output: 10.006037676633047
Mean of test-time output: 10.001829294303864
Fraction of train-time output set to zero: 0.249652
Fraction of test-time output set to zero: 0.0
```

## Dropout backward pass

Implement the backward pass, `dropout_backward`, in `nnd1/layers.py`. After that, test your gradients by running the following cell:



In [4]:

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

dropout_param = {'mode': 'train', 'p': 0.8, 'seed': 123}
out, cache = dropout_forward(x, dropout_param)
dx = dropout_backward(dout, cache)
dx_num = eval_numerical_gradient_array(lambda xx: dropout_forward(xx, dropout_param)[0], x, dout)

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

dx relative error:  5.4456100419898414e-11
```

## Implement a fully connected neural network with dropout layers

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

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

Check your implementation by running the following code. Our  $W_1$  gradient relative error is on the order of  $1e-6$  (the largest of all the relative errors).

In [5]:

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

for dropout in [0, 0.25, 0.5]:
    print('Running check with dropout = ', dropout)
    model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
                              weight_scale=5e-2, dtype=np.float64,
                              dropout=dropout, seed=123)

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

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

```
Running check with dropout = 0
Initial loss: 2.3051948273987857
W1 relative error: 2.5272575344376073e-07
W2 relative error: 1.5034484929313676e-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: 2.3126468345657742
W1 relative error: 1.483854795975875e-08
W2 relative error: 2.3427832149940254e-10
W3 relative error: 3.564454999162522e-08
b1 relative error: 1.5292167232408546e-09
b2 relative error: 1.842268868410678e-10
b3 relative error: 1.4026015558098908e-10
```

```
Running check with dropout = 0.5
Initial loss: 2.302437587710995
W1 relative error: 4.553387957138422e-08
W2 relative error: 2.974218050584597e-08
W3 relative error: 4.3413247403122424e-07
b1 relative error: 1.872462967441693e-08
b2 relative error: 5.045591219274328e-09
b3 relative error: 8.009887154529434e-11
```

## Dropout as a regularizer

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

In [6]:

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

num_train = 500
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

solvers = {}
dropout_choices = [0, 0.6]
for dropout in dropout_choices:
    model = FullyConnectedNet([100, 100, 100], dropout=dropout)

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

```
(Iteration 1 / 125) loss: 2.300804
(Epoch 0 / 25) train acc: 0.220000; val_acc: 0.168000
(Epoch 1 / 25) train acc: 0.188000; val_acc: 0.147000
(Epoch 2 / 25) train acc: 0.266000; val_acc: 0.200000
(Epoch 3 / 25) train acc: 0.338000; val_acc: 0.262000
(Epoch 4 / 25) train acc: 0.378000; val_acc: 0.278000
(Epoch 5 / 25) train acc: 0.428000; val_acc: 0.297000
(Epoch 6 / 25) train acc: 0.468000; val_acc: 0.323000
(Epoch 7 / 25) train acc: 0.494000; val_acc: 0.287000
(Epoch 8 / 25) train acc: 0.566000; val_acc: 0.328000
(Epoch 9 / 25) train acc: 0.572000; val_acc: 0.322000
(Epoch 10 / 25) train acc: 0.622000; val_acc: 0.324000
(Epoch 11 / 25) train acc: 0.670000; val_acc: 0.279000
(Epoch 12 / 25) train acc: 0.710000; val_acc: 0.338000
(Epoch 13 / 25) train acc: 0.746000; val_acc: 0.319000
(Epoch 14 / 25) train acc: 0.792000; val_acc: 0.307000
(Epoch 15 / 25) train acc: 0.834000; val_acc: 0.297000
(Epoch 16 / 25) train acc: 0.876000; val_acc: 0.327000
(Epoch 17 / 25) train acc: 0.886000; val_acc: 0.320000
(Epoch 18 / 25) train acc: 0.918000; val_acc: 0.314000
(Epoch 19 / 25) train acc: 0.922000; val_acc: 0.290000
(Epoch 20 / 25) train acc: 0.944000; val_acc: 0.306000
(Iteration 101 / 125) loss: 0.156105
(Epoch 21 / 25) train acc: 0.968000; val_acc: 0.302000
(Epoch 22 / 25) train acc: 0.978000; val_acc: 0.302000
(Epoch 23 / 25) train acc: 0.976000; val_acc: 0.289000
(Epoch 24 / 25) train acc: 0.986000; val_acc: 0.285000
(Epoch 25 / 25) train acc: 0.978000; val_acc: 0.311000
(Iteration 1 / 125) loss: 2.301328
(Epoch 0 / 25) train acc: 0.154000; val_acc: 0.143000
(Epoch 1 / 25) train acc: 0.214000; val_acc: 0.195000
(Epoch 2 / 25) train acc: 0.252000; val_acc: 0.217000
(Epoch 3 / 25) train acc: 0.276000; val_acc: 0.200000
(Epoch 4 / 25) train acc: 0.308000; val_acc: 0.254000
(Epoch 5 / 25) train acc: 0.316000; val_acc: 0.241000
(Epoch 6 / 25) train acc: 0.322000; val_acc: 0.282000
(Epoch 7 / 25) train acc: 0.354000; val_acc: 0.273000
(Epoch 8 / 25) train acc: 0.364000; val_acc: 0.276000
(Epoch 9 / 25) train acc: 0.408000; val_acc: 0.282000
(Epoch 10 / 25) train acc: 0.454000; val_acc: 0.302000
(Epoch 11 / 25) train acc: 0.472000; val_acc: 0.297000
(Epoch 12 / 25) train acc: 0.498000; val_acc: 0.318000
(Epoch 13 / 25) train acc: 0.510000; val_acc: 0.309000
(Epoch 14 / 25) train acc: 0.534000; val_acc: 0.315000
(Epoch 15 / 25) train acc: 0.546000; val_acc: 0.331000
(Epoch 16 / 25) train acc: 0.584000; val_acc: 0.302000
(Epoch 17 / 25) train acc: 0.626000; val_acc: 0.332000
(Epoch 18 / 25) train acc: 0.614000; val_acc: 0.327000
(Epoch 19 / 25) train acc: 0.626000; val_acc: 0.325000
(Epoch 20 / 25) train acc: 0.656000; val_acc: 0.338000
(Iteration 101 / 125) loss: 1.299273
(Epoch 21 / 25) train acc: 0.676000; val_acc: 0.329000
(Epoch 22 / 25) train acc: 0.682000; val_acc: 0.324000
(Epoch 23 / 25) train acc: 0.730000; val_acc: 0.344000
(Epoch 24 / 25) train acc: 0.740000; val_acc: 0.321000
(Epoch 25 / 25) train acc: 0.770000; val_acc: 0.332000
```

