

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

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

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

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

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

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

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

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



```

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

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

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

Building upon your HW #3 implementation

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

- `affine_forward` in `nndl/layers.py`
- `affine_backward` in `nndl/layers.py`
- `relu_forward` in `nndl/layers.py`
- `relu_backward` in `nndl/layers.py`
- `affine_relu_forward` in `nndl/layer_utils.py`
- `affine_relu_backward` in `nndl/layer_utils.py`

- The FullyConnectedNet class in nndl/fc_net.py

Test all functions you copy and pasted

```
In [5]: 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: 9.769847728806635e-10

If affine_backward is working, error should be less than 1e-9::
dx error: 2.363967198662533e-10
dw error: 1.6097343635075995e-10
db error: 1.580234993153846e-10

If relu_forward function is working, difference should be around 1e-8:
difference: 4.99999798022158e-08

If relu_forward function is working, error should be less than 1e-9:
dx error: 3.275625385633504e-12

If affine_relu_forward and affine_relu_backward are working, error should be less than 1e-9::
dx error: 3.2582542651794285e-10
dw error: 6.005113770498464e-10
db error: 2.548002399814824e-11

```
Running check with reg = 0
Initial loss: 2.303374998448424
W1 relative error: 6.857310860468233e-07
W2 relative error: 7.676133530383611e-07
W3 relative error: 3.21777696148275e-07
b1 relative error: 2.516268473877123e-09
b2 relative error: 5.8869748840310205e-09
b3 relative error: 1.3958093070769016e-10
Running check with reg = 3.14
Initial loss: 6.805991217574755
W1 relative error: 1.1150231465873393e-08
W2 relative error: 1.079142654594727e-07
W3 relative error: 2.63566134196156e-08
b1 relative error: 2.8315564777658047e-08
b2 relative error: 2.575581561894071e-09
b3 relative error: 2.4802815598088873e-10
```

Training a larger model

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

SGD + momentum

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

```
In [6]: from nnndl.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 [7]: from nnndl.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 [8]: num_train = 4000
```

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

solvers = {}

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

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

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

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

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

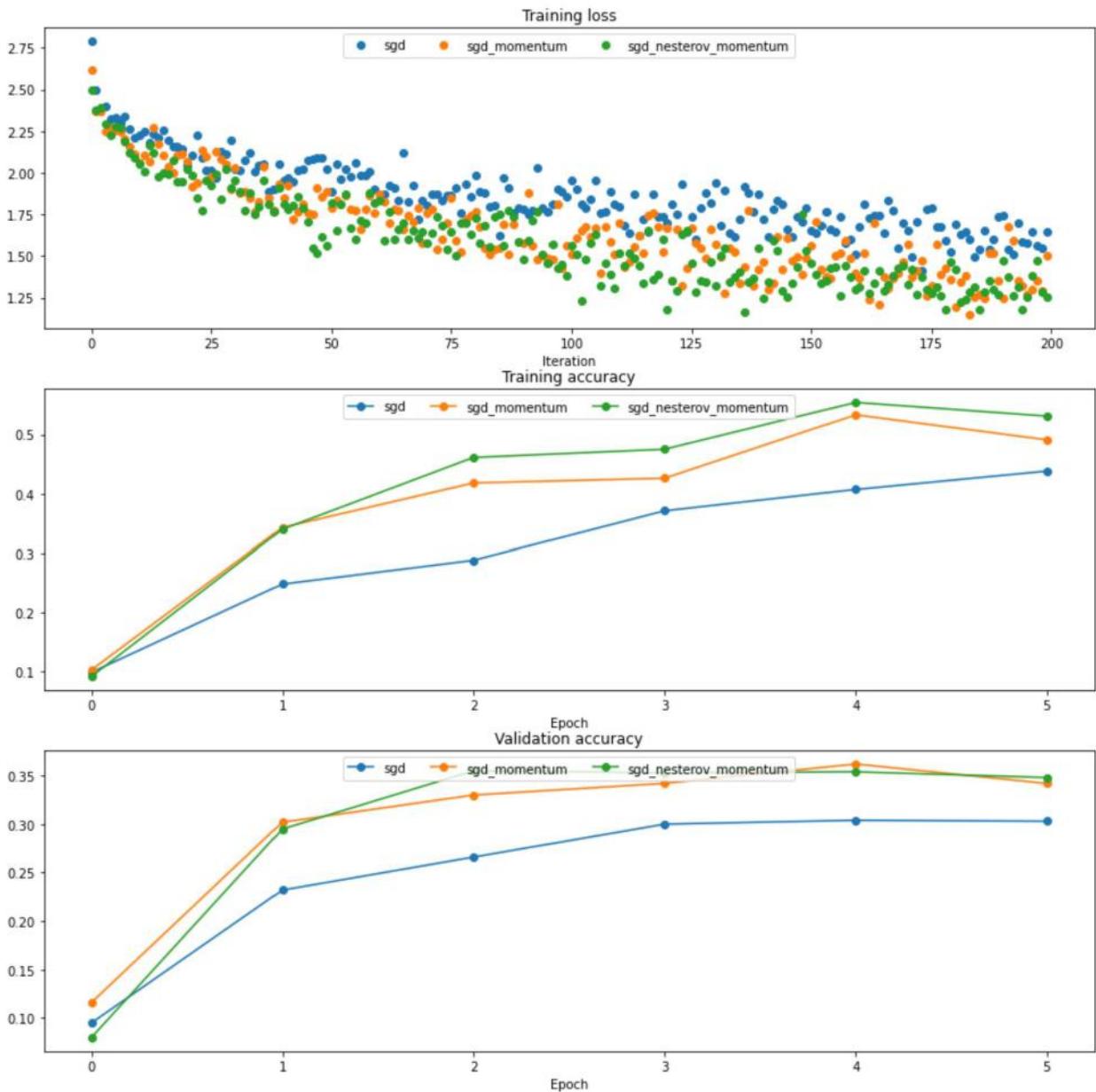
for update_rule, solver in solvers.items():
    plt.subplot(3, 1, 1)
    plt.plot(solver.loss_history, 'o', label=update_rule)

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

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

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

Optimizing with sgd
Optimizing with sgd_momentum
Optimizing with sgd_nesterov_momentum



RMSProp

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

In [9]:

```
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 [11]:

```
# 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 [12]:

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

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

    solver = Solver(model, small_data,
                    num_epochs=5, batch_size=100,
                    update_rule=update_rule,
                    optim_config={
                        'learning_rate': learning_rates[update_rule]
                    },
```