In [7]:

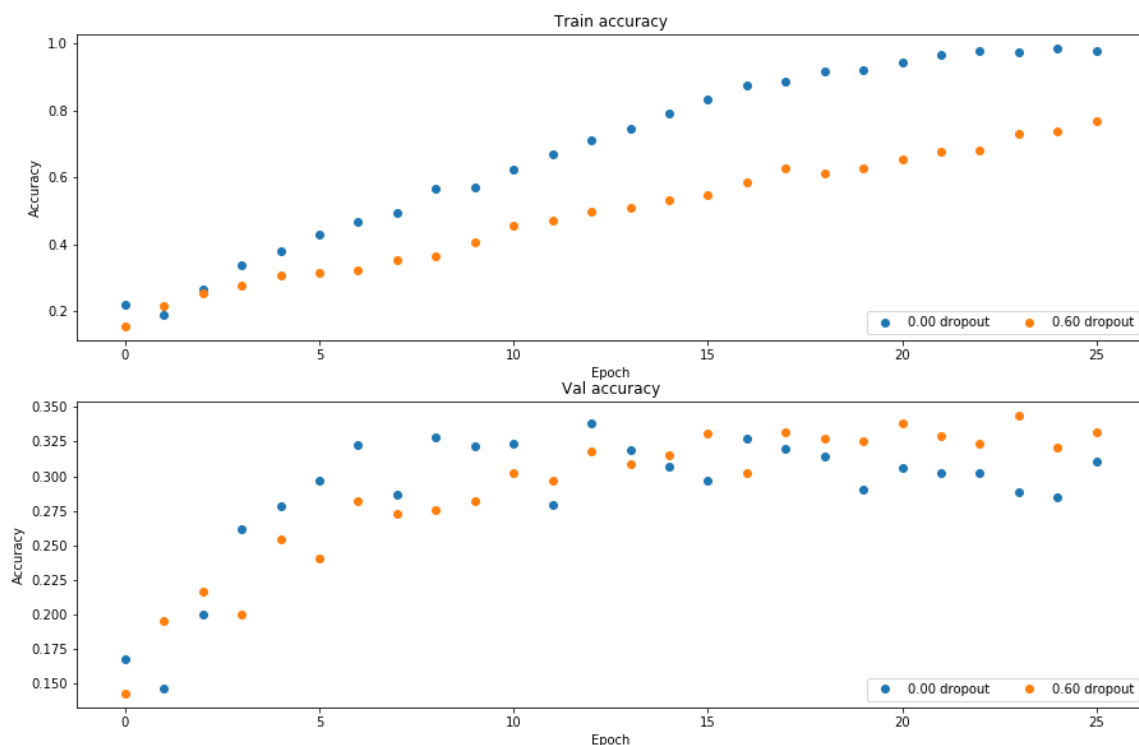
```
# Plot train and validation accuracies of the two models

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

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

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

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



## Question

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

## Answer:

### Final part of the assignment

Get over 60% validation accuracy on CIFAR-10 by using the layers you have implemented. You will be graded according to the following equation:

$\min(\text{floor}((X - 32\%) / 28\%, 1)$  where if you get 60% or higher validation accuracy, you get full points.

In [10]:

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

model = FullyConnectedNet([1000, 1000, 1000, 1000], dropout=0.65, use_batchnorm=True)

solver = Solver(model, data, num_epochs=25, batch_size=305,
                update_rule = 'adam',
                optim_config={'learning_rate': 1e-3,}, lr_decay=0.9, verbose=True, print_
every=100)

solver.train()

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