```
        verbose=False)
solvers[update_rule] = solver
solver.train()
print

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

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

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

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

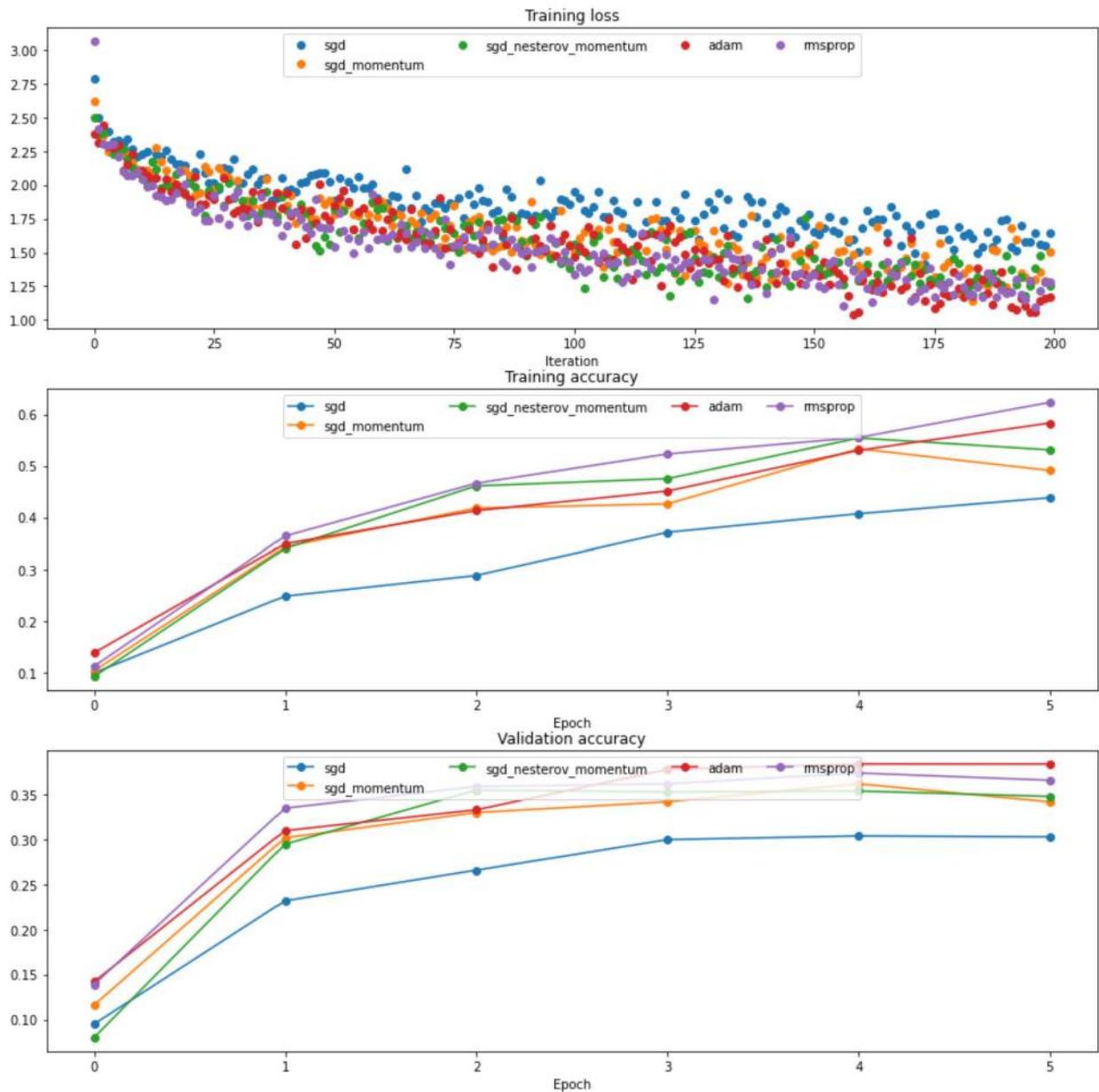
    plt.subplot(3, 1, 2)
    plt.plot(solver.train_acc_history, '-o', label=update_rule)

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

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

Optimizing with adam

Optimizing with rmsprop



Easier optimization

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

In [13]:

```

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: 2.303629
(Epoch 0 / 10) train acc: 0.150000; val_acc: 0.164000
(Iteration 51 / 4900) loss: 1.774466
(Iteration 101 / 4900) loss: 1.864332
(Iteration 151 / 4900) loss: 1.556623
(Iteration 201 / 4900) loss: 1.557418
(Iteration 251 / 4900) loss: 1.742280
(Iteration 301 / 4900) loss: 1.905783
(Iteration 351 / 4900) loss: 1.637200
(Iteration 401 / 4900) loss: 1.526135
(Iteration 451 / 4900) loss: 1.521564
(Epoch 1 / 10) train acc: 0.428000; val_acc: 0.441000
(Iteration 501 / 4900) loss: 1.561486
(Iteration 551 / 4900) loss: 1.554491
(Iteration 601 / 4900) loss: 1.543946
(Iteration 651 / 4900) loss: 1.385529
(Iteration 701 / 4900) loss: 1.425173
(Iteration 751 / 4900) loss: 1.540616
(Iteration 801 / 4900) loss: 1.357907
(Iteration 851 / 4900) loss: 1.626862
(Iteration 901 / 4900) loss: 1.466829
(Iteration 951 / 4900) loss: 1.544114
(Epoch 2 / 10) train acc: 0.488000; val_acc: 0.484000
(Iteration 1001 / 4900) loss: 1.666742
(Iteration 1051 / 4900) loss: 1.526628
(Iteration 1101 / 4900) loss: 1.539166
(Iteration 1151 / 4900) loss: 1.446044
(Iteration 1201 / 4900) loss: 1.414799
(Iteration 1251 / 4900) loss: 1.342977
(Iteration 1301 / 4900) loss: 1.575358
(Iteration 1351 / 4900) loss: 1.419483
(Iteration 1401 / 4900) loss: 1.440173
(Iteration 1451 / 4900) loss: 1.176719
(Epoch 3 / 10) train acc: 0.516000; val_acc: 0.502000
(Iteration 1501 / 4900) loss: 1.335112
(Iteration 1551 / 4900) loss: 1.293916
(Iteration 1601 / 4900) loss: 1.470218
(Iteration 1651 / 4900) loss: 1.232994
(Iteration 1701 / 4900) loss: 1.411225
(Iteration 1751 / 4900) loss: 1.422639
(Iteration 1801 / 4900) loss: 1.352091
(Iteration 1851 / 4900) loss: 1.324811
(Iteration 1901 / 4900) loss: 1.329972
(Iteration 1951 / 4900) loss: 1.233757
(Epoch 4 / 10) train acc: 0.545000; val_acc: 0.476000
(Iteration 2001 / 4900) loss: 1.205185
(Iteration 2051 / 4900) loss: 1.470377
(Iteration 2101 / 4900) loss: 1.162367
(Iteration 2151 / 4900) loss: 1.357951
(Iteration 2201 / 4900) loss: 1.232868
(Iteration 2251 / 4900) loss: 1.272105
(Iteration 2301 / 4900) loss: 1.074457
(Iteration 2351 / 4900) loss: 1.220854
(Iteration 2401 / 4900) loss: 0.974765
(Epoch 5 / 10) train acc: 0.562000; val_acc: 0.529000
(Iteration 2451 / 4900) loss: 1.174441
(Iteration 2501 / 4900) loss: 1.105419
(Iteration 2551 / 4900) loss: 1.082825
(Iteration 2601 / 4900) loss: 1.200558
(Iteration 2651 / 4900) loss: 1.071913
(Iteration 2701 / 4900) loss: 1.170630
(Iteration 2751 / 4900) loss: 1.241453
(Iteration 2801 / 4900) loss: 1.030365
(Iteration 2851 / 4900) loss: 0.973189
(Iteration 2901 / 4900) loss: 1.332866
(Epoch 6 / 10) train acc: 0.599000; val_acc: 0.517000
(Iteration 2951 / 4900) loss: 1.107921
(Iteration 3001 / 4900) loss: 1.103678
```

```
(Iteration 3051 / 4900) loss: 1.173848
(Iteration 3101 / 4900) loss: 1.237449
(Iteration 3151 / 4900) loss: 1.363807
(Iteration 3201 / 4900) loss: 1.139764
(Iteration 3251 / 4900) loss: 1.208968
(Iteration 3301 / 4900) loss: 0.998561
(Iteration 3351 / 4900) loss: 0.847466
(Iteration 3401 / 4900) loss: 1.143264
(Epoch 7 / 10) train acc: 0.624000; val_acc: 0.529000
(Iteration 3451 / 4900) loss: 1.009258
(Iteration 3501 / 4900) loss: 0.850020
(Iteration 3551 / 4900) loss: 0.853643
(Iteration 3601 / 4900) loss: 0.847587
(Iteration 3651 / 4900) loss: 1.089864
(Iteration 3701 / 4900) loss: 0.983375
(Iteration 3751 / 4900) loss: 1.192917
(Iteration 3801 / 4900) loss: 1.186817
(Iteration 3851 / 4900) loss: 0.964879
(Iteration 3901 / 4900) loss: 1.132026
(Epoch 8 / 10) train acc: 0.670000; val_acc: 0.523000
(Iteration 3951 / 4900) loss: 1.224819
(Iteration 4001 / 4900) loss: 0.842471
(Iteration 4051 / 4900) loss: 0.991945
(Iteration 4101 / 4900) loss: 1.026309
(Iteration 4151 / 4900) loss: 0.979127
(Iteration 4201 / 4900) loss: 1.057147
(Iteration 4251 / 4900) loss: 0.888057
(Iteration 4301 / 4900) loss: 1.153214
(Iteration 4351 / 4900) loss: 0.899978
(Iteration 4401 / 4900) loss: 1.044466
(Epoch 9 / 10) train acc: 0.689000; val_acc: 0.538000
(Iteration 4451 / 4900) loss: 1.143102
(Iteration 4501 / 4900) loss: 1.131958
(Iteration 4551 / 4900) loss: 0.922324
(Iteration 4601 / 4900) loss: 1.114319
(Iteration 4651 / 4900) loss: 0.966344
(Iteration 4701 / 4900) loss: 0.886925
(Iteration 4751 / 4900) loss: 1.032075
(Iteration 4801 / 4900) loss: 1.079755
(Iteration 4851 / 4900) loss: 0.713890
(Epoch 10 / 10) train acc: 0.684000; val_acc: 0.548000
```

In [14]:

```
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: 0.548

Test set accuracy: 0.526

In []:

```
#py files are appended in the end
```

ECE C147/247 HW4 Q2: Batch Normalization

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

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

In [1]:

```
## Import and setups

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

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

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

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



```

In [2]:

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

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

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

Batchnorm forward pass

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

In [10]:

```
# 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: [-31.54775839 -2.13022673 7.25127046]
stds: [34.27351436 28.88761231 26.67249724]
After batch normalization (gamma=1, beta=0)
mean: [1.07691633e-16 2.73392420e-17 -2.25375274e-16]
std: [1. 0.99999999 0.99999999]
After batch normalization (nontrivial gamma, beta)
means: [11. 12. 13.]
stds: [1. 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 [12]:

```

# 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: [0.07545855 -0.0280366 0.00382865]
stds: [1.07523389 0.88768455 1.00772173]

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 [29]: # 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: 3.5297368619555174e-10
dgamma error: 7.31144936081657e-12
dbeta error: 3.55494649413544e-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 [39]:

```
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: 2.319241634671405
W1 relative error: 1.6918241008909617e-05
W2 relative error: 4.792639563597038e-06
W3 relative error: 1.675576175771179e-09
b1 relative error: 2.220446049250313e-08
b2 relative error: 4.163336342344337e-09
b3 relative error: 1.6114415179444514e-10
beta1 relative error: 5.5828517433487906e-09
beta2 relative error: 5.346903076489941e-09
gamma1 relative error: 5.877709094343308e-09
gamma2 relative error: 4.9197460741167614e-09

```

```

Running check with reg = 3.14
Initial loss: 7.399125432521264
W1 relative error: 0.00014917960368898732
W2 relative error: 1.6363784203472635e-05
W3 relative error: 2.3982468429133698e-08
b1 relative error: 5.551115123125783e-09
b2 relative error: 4.6629367034256575e-07
b3 relative error: 2.055117372645774e-10
beta1 relative error: 3.3494040953666204e-08
beta2 relative error: 4.185378026215697e-09
gamma1 relative error: 3.2529576716641514e-08
gamma2 relative error: 7.05719777990773e-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 [42]:

```

# Try training a very deep net with batchnorm
hidden_dims = [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: 2.297757
(Epoch 0 / 10) train acc: 0.101000; val_acc: 0.109000
(Epoch 1 / 10) train acc: 0.335000; val_acc: 0.267000
(Epoch 2 / 10) train acc: 0.415000; val_acc: 0.307000