```
(Iteration 1 / 4000) loss: 2.322243
(Epoch 0 / 25) train acc: 0.245000; val_acc: 0.220000
(Iteration 101 / 4000) loss: 1.522808
(Epoch 1 / 25) train acc: 0.506000; val_acc: 0.475000
(Iteration 201 / 4000) loss: 1.418678
(Iteration 301 / 4000) loss: 1.303142
(Epoch 2 / 25) train acc: 0.514000; val_acc: 0.517000
(Iteration 401 / 4000) loss: 1.298947
(Epoch 3 / 25) train acc: 0.579000; val_acc: 0.539000
(Iteration 501 / 4000) loss: 1.218159
(Iteration 601 / 4000) loss: 1.209194
(Epoch 4 / 25) train acc: 0.593000; val_acc: 0.544000
(Iteration 701 / 4000) loss: 1.086574
(Epoch 5 / 25) train acc: 0.614000; val_acc: 0.524000
(Iteration 801 / 4000) loss: 1.238556
(Iteration 901 / 4000) loss: 1.094975
(Epoch 6 / 25) train acc: 0.657000; val_acc: 0.556000
(Iteration 1001 / 4000) loss: 1.072980
(Iteration 1101 / 4000) loss: 0.882316
(Epoch 7 / 25) train acc: 0.658000; val_acc: 0.562000
(Iteration 1201 / 4000) loss: 0.970731
(Epoch 8 / 25) train acc: 0.683000; val_acc: 0.565000
(Iteration 1301 / 4000) loss: 1.148996
(Iteration 1401 / 4000) loss: 1.062047
(Epoch 9 / 25) train acc: 0.696000; val_acc: 0.569000
(Iteration 1501 / 4000) loss: 1.031162
(Epoch 10 / 25) train acc: 0.740000; val_acc: 0.587000
(Iteration 1601 / 4000) loss: 0.940797
(Iteration 1701 / 4000) loss: 0.940836
(Epoch 11 / 25) train acc: 0.735000; val_acc: 0.584000
(Iteration 1801 / 4000) loss: 0.872584
(Iteration 1901 / 4000) loss: 0.815859
(Epoch 12 / 25) train acc: 0.765000; val_acc: 0.587000
(Iteration 2001 / 4000) loss: 0.760623
(Epoch 13 / 25) train acc: 0.750000; val_acc: 0.567000
(Iteration 2101 / 4000) loss: 0.830717
(Iteration 2201 / 4000) loss: 0.723647
(Epoch 14 / 25) train acc: 0.774000; val_acc: 0.592000
(Iteration 2301 / 4000) loss: 0.921479
(Epoch 15 / 25) train acc: 0.779000; val_acc: 0.584000
(Iteration 2401 / 4000) loss: 0.889501
(Iteration 2501 / 4000) loss: 0.726995
(Epoch 16 / 25) train acc: 0.793000; val_acc: 0.598000
(Iteration 2601 / 4000) loss: 0.720859
(Iteration 2701 / 4000) loss: 0.906341
(Epoch 17 / 25) train acc: 0.779000; val_acc: 0.592000
(Iteration 2801 / 4000) loss: 0.709677
(Epoch 18 / 25) train acc: 0.831000; val_acc: 0.595000
(Iteration 2901 / 4000) loss: 0.692159
(Iteration 3001 / 4000) loss: 0.725254
(Epoch 19 / 25) train acc: 0.828000; val_acc: 0.581000
(Iteration 3101 / 4000) loss: 0.719448
(Epoch 20 / 25) train acc: 0.840000; val_acc: 0.591000
(Iteration 3201 / 4000) loss: 0.708727
(Iteration 3301 / 4000) loss: 0.572546
(Epoch 21 / 25) train acc: 0.861000; val_acc: 0.590000
(Iteration 3401 / 4000) loss: 0.745503
(Iteration 3501 / 4000) loss: 0.605202
(Epoch 22 / 25) train acc: 0.865000; val_acc: 0.576000
(Iteration 3601 / 4000) loss: 0.712708
(Epoch 23 / 25) train acc: 0.875000; val_acc: 0.582000
```



```
(Iteration 3701 / 4000) loss: 0.700613  
(Iteration 3801 / 4000) loss: 0.712505  
(Epoch 24 / 25) train acc: 0.876000; val_acc: 0.583000  
(Iteration 3901 / 4000) loss: 0.598972  
(Epoch 25 / 25) train acc: 0.863000; val_acc: 0.580000
```

```

1  import numpy as np
2  import pdb
3
4  """
5  This code was originally written for CS 231n at Stanford University
6  (cs231n.stanford.edu). It has been modified in various areas for use in the
7  ECE 239AS class at UCLA. This includes the descriptions of what code to
8  implement as well as some slight potential changes in variable names to be
9  consistent with class nomenclature. We thank Justin Johnson & Serena Yeung for
10 permission to use this code. To see the original version, please visit
11 cs231n.stanford.edu.
12 """
13
14 def affine_forward(x, w, b):
15     """
16     Computes the forward pass for an affine (fully-connected) layer.
17
18     The input x has shape (N, d_1, ..., d_k) and contains a minibatch of N
19     examples, where each example x[i] has shape (d_1, ..., d_k). We will
20     reshape each input into a vector of dimension D = d_1 * ... * d_k, and
21     then transform it to an output vector of dimension M.
22
23     Inputs:
24     - x: A numpy array containing input data, of shape (N, d_1, ..., d_k)
25     - w: A numpy array of weights, of shape (D, M)
26     - b: A numpy array of biases, of shape (M,)
27
28     Returns a tuple of:
29     - out: output, of shape (N, M)
30     - cache: (x, w, b)
31     """
32
33     # ===== #
34     # YOUR CODE HERE:
35     # Calculate the output of the forward pass. Notice the dimensions
36     # of w are D x M, which is the transpose of what we did in earlier
37     # assignments.
38     # ===== #
39
40     out = x.reshape(x.shape[0], w.shape[0]).dot(w) + b
41
42     # ===== #
43     # END YOUR CODE HERE
44     # ===== #
45
46     cache = (x, w, b)
47     return out, cache
48
49
50 def affine_backward(dout, cache):
51     """
52     Computes the backward pass for an affine layer.
53
54     Inputs:
55     - dout: Upstream derivative, of shape (N, M)
56     - cache: Tuple of:
57       - x: Input data, of shape (N, d_1, ... d_k)
58       - w: Weights, of shape (D, M)
59
60     Returns a tuple of:
61     - dx: Gradient with respect to x, of shape (N, d1, ..., d_k)
62     - dw: Gradient with respect to w, of shape (D, M)
63     - db: Gradient with respect to b, of shape (M,)
64     """
65     x, w, b = cache
66     dx, dw, db = None, None, None
67

```