```

```
(Epoch 3 / 10) train acc: 0.444000; val_acc: 0.295000
(Epoch 4 / 10) train acc: 0.483000; val_acc: 0.317000
(Epoch 5 / 10) train acc: 0.591000; val_acc: 0.300000
(Epoch 6 / 10) train acc: 0.603000; val_acc: 0.308000
(Epoch 7 / 10) train acc: 0.686000; val_acc: 0.317000
(Epoch 8 / 10) train acc: 0.716000; val_acc: 0.304000
(Epoch 9 / 10) train acc: 0.809000; val_acc: 0.338000
(Epoch 10 / 10) train acc: 0.795000; val_acc: 0.339000
(Iteration 1 / 200) loss: 2.302021
(Epoch 0 / 10) train acc: 0.144000; val_acc: 0.124000
(Epoch 1 / 10) train acc: 0.268000; val_acc: 0.235000
(Epoch 2 / 10) train acc: 0.299000; val_acc: 0.277000
(Epoch 3 / 10) train acc: 0.331000; val_acc: 0.266000
(Epoch 4 / 10) train acc: 0.329000; val_acc: 0.260000
(Epoch 5 / 10) train acc: 0.396000; val_acc: 0.277000
(Epoch 6 / 10) train acc: 0.449000; val_acc: 0.314000
(Epoch 7 / 10) train acc: 0.482000; val_acc: 0.310000
(Epoch 8 / 10) train acc: 0.512000; val_acc: 0.287000
(Epoch 9 / 10) train acc: 0.558000; val_acc: 0.305000
(Epoch 10 / 10) train acc: 0.540000; val_acc: 0.263000
```

In [43]:

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

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

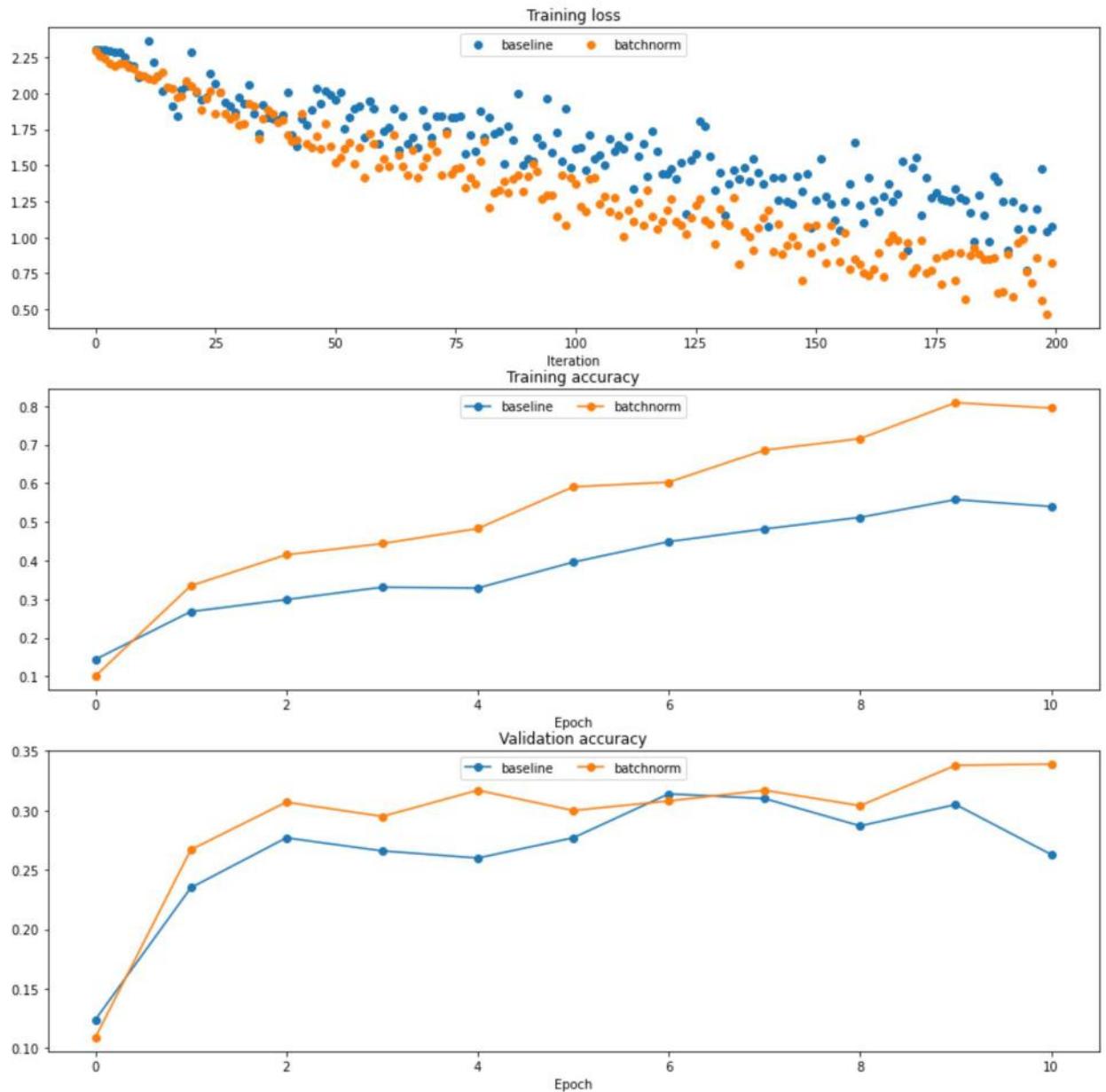
plt.subplot(3, 1, 3)
plt.title('Validation accuracy')
plt.xlabel('Epoch')

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

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

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

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



Batchnorm and initialization

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

In [45]:

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

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

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

    bn_solver = Solver(bn_model, small_data,
```

```

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

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

```

```

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

```

In [46]:

```

# 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()

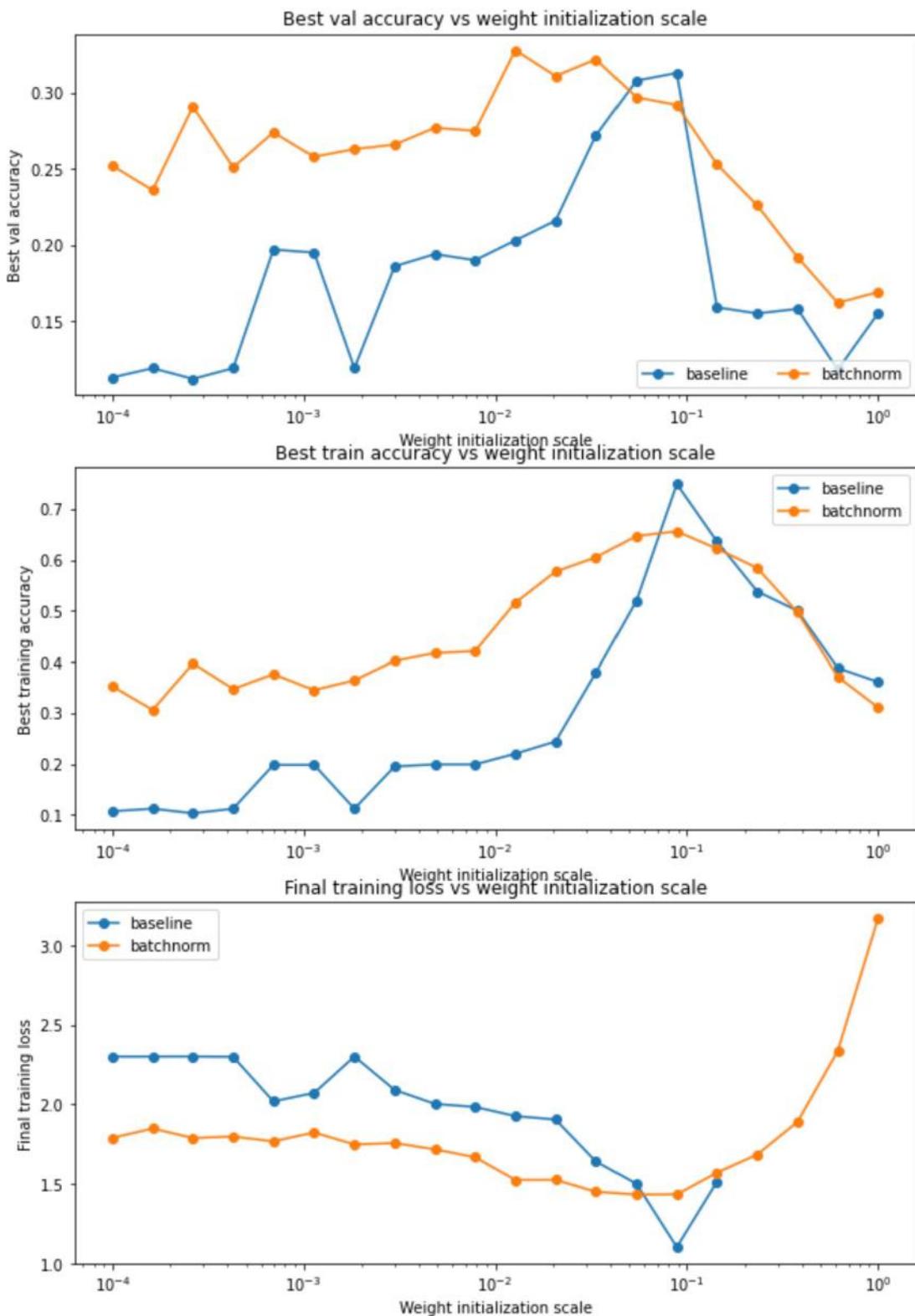
```

```

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:

1. Batchnorm makes the model more robust against weight initialization. We see that training and validation accuracy fluctuate more for baseline model than that for batchnorm model.
2. We also see that there is an optimum weight scale for which both baseline and batchnorm model performs good than other weight scale.

```
In [ ]: #ALL py files are appended in the end
```

ECE C147/247 HW4 Q3: Dropout

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

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

In [2]:

```
## Import and setups

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

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

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

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



```

In [3]:

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

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

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

Dropout forward pass

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

In [4]:

```
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: 10.000668613713762
Mean of train-time output: 9.985453927668857
Mean of test-time output: 10.000668613713762
Fraction of train-time output set to zero: 0.301104
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.6
Mean of input: 10.000668613713762
Mean of train-time output: 10.00106528141077
Mean of test-time output: 10.000668613713762
Fraction of train-time output set to zero: 0.600112
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.75
Mean of input: 10.000668613713762
Mean of train-time output: 10.00447872476764
Mean of test-time output: 10.000668613713762
Fraction of train-time output set to zero: 0.749832
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]:
x = np.random.randn(10, 10) + 10
dout = np.random.randn(*x.shape)

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

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

dx relative error: 1.892902803621409e-11
```

Implement a fully connected neural network with dropout layers

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

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

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

```
In [6]:
N, D, H1, H2, C = 2, 15, 20, 30, 10
X = np.random.randn(N, D)
y = np.random.randint(C, size=(N,))

for 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: 2.5272575351926774e-07
W2 relative error: 1.5034484923408755e-05
W3 relative error: 2.753446897744538e-07
b1 relative error: 2.93695746140462e-06
b2 relative error: 5.0513392951881574e-08
b3 relative error: 1.1740467838205477e-10

Running check with dropout = 0.25
Initial loss: 2.29898614757146
W1 relative error: 9.737728978465634e-07
W2 relative error: 2.4340451942780774e-08
W3 relative error: 3.042456581180938e-08
b1 relative error: 2.005618875946033e-08
b2 relative error: 1.897778283870511e-09
b3 relative error: 1.302003889798156e-10

Running check with dropout = 0.5
Initial loss: 2.302437587710995
W1 relative error: 4.553388018985473e-08
W2 relative error: 2.9742181064823555e-08
W3 relative error: 4.3413246357224405e-07
b1 relative error: 1.872463027156301e-08
b2 relative error: 5.045590828136527e-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.

```

In [7]: # Train two identical nets, one with dropout and one without

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

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

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

```

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

In [8]:

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

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

plt.subplot(3, 1, 1)
for dropout in dropout_choices:
    plt.plot(solvers[dropout].train_acc_history, 'o', label='%.2f dropout' % dropout)
```

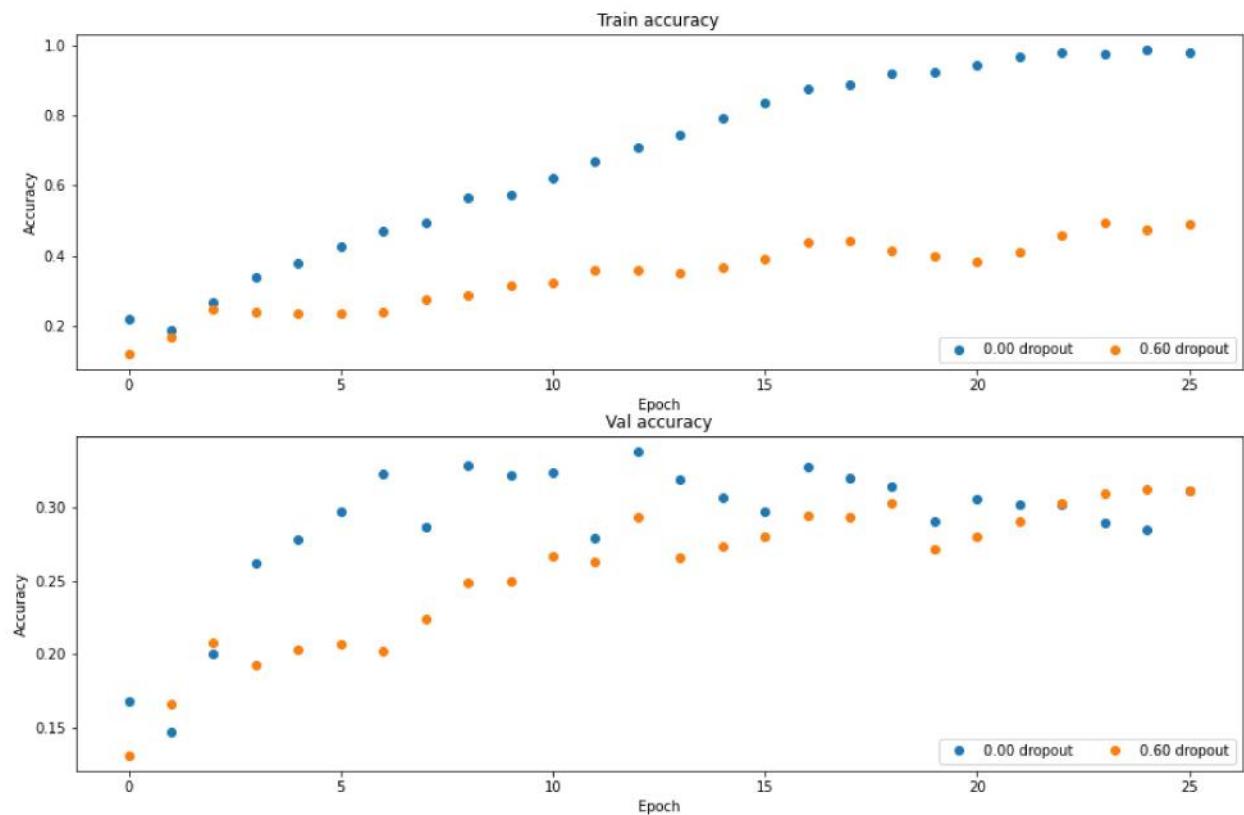
```

plt.title('Train accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend(ncol=2, loc='lower right')

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

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