```

68 # ===== #
69 # YOUR CODE HERE:
70 #   Calculate the gradients for the backward pass.
71 # ===== #
72
73 dx = dout.dot(w.T).reshape(x.shape)
74 db = np.sum(dout, axis=0)
75 dw = x.reshape(x.shape[0], w.shape[0]).T.dot(dout)
76
77 # ===== #
78 # END YOUR CODE HERE
79 # ===== #
80
81 return dx, dw, db
82
83 def relu_forward(x):
84     """
85     Computes the forward pass for a layer of rectified linear units (ReLU).
86
87     Input:
88     - x: Inputs, of any shape
89
90     Returns a tuple of:
91     - out: Output, of the same shape as x
92     - cache: x
93     """
94     # ===== #
95     # YOUR CODE HERE:
96     #   Implement the ReLU forward pass.
97     # ===== #
98
99     out = np.maximum(0, x)
100
101     # ===== #
102     # END YOUR CODE HERE
103     # ===== #
104
105     cache = x
106     return out, cache
107
108
109 def relu_backward(dout, cache):
110     """
111     Computes the backward pass for a layer of rectified linear units (ReLU).
112
113     Input:
114     - dout: Upstream derivatives, of any shape
115     - cache: Input x, of same shape as dout
116
117     Returns:
118     - dx: Gradient with respect to x
119     """
120     x = cache
121
122     # ===== #
123     # YOUR CODE HERE:
124     #   Implement the ReLU backward pass
125     # ===== #
126
127     dx = dout
128     dx[cache < 0] = 0
129
130     # ===== #
131     # END YOUR CODE HERE
132     # ===== #
133
134     return dx

```

```

135 def batchnorm_forward(x, gamma, beta, bn_param):
136     """
137     Forward pass for batch normalization.
138
139     During training the sample mean and (uncorrected) sample variance are
140     computed from minibatch statistics and used to normalize the incoming data.
141     During training we also keep an exponentially decaying running mean of the mean
142     and variance of each feature, and these averages are used to normalize data
143     at test-time.
144
145     At each timestep we update the running averages for mean and variance using
146     an exponential decay based on the momentum parameter:
147
148     running_mean = momentum * running_mean + (1 - momentum) * sample_mean
149     running_var = momentum * running_var + (1 - momentum) * sample_var
150
151     Note that the batch normalization paper suggests a different test-time
152     behavior: they compute sample mean and variance for each feature using a
153     large number of training images rather than using a running average. For
154     this implementation we have chosen to use running averages instead since
155     they do not require an additional estimation step; the torch7 implementation
156     of batch normalization also uses running averages.
157
158     Input:
159     - x: Data of shape (N, D)
160     - gamma: Scale parameter of shape (D,)
161     - beta: Shift parameter of shape (D,)
162     - bn_param: Dictionary with the following keys:
163       - mode: 'train' or 'test'; required
164       - eps: Constant for numeric stability
165       - momentum: Constant for running mean / variance.
166       - running_mean: Array of shape (D,) giving running mean of features
167       - running_var: Array of shape (D,) giving running variance of features
168
169     Returns a tuple of:
170     - out: of shape (N, D)
171     - cache: A tuple of values needed in the backward pass
172     """
173     mode = bn_param['mode']
174     eps = bn_param.get('eps', 1e-5)
175     momentum = bn_param.get('momentum', 0.9)
176
177     N, D = x.shape
178     running_mean = bn_param.get('running_mean', np.zeros(D, dtype=x.dtype))
179     running_var = bn_param.get('running_var', np.zeros(D, dtype=x.dtype))
180
181     out, cache = None, None
182     if mode == 'train':
183
184         # ===== #
185         # YOUR CODE HERE:
186         #   A few steps here:
187         #   (1) Calculate the running mean and variance of the minibatch.
188         #   (2) Normalize the activations with the batch mean and variance.
189         #   (3) Scale and shift the normalized activations. Store this
190         #       as the variable 'out'
191         #   (4) Store any variables you may need for the backward pass in
192         #       the 'cache' variable.
193         # ===== #
194
195         mean = np.mean(x, axis=0)
196         var = np.var(x, axis=0)
197
198         out = (x - mean) / np.sqrt(var + eps) * gamma + beta
199         cache = (x, mean, var, gamma, beta, eps)
200
201         running_mean = momentum * running_mean + (1 - momentum) * mean

```