```



Question

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

Answer:

Yes, dropout is performing regularization as we can see the both the models have similar validation accuracy but model with no dropout have very high training accuracy i.e. model without dropout is overfitting.

Final part of the assignment

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

```
min(floor((X - 32%) / 28%, 1)
```

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

In [15]:

```
# ===== #
# YOUR CODE HERE:
```

```
# Implement a FC-net that achieves at Least 55% validation accuracy
# on CIFAR-10.

# ===== #
hidden_dims = [200, 400, 800, 800, 400]
model = FullyConnectedNet(hidden_dims, dropout=0.2, use_batchnorm=True, weight_scale=1e-2)

solver = Solver(model, data,
                num_epochs=25, batch_size=256,
                update_rule='adam',
                optim_config={
                    'learning_rate': 5e-3,
                },
                lr_decay=0.95,
                verbose=True, print_every=100)
solver.train()

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

```
(Iteration 1 / 4775) loss: 2.302327
(Epoch 0 / 25) train acc: 0.113000; val_acc: 0.097000
(Iteration 101 / 4775) loss: 1.779209
(Epoch 1 / 25) train acc: 0.436000; val_acc: 0.468000
(Iteration 201 / 4775) loss: 1.513056
(Iteration 301 / 4775) loss: 1.578397
(Epoch 2 / 25) train acc: 0.525000; val_acc: 0.457000
(Iteration 401 / 4775) loss: 1.602773
(Iteration 501 / 4775) loss: 1.488037
(Epoch 3 / 25) train acc: 0.521000; val_acc: 0.509000
(Iteration 601 / 4775) loss: 1.590357
(Iteration 701 / 4775) loss: 1.290636
(Epoch 4 / 25) train acc: 0.552000; val_acc: 0.540000
(Iteration 801 / 4775) loss: 1.365004
(Iteration 901 / 4775) loss: 1.337844
(Epoch 5 / 25) train acc: 0.577000; val_acc: 0.543000
(Iteration 1001 / 4775) loss: 1.164903
(Iteration 1101 / 4775) loss: 1.314970
(Epoch 6 / 25) train acc: 0.620000; val_acc: 0.564000
(Iteration 1201 / 4775) loss: 1.256558
(Iteration 1301 / 4775) loss: 1.158617
(Epoch 7 / 25) train acc: 0.587000; val_acc: 0.548000
(Iteration 1401 / 4775) loss: 1.183747
(Iteration 1501 / 4775) loss: 1.192084
(Epoch 8 / 25) train acc: 0.621000; val_acc: 0.564000
(Iteration 1601 / 4775) loss: 1.161690
(Iteration 1701 / 4775) loss: 1.077603
(Epoch 9 / 25) train acc: 0.653000; val_acc: 0.584000
(Iteration 1801 / 4775) loss: 1.049501
(Iteration 1901 / 4775) loss: 1.013423
(Epoch 10 / 25) train acc: 0.673000; val_acc: 0.559000
(Iteration 2001 / 4775) loss: 1.204154
(Iteration 2101 / 4775) loss: 1.177057
(Epoch 11 / 25) train acc: 0.659000; val_acc: 0.565000
(Iteration 2201 / 4775) loss: 1.151376
(Epoch 12 / 25) train acc: 0.703000; val_acc: 0.582000
(Iteration 2301 / 4775) loss: 1.034299
(Iteration 2401 / 4775) loss: 1.110599
(Epoch 13 / 25) train acc: 0.704000; val_acc: 0.588000
(Iteration 2501 / 4775) loss: 1.062121
(Iteration 2601 / 4775) loss: 0.994460
(Epoch 14 / 25) train acc: 0.725000; val_acc: 0.590000
(Iteration 2701 / 4775) loss: 1.092445
(Iteration 2801 / 4775) loss: 1.001792
(Epoch 15 / 25) train acc: 0.704000; val_acc: 0.584000
(Iteration 2901 / 4775) loss: 0.977682
(Iteration 3001 / 4775) loss: 0.877472
(Epoch 16 / 25) train acc: 0.727000; val_acc: 0.577000
(Iteration 3101 / 4775) loss: 0.887661
(Iteration 3201 / 4775) loss: 0.853887
(Epoch 17 / 25) train acc: 0.748000; val_acc: 0.583000
```

```
(Iteration 3301 / 4775) loss: 0.907220
(Iteration 3401 / 4775) loss: 0.847191
(Epoch 18 / 25) train acc: 0.752000; val_acc: 0.585000
(Iteration 3501 / 4775) loss: 0.808808
(Iteration 3601 / 4775) loss: 0.825502
(Epoch 19 / 25) train acc: 0.777000; val_acc: 0.593000
(Iteration 3701 / 4775) loss: 0.812555
(Iteration 3801 / 4775) loss: 0.804781
(Epoch 20 / 25) train acc: 0.786000; val_acc: 0.581000
(Iteration 3901 / 4775) loss: 0.853470
(Iteration 4001 / 4775) loss: 0.884780
(Epoch 21 / 25) train acc: 0.786000; val_acc: 0.584000
(Iteration 4101 / 4775) loss: 0.791657
(Iteration 4201 / 4775) loss: 0.919025
(Epoch 22 / 25) train acc: 0.802000; val_acc: 0.574000
(Iteration 4301 / 4775) loss: 0.790375
(Epoch 23 / 25) train acc: 0.807000; val_acc: 0.586000
(Iteration 4401 / 4775) loss: 0.818866
(Iteration 4501 / 4775) loss: 0.964608
(Epoch 24 / 25) train acc: 0.809000; val_acc: 0.580000
(Iteration 4601 / 4775) loss: 0.869007
(Iteration 4701 / 4775) loss: 0.870042
(Epoch 25 / 25) train acc: 0.818000; val_acc: 0.600000
```

In []: *#all py files are attached in the end*