```

202     running_var = momentum * running_var + (1 - momentum) * var
203
204     # ===== #
205     # END YOUR CODE HERE
206     # ===== #
207
208     elif mode == 'test':
209
210         # ===== #
211         # YOUR CODE HERE:
212         # Calculate the testing time normalized activations. Normalize using
213         # the running mean and variance, and then scale and shift appropriately.
214         # Store the output as 'out'.
215         # ===== #
216
217         mean = running_mean
218         variance = running_var
219
220         out = (x - mean)/(np.sqrt(variance + eps))* gamma + beta
221
222         # ===== #
223         # END YOUR CODE HERE
224         # ===== #
225
226     else:
227         raise ValueError('Invalid forward batchnorm mode "%s"' % mode)
228
229     # Store the updated running means back into bn_param
230     bn_param['running_mean'] = running_mean
231     bn_param['running_var'] = running_var
232
233     return out, cache
234
235 def batchnorm_backward(dout, cache):
236     """
237     Backward pass for batch normalization.
238
239     For this implementation, you should write out a computation graph for
240     batch normalization on paper and propagate gradients backward through
241     intermediate nodes.
242
243     Inputs:
244     - dout: Upstream derivatives, of shape (N, D)
245     - cache: Variable of intermediates from batchnorm_forward.
246
247     Returns a tuple of:
248     - dx: Gradient with respect to inputs x, of shape (N, D)
249     - dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
250     - dbeta: Gradient with respect to shift parameter beta, of shape (D,)
251     """
252     dx, dgamma, dbeta = None, None, None
253
254     # ===== #
255     # YOUR CODE HERE:
256     # Implement the batchnorm backward pass, calculating dx, dgamma, and dbeta.
257     # ===== #
258     N, D = dout.shape
259
260     x, x_mean, x_var, gamma, beta, eps = cache
261
262     dgamma = np.sum(dout * (x - x_mean) / np.sqrt(x_var + eps), axis=0)
263     dbeta = np.sum(dout, axis=0)
264     ddiv = dout * gamma
265
266
267     dnumerator = ddiv / np.sqrt(x_var + eps)
268     ddenominator = - ddiv * (x - x_mean) / (x_var + eps)

```

```

269 dxmean = -dnumerator
270
271 dx = (1.0/N) * np.ones((N, N)).dot(dxmean)
272 dx += dnumerator
273 dvar = 0.5 * ((x_var + eps)**(-0.5)) * ddenominator
274 dtemp = 2 * (x - x_mean) * ((1.0/N) * np.ones((N, N)).dot(dvar))
275
276 dx += dtemp
277 dx += -np.mean(dtemp, axis=0)
278
279 # ===== #
280 # END YOUR CODE HERE
281 # ===== #
282
283 return dx, dgamma, dbeta
284
285 def dropout_forward(x, dropout_param):
286     """
287     Performs the forward pass for (inverted) dropout.
288
289     Inputs:
290     - x: Input data, of any shape
291     - dropout_param: A dictionary with the following keys:
292       - p: Dropout parameter. We drop each neuron output with probability p.
293       - mode: 'test' or 'train'. If the mode is train, then perform dropout;
294         if the mode is test, then just return the input.
295       - seed: Seed for the random number generator. Passing seed makes this
296         function deterministic, which is needed for gradient checking but not in
297         real networks.
298
299     Outputs:
300     - out: Array of the same shape as x.
301     - cache: A tuple (dropout_param, mask). In training mode, mask is the dropout
302       mask that was used to multiply the input; in test mode, mask is None.
303     """
304     p, mode = dropout_param['p'], dropout_param['mode']
305     if 'seed' in dropout_param:
306         np.random.seed(dropout_param['seed'])
307
308     mask = None
309     out = None
310
311     if mode == 'train':
312         # ===== #
313         # YOUR CODE HERE:
314         # Implement the inverted dropout forward pass during training time.
315         # Store the masked and scaled activations in out, and store the
316         # dropout mask as the variable mask.
317         # ===== #
318
319         mask = (np.random.rand(*x.shape)<p) / p
320         out = x*mask
321
322         # ===== #
323         # END YOUR CODE HERE
324         # ===== #
325
326     elif mode == 'test':
327
328         # ===== #
329         # YOUR CODE HERE:
330         # Implement the inverted dropout forward pass during test time.
331         # ===== #
332
333         out = x
334         # ===== #
335         # END YOUR CODE HERE

```

```

336     # ===== #
337
338     cache = (dropout_param, mask)
339     out = out.astype(x.dtype, copy=False)
340
341     return out, cache
342
343 def dropout_backward(dout, cache):
344     """
345     Perform the backward pass for (inverted) dropout.
346
347     Inputs:
348     - dout: Upstream derivatives, of any shape
349     - cache: (dropout_param, mask) from dropout_forward.
350     """
351     dropout_param, mask = cache
352     mode = dropout_param['mode']
353
354     dx = None
355     if mode == 'train':
356         # ===== #
357         # YOUR CODE HERE:
358         #   Implement the inverted dropout backward pass during training time.
359         # ===== #
360         dx = dout * mask
361         # ===== #
362         # END YOUR CODE HERE
363         # ===== #
364     elif mode == 'test':
365         # ===== #
366         # YOUR CODE HERE:
367         #   Implement the inverted dropout backward pass during test time.
368         # ===== #
369         dx = dout
370         # ===== #
371         # END YOUR CODE HERE
372         # ===== #
373     return dx
374
375 def svm_loss(x, y):
376     """
377     Computes the loss and gradient using for multiclass SVM classification.
378
379     Inputs:
380     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
381         for the ith input.
382     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
383         0 <= y[i] < C
384
385     Returns a tuple of:
386     - loss: Scalar giving the loss
387     - dx: Gradient of the loss with respect to x
388     """
389     N = x.shape[0]
390     correct_class_scores = x[np.arange(N), y]
391     margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
392     margins[np.arange(N), y] = 0
393     loss = np.sum(margins) / N
394     num_pos = np.sum(margins > 0, axis=1)
395     dx = np.zeros_like(x)
396     dx[margins > 0] = 1
397     dx[np.arange(N), y] -= num_pos
398     dx /= N
399     return loss, dx
400
401
402 def softmax_loss(x, y):

```