```

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

```

```

72 # END YOUR CODE HERE
73 # ===== #
74
75 config['velocity'] = v
76
77 return next_w, config
78
79 def sgd_nesterov_momentum(w, dw, config=None):
80 """
81     Performs stochastic gradient descent with Nesterov momentum.
82
83     config format:
84     - learning_rate: Scalar learning rate.
85     - momentum: Scalar between 0 and 1 giving the momentum value.
86         Setting momentum = 0 reduces to sgd.
87     - velocity: A numpy array of the same shape as w and dw used to store a moving
88         average of the gradients.
89 """
90 if config is None: config = {}
91 config.setdefault('learning_rate', 1e-2)
92 config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
93 v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.
94
95 # ===== #
96 # YOUR CODE HERE:
97 #     Implement the momentum update formula. Return the updated weights
98 #     as next_w, and the updated velocity as v.
99 # ===== #
100 v_new = config['momentum'] * v - config['learning_rate'] * dw
101 next_w = w+v_new + config['momentum'] *(v_new-v)
102 v = v_new
103 # ===== #
104 # END YOUR CODE HERE
105 # ===== #
106
107 config['velocity'] = v
108
109 return next_w, config
110
111 def rmsprop(w, dw, config=None):
112 """
113     Uses the RMSProp update rule, which uses a moving average of squared gradient
114     values to set adaptive per-parameter learning rates.
115
116     config format:
117     - learning_rate: Scalar learning rate.
118     - decay_rate: Scalar between 0 and 1 giving the decay rate for the squared
119         gradient cache.
120     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
121     - beta: Moving average of second moments of gradients.
122 """
123 if config is None: config = {}
124 config.setdefault('learning_rate', 1e-2)
125 config.setdefault('decay_rate', 0.99)
126 config.setdefault('epsilon', 1e-8)
127 config.setdefault('a', np.zeros_like(w))
128
129 next_w = None
130
131 # ===== #
132 # YOUR CODE HERE:
133 #     Implement RMSProp. Store the next value of w as next_w. You need
134 #     to also store in config['a'] the moving average of the second
135 #     moment gradients, so they can be used for future gradients. Concretely,
136 #     config['a'] corresponds to "a" in the lecture notes.
137 # ===== #
138 config['a'] = config['decay_rate']*config['a'] + (1-config['decay_rate'])*dw*dw
139 next_w = w - (config['learning_rate']*dw)/(np.sqrt(config['a'])+config['epsilon'])
140 # ===== #
141 # END YOUR CODE HERE
142 # ===== #

```

```

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

```

```

1 import numpy as np
2
3
4 def affine_forward(x, w, b):
5     """
6         Computes the forward pass for an affine (fully-connected) layer.
7
8         The input x has shape (N, d_1, ..., d_k) and contains a minibatch of N
9         examples, where each example x[i] has shape (d_1, ..., d_k). We will
10        reshape each input into a vector of dimension D = d_1 * ... * d_k, and
11        then transform it to an output vector of dimension M.
12
13    Inputs:
14    - x: A numpy array containing input data, of shape (N, d_1, ..., d_k)
15    - w: A numpy array of weights, of shape (D, M)
16    - b: A numpy array of biases, of shape (M,)
17
18    Returns a tuple of:
19    - out: output, of shape (N, M)
20    - cache: (x, w, b)
21    """
22
23    # ===== #
24    # YOUR CODE HERE:
25    #   Calculate the output of the forward pass. Notice the dimensions
26    #   of w are D x M, which is the transpose of what we did in earlier
27    #   assignments.
28    # ===== #
29
30    out = x.reshape(x.shape[0],-1).dot(w)+b
31
32    # ===== #
33    # END YOUR CODE HERE
34    # ===== #
35
36    cache = (x, w, b)
37    return out, cache
38
39
40 def affine_backward(dout, cache):
41     """
42         Computes the backward pass for an affine layer.
43
44     Inputs:
45     - dout: Upstream derivative, of shape (N, M)
46     - cache: Tuple of:
47         - x: A numpy array containing input data, of shape (N, d_1, ..., d_k)
48         - w: A numpy array of weights, of shape (D, M)
49         - b: A numpy array of biases, of shape (M,)
50
51     Returns a tuple of:
52     - dx: Gradient with respect to x, of shape (N, d_1, ..., d_k)
53     - dw: Gradient with respect to w, of shape (D, M)
54     - db: Gradient with respect to b, of shape (M,)
55     """
56     x, w, b = cache
57     dx, dw, db = None, None, None
58
59    # ===== #
60    # YOUR CODE HERE:
61    #   Calculate the gradients for the backward pass.
62    # Notice:
63    #   dout is N x M
64    #   dx should be N x d1 x ... x dk; it relates to dout through multiplication with
65    #   w, which is D x M
66    #   dw should be D x M; it relates to dout through multiplication with x, which is
67    #   N x D after reshaping
68    #   db should be M; it is just the sum over dout examples
69    # ===== #
70    dx = dout.dot(w.T).reshape(x.shape)
71    dw = x.reshape(x.shape[0],-1).T.dot(dout).reshape(w.shape)
72    db = np.sum(dout, axis=0)

```

```

71
72     # ===== #
73     # END YOUR CODE HERE
74     # ===== #
75
76     return dx, dw, db
77
78 def relu_forward(x):
79     """
80         Computes the forward pass for a layer of rectified linear units (ReLUs).
81
82         Input:
83         - x: Inputs, of any shape
84
85         Returns a tuple of:
86         - out: Output, of the same shape as x
87         - cache: x
88     """
89     # ===== #
90     # YOUR CODE HERE:
91     #     Implement the ReLU forward pass.
92     # ===== #
93     out = np.maximum(x,0)
94     # ===== #
95     # END YOUR CODE HERE
96     # ===== #
97
98     cache = x
99     return out, cache
100
101
102 def relu_backward(dout, cache):
103     """
104         Computes the backward pass for a layer of rectified linear units (ReLUs).
105
106         Input:
107         - dout: Upstream derivatives, of any shape
108         - cache: Input x, of same shape as dout
109
110        Returns:
111        - dx: Gradient with respect to x
112    """
113     x = cache
114
115     # ===== #
116     # YOUR CODE HERE:
117     #     Implement the ReLU backward pass
118     # ===== #
119     dx = dout
120     dx[x<0]=0
121
122     # ===== #
123     # END YOUR CODE HERE
124     # ===== #
125
126     return dx
127
128 def batchnorm_forward(x, gamma, beta, bn_param):
129     """
130         Forward pass for batch normalization.
131
132         During training the sample mean and (uncorrected) sample variance are
133         computed from minibatch statistics and used to normalize the incoming data.
134         During training we also keep an exponentially decaying running mean of the mean
135         and variance of each feature, and these averages are used to normalize data
136         at test-time.
137
138         At each timestep we update the running averages for mean and variance using
139         an exponential decay based on the momentum parameter:
140
141         running_mean = momentum * running_mean + (1 - momentum) * sample_mean
142         running_var = momentum * running_var + (1 - momentum) * sample_var

```

```

143
144 Note that the batch normalization paper suggests a different test-time
145 behavior: they compute sample mean and variance for each feature using a
146 large number of training images rather than using a running average. For
147 this implementation we have chosen to use running averages instead since
148 they do not require an additional estimation step; the torch7 implementation
149 of batch normalization also uses running averages.
150
151 Input:
152 - x: Data of shape (N, D)
153 - gamma: Scale parameter of shape (D,)
154 - beta: Shift paremeter of shape (D,)
155 - bn_param: Dictionary with the following keys:
156     - mode: 'train' or 'test'; required
157     - eps: Constant for numeric stability
158     - momentum: Constant for running mean / variance.
159     - running_mean: Array of shape (D,) giving running mean of features
160     - running_var Array of shape (D,) giving running variance of features
161
162 Returns a tuple of:
163 - out: of shape (N, D)
164 - cache: A tuple of values needed in the backward pass
165 """
166 mode = bn_param['mode']
167 eps = bn_param.get('eps', 1e-5)
168 momentum = bn_param.get('momentum', 0.9)
169
170 N, D = x.shape
171 running_mean = bn_param.get('running_mean', np.zeros(D, dtype=x.dtype))
172 running_var = bn_param.get('running_var', np.zeros(D, dtype=x.dtype))
173
174 out, cache = None, None
175 if mode == 'train':
176
177     # ===== #
178     # YOUR CODE HERE:
179     # A few steps here:
180     #     (1) Calculate the running mean and variance of the minibatch.
181     #     (2) Normalize the activations with the running mean and variance.
182     #     (3) Scale and shift the normalized activations. Store this
183     #         as the variable 'out'
184     #     (4) Store any variables you may need for the backward pass in
185     #         the 'cache' variable.
186     # ===== #
187     cur_mean = np.mean(x, axis=0)
188     cur_var = np.var(x, axis=0)
189     running_mean = momentum * running_mean + (1 - momentum) * cur_mean
190     running_var = momentum * running_var + (1 - momentum) * cur_var
191     x_n = (x - cur_mean) / np.sqrt(cur_var + eps)
192     out = gamma * x_n + beta
193     cache = x, x_n, cur_mean, cur_var, eps, gamma
194     # ===== #
195     # END YOUR CODE HERE
196     # ===== #
197
198 elif mode == 'test':
199
200     # ===== #
201     # YOUR CODE HERE:
202     #     Calculate the testing time normalized activation. Normalize using
203     #     the running mean and variance, and then scale and shift appropriately.
204     #     Store the output as 'out'.
205     # ===== #
206
207     x = (x - running_mean) / np.sqrt(running_var + eps)
208     out = gamma * x + beta
209
210     # ===== #
211     # END YOUR CODE HERE
212     # ===== #
213
214 else:

```

```

215     raise ValueError('Invalid forward batchnorm mode "%s"' % mode)
216
217     # Store the updated running means back into bn_param
218     bn_param['running_mean'] = running_mean
219     bn_param['running_var'] = running_var
220
221     return out, cache
222
223 def batchnorm_backward(dout, cache):
224     """
225     Backward pass for batch normalization.
226
227     For this implementation, you should write out a computation graph for
228     batch normalization on paper and propagate gradients backward through
229     intermediate nodes.
230
231     Inputs:
232     - dout: Upstream derivatives, of shape (N, D)
233     - cache: Variable of intermediates from batchnorm_forward.
234
235     Returns a tuple of:
236     - dx: Gradient with respect to inputs x, of shape (N, D)
237     - dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
238     - dbeta: Gradient with respect to shift parameter beta, of shape (D,)
239     """
240     dx, dgamma, dbeta = None, None, None
241
242     # ===== #
243     # YOUR CODE HERE:
244     # Implement the batchnorm backward pass, calculating dx, dgamma, and dbeta.
245     # ===== #
246     x, x_n, mean, var, eps, gamma = cache
247     dbeta = np.sum(dout, axis=0)
248     dgamma = np.sum(dout*x_n, axis=0)
249     dxh = dout*gamma
250     std = np.sqrt(var+eps)
251     du = -np.sum((1/std)*dxh, axis=0)
252     dv = -np.sum((1/(2*((std)**3)))*(x-mean)*dxh, axis=0)
253     dx = dxh*(1/std) + du/x.shape[0] + (2*(x-mean)/x.shape[0])*dv
254     # ===== #
255     # END YOUR CODE HERE
256     # ===== #
257
258     return dx, dgamma, dbeta
259
260 def dropout_forward(x, dropout_param):
261     """
262     Performs the forward pass for (inverted) dropout.
263
264     Inputs:
265     - x: Input data, of any shape
266     - dropout_param: A dictionary with the following keys:
267         - p: Dropout parameter. We drop each neuron output with probability p.
268         - mode: 'test' or 'train'. If the mode is train, then perform dropout;
269             if the mode is test, then just return the input.
270         - seed: Seed for the random number generator. Passing seed makes this
271             function deterministic, which is needed for gradient checking but not in
272             real networks.
273
274     Outputs:
275     - out: Array of the same shape as x.
276     - cache: A tuple (dropout_param, mask). In training mode, mask is the dropout
277         mask that was used to multiply the input; in test mode, mask is None.
278     """
279     p, mode = dropout_param['p'], dropout_param['mode']
280     if 'seed' in dropout_param:
281         np.random.seed(dropout_param['seed'])
282
283     mask = None
284     out = None
285
286     if mode == 'train':

```

```

287 # ===== #
288 # YOUR CODE HERE:
289 #   Implement the inverted dropout forward pass during training time.
290 #   Store the masked and scaled activations in out, and store the
291 #   dropout mask as the variable mask.
292 # ===== #
293
294 mask = (np.random.rand(*x.shape) < (1-p)) / (1-p)
295 out = x * mask
296
297 # ===