```

403 """
404 Computes the loss and gradient for softmax classification.
405
406 Inputs:
407 - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
408     for the ith input.
409 - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
410     0 <= y[i] < C
411
412 Returns a tuple of:
413 - loss: Scalar giving the loss
414 - dx: Gradient of the loss with respect to x
415 """
416
417 probs = np.exp(x - np.max(x, axis=1, keepdims=True))
418 probs /= np.sum(probs, axis=1, keepdims=True)
419 N = x.shape[0]
420 loss = -np.sum(np.log(probs[np.arange(N), y])) / N
421 dx = probs.copy()
422 dx[np.arange(N), y] -= 1
423 dx /= N
424 return loss, dx
425

```



```

1  from .layers import *
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 def affine_relu_forward(x, w, b):
14     """
15     Convenience layer that performs an affine transform followed by a ReLU
16
17     Inputs:
18     - x: Input to the affine layer
19     - w, b: Weights for the affine layer
20
21     Returns a tuple of:
22     - out: Output from the ReLU
23     - cache: Object to give to the backward pass
24     """
25     a, fc_cache = affine_forward(x, w, b)
26     out, relu_cache = relu_forward(a)
27     cache = (fc_cache, relu_cache)
28     return out, cache
29
30 def affine_relu_backward(dout, cache):
31     """
32     Backward pass for the affine-relu convenience layer
33     """
34     fc_cache, relu_cache = cache
35     da = relu_backward(dout, relu_cache)
36     dx, dw, db = affine_backward(da, fc_cache)
37     return dx, dw, db
38
39 def affine_batchnorm_relu(X, w, b, gamma, beta, bn_param):
40     a, fc_cache = affine_forward(X, w, b)
41     batch_out, batch_cache = batchnorm_forward(a, gamma, beta, bn_param)
42     out, relu_cache = relu_forward(batch_out)
43     cache = (fc_cache, batch_cache, relu_cache)
44     return out, cache
45
46 def affine_batchnorm_relu_backwards(dout, cache):
47     fc_cache, batch_cache, relu_cache = cache
48     da = relu_backward(dout, relu_cache)
49     dy, dgamma, dbeta = batchnorm_backward(da, batch_cache)
50     dx, dw, db = affine_backward(dy, fc_cache)
51     return dx, dw, db, dgamma, dbeta

```

```

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     next_w = None
74     # ===== #
75     # YOUR CODE HERE:
76     # Implement the momentum update formula. Return the updated weights
77     # as next_w, and store the updated velocity as v.
78     # ===== #
79
80     v = config['momentum'] * v - config['learning_rate'] * dw
81     next_w = w + v
82
83     # ===== #
84     # END YOUR CODE HERE
85     # ===== #
86
87     config['velocity'] = v
88
89     return next_w, config
90
91 def sgd_nesterov_momentum(w, dw, config=None):
92     """
93     Performs stochastic gradient descent with Nesterov momentum.
94
95     config format:
96     - learning_rate: Scalar learning rate.
97     - momentum: Scalar between 0 and 1 giving the momentum value.
98       Setting momentum = 0 reduces to sgd.
99     - velocity: A numpy array of the same shape as w and dw used to store a moving
100       average of the gradients.
101     """
102     if config is None: config = {}
103     config.setdefault('learning_rate', 1e-2)
104     config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
105     v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.
106
107     # ===== #
108     # YOUR CODE HERE:
109     # Implement the momentum update formula. Return the updated weights
110     # as next_w, and store the updated velocity as v.
111     # ===== #
112     v_prev = v # back this up
113     v = config['momentum'] * v - config['learning_rate'] * dw # velocity update stays the
114     same
115     next_w = w + v + config['momentum'] * (v - v_prev) # position update changes form
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.

```

```

134 - beta: Moving average of second moments of gradients.
135 """
136 if config is None: config = {}
137 config.setdefault('learning_rate', 1e-2)
138 config.setdefault('decay_rate', 0.99)
139 config.setdefault('epsilon', 1e-8)
140 config.setdefault('a', np.zeros_like(w))
141
142 next_w = None
143
144 # ===== #
145 # YOUR CODE HERE:
146 # Implement RMSProp. Store the next value of w as next_w. You need
147 # to also store in config['a'] the moving average of the second
148 # moment gradients, so they can be used for future gradients. Concretely,
149 # config['a'] corresponds to "a" in the lecture notes.
150 # ===== #
151
152 cache, decay_rate, eps, learning_rate = config['a'], config['decay_rate'], config[
153 'epsilon'], config['learning_rate']
154 cache = decay_rate * cache + (1 - decay_rate) * dw**2
155 next_w = -learning_rate * dw / (np.sqrt(cache) + eps) + w
156 config['a'] = cache
157 # ===== #
158 # END YOUR CODE HERE
159 # ===== #
160
161 return next_w, config
162
163 def adam(w, dw, config=None):
164     """
165     Uses the Adam update rule, which incorporates moving averages of both the
166     gradient and its square and a bias correction term.
167
168     config format:
169     - learning_rate: Scalar learning rate.
170     - beta1: Decay rate for moving average of first moment of gradient.
171     - beta2: Decay rate for moving average of second moment of gradient.
172     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
173     - m: Moving average of gradient.
174     - v: Moving average of squared gradient.
175     - t: Iteration number.
176     """
177     if config is None: config = {}
178     config.setdefault('learning_rate', 1e-3)
179     config.setdefault('beta1', 0.9)
180     config.setdefault('beta2', 0.999)
181     config.setdefault('epsilon', 1e-8)
182     config.setdefault('v', np.zeros_like(w))
183     config.setdefault('a', np.zeros_like(w))
184     config.setdefault('t', 0)
185
186     next_w = None
187
188     # ===== #
189     # YOUR CODE HERE:
190     # Implement Adam. Store the next value of w as next_w. You need
191     # to also store in config['a'] the moving average of the second
192     # moment gradients, and in config['v'] the moving average of the
193     # first moments. Finally, store in config['t'] the increasing time.
194     # ===== #
195
196     learning_rate, beta1, beta2, eps, a, v, t = config['learning_rate'], config['beta1'],
197     config['beta2'], config['epsilon'], config['a'], config['v'], config['t']
198     t = t + 1

```

```
199 v = beta1* v + (1-beta1)*dw
200 a = beta2*a + (1-beta2)* (dw**2)
201
202
203 ab = a / (1 - beta2**t)
204 vb = v / (1 - beta1**t)
205
206 next_w = -learning_rate * vb / (np.sqrt(ab) + eps) + w
207
208 config['a'], config['v'], config['t'] = a, v, t
209
210 # ===== #
211 # END YOUR CODE HERE
212 # ===== #
213
214 return next_w, config
215
216
217
218
219
220
```

```

1  import numpy as np
2  import pdb
3
4  from .layers import *
5  from .layer_utils import *
6
7  """
8  This code was originally written for CS 231n at Stanford University
9  (cs231n.stanford.edu). It has been modified in various areas for use in the
10 ECE 239AS class at UCLA. This includes the descriptions of what code to
11 implement as well as some slight potential changes in variable names to be
12 consistent with class nomenclature. We thank Justin Johnson & Serena Yeung for
13 permission to use this code. To see the original version, please visit
14 cs231n.stanford.edu.
15 """
16
17 class FullyConnectedNet(object):
18     """
19     A fully-connected neural network with an arbitrary number of hidden layers,
20     ReLU nonlinearities, and a softmax loss function. This will also implement
21     dropout and batch normalization as options. For a network with L layers,
22     the architecture will be
23
24     {affine - [batch norm] - relu - [dropout]} x (L - 1) - affine - softmax
25
26     where batch normalization and dropout are optional, and the {...} block is
27     repeated L - 1 times.
28
29     Similar to the TwoLayerNet above, learnable parameters are stored in the
30     self.params dictionary and will be learned using the Solver class.
31     """
32
33     def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,
34                 dropout=0, use_batchnorm=False, reg=0.0,
35                 weight_scale=1e-2, dtype=np.float32, seed=None):
36         """
37         Initialize a new FullyConnectedNet.
38
39         Inputs:
40         - hidden_dims: A list of integers giving the size of each hidden layer.
41         - input_dim: An integer giving the size of the input.
42         - num_classes: An integer giving the number of classes to classify.
43         - dropout: Scalar between 0 and 1 giving dropout strength. If dropout=0 then
44           the network should not use dropout at all.
45         - use_batchnorm: Whether or not the network should use batch normalization.
46         - reg: Scalar giving L2 regularization strength.
47         - weight_scale: Scalar giving the standard deviation for random
48           initialization of the weights.
49         - dtype: A numpy datatype object; all computations will be performed using
50           this datatype. float32 is faster but less accurate, so you should use
51           float64 for numeric gradient checking.
52         - seed: If not None, then pass this random seed to the dropout layers. This
53           will make the dropout layers deterministic so we can gradient check the
54           model.
55         """
56         self.use_batchnorm = use_batchnorm
57         self.use_dropout = dropout > 0
58         self.reg = reg
59         self.num_layers = 1 + len(hidden_dims)
60         self.dtype = dtype
61         self.params = {}
62
63         # ===== #
64         # YOUR CODE HERE:
65         # Initialize all parameters of the network in the self.params dictionary.
66         # The weights and biases of layer 1 are W1 and b1; and in general the
67         # weights and biases of layer i are Wi and bi. The

```

```

68 # biases are initialized to zero and the weights are initialized
69 # so that each parameter has mean 0 and standard deviation weight_scale.
70 # ===== #
71
72 dims = []
73 dims = [input_dim] + hidden_dims + [num_classes]
74 for i in np.arange(self.num_layers):
75     self.params['b%d' % (i+1)] = np.zeros(dims[i + 1])
76     self.params['W%d' % (i+1)] = np.random.randn(dims[i], dims[i + 1]) * weight_scale
77     if self.use_batchnorm and i != self.num_layers - 1:
78         self.params['gamma%d' % (i + 1)] = np.ones(dims[i + 1])
79         self.params['beta%d' % (i + 1)] = np.zeros(dims[i + 1])
80
81 # ===== #
82 # END YOUR CODE HERE
83 # ===== #
84
85 # When using dropout we need to pass a dropout_param dictionary to each
86 # dropout layer so that the layer knows the dropout probability and the mode
87 # (train / test). You can pass the same dropout_param to each dropout layer.
88 self.dropout_param = {}
89 if self.use_dropout:
90     self.dropout_param = {'mode': 'train', 'p': dropout}
91     if seed is not None:
92         self.dropout_param['seed'] = seed
93
94 # With batch normalization we need to keep track of running means and
95 # variances, so we need to pass a special bn_param object to each batch
96 # normalization layer. You should pass self.bn_params[0] to the forward pass
97 # of the first batch normalization layer, self.bn_params[1] to the forward
98 # pass of the second batch normalization layer, etc.
99 self.bn_params = []
100 if self.use_batchnorm:
101     self.bn_params = [{'mode': 'train'} for i in np.arange(self.num_layers - 1)]
102
103 # Cast all parameters to the correct datatype
104 for k, v in self.params.items():
105     self.params[k] = v.astype(dtype)
106
107
108 def loss(self, X, y=None):
109     """
110     Compute loss and gradient for the fully-connected net.
111
112     Input / output: Same as TwoLayerNet above.
113     """
114     X = X.astype(self.dtype)
115     mode = 'test' if y is None else 'train'
116
117     # Set train/test mode for batchnorm params and dropout param since they
118     # behave differently during training and testing.
119     if self.dropout_param is not None:
120         self.dropout_param['mode'] = mode
121     if self.use_batchnorm:
122         for bn_param in self.bn_params:
123             bn_param[mode] = mode
124
125     scores = None
126
127     # ===== #
128     # YOUR CODE HERE:
129     # Implement the forward pass of the FC net and store the output
130     # scores as the variable "scores".
131     # ===== #
132
133     layer = {}
134     layer[0] = X

```

```

135 cache_layer = {}
136
137 dropout_caches = {}
138
139 if not self.use_batchnorm:
140     for i in np.arange(1, self.num_layers):
141
142
143         layer[i], cache_layer[i] = affine_relu_forward(layer[i - 1],
144                                                         self.params['W%d' % i],
145                                                         self.params['b%d' % i])
146
147         if self.use_dropout:
148             layer[i], dropout_caches[i] = dropout_forward(layer[i], self.
149                                                         dropout_param)
150
151         Weight_out = 'W%d' % self.num_layers
152         bais_out = 'b%d' % self.num_layers
153         scores, cache_scores = affine_forward(layer[self.num_layers - 1], self.params[
154             Weight_out], self.params[bais_out])
155
156     else:
157         for i in np.arange(1, self.num_layers):
158
159
160             layer[i], cache_layer[i] = affine_batchnorm_relu(layer[i - 1], self.params[
161                 "W%d" % i], self.params['b%d'%i], self.params['gamma%d' % i], self.params[
162                 'beta%d' % i], bn_param=self.bn_params[i - 1])
163
164             if self.use_dropout:
165                 layer[i], dropout_caches[i] = dropout_forward(layer[i], self.
166                                                         dropout_param)
167
168             Weight_out = 'W%d' % self.num_layers
169             bais_out = 'b%d' % self.num_layers
170             scores, cache_scores = affine_forward(layer[self.num_layers - 1], self.params[
171                 Weight_out], self.params[bais_out])
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193

```

```

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

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

loss, grads = 0.0, {}
# ===== #
# YOUR CODE HERE:
#   Implement the backwards pass of the FC net and store the gradients
#   in the grads dict, so that grads[k] is the gradient of self.params[k]
#   Be sure your L2 regularization includes a 0.5 factor.
# ===== #

loss, scores_delta = softmax_loss(scores, y)

for i in np.arange(1, self.num_layers + 1):
    loss += 0.5 * self.reg * np.sum(self.params['W%d' % i]**2) #add regularization
    for each layer

dx = {}

dx[self.num_layers], grads[Weight_out], grads[bais_out] = affine_backward(
    scores_delta, cache_scores)

grads[Weight_out] += self.reg * self.params[Weight_out] #add regularization to output

for i in reversed(np.arange(1, self.num_layers)):
    if not self.use_batchnorm:
        if self.use_dropout:
            dx[i] = dropout_backward(dx[i + 1], dropout_caches[i])

```



```

194         dx[i], grads['W%d' % i], grads['b%d' % i] = affine_relu_backward(dx[i],
195                                     cache_layer[i])
196     else:
197         dx[i], grads['W%d' % i], grads['b%d' % i] = affine_relu_backward(dx[i +
198                                     1], cache_layer[i])
199
200     else:
201         if self.use_dropout:
202             dx[i] = dropout_backward(dx[i + 1], dropout_caches[i])
203             dx[i], grads["W%d" % i], grads['b%d'%i], grads['gamma%d' % (i)], grads[
204                 'beta%d' % (i)] = affine_batchnorm_relu_backwards(dx[i + 1], cache_layer
205                             [i])
206
207             else:
208                 dx[i], grads["W%d" % i], grads['b%d'%i], grads['gamma%d' % (i)], grads[
209                     'beta%d' % (i)] = affine_batchnorm_relu_backwards(dx[i + 1], cache_layer
210                             [i])
211
212         grads['W%d' % i] += self.reg * self.params['W%d' % i]
213
214     # ===== #
215     # END YOUR CODE HERE
216     # ===== #
217
218     return loss, grads

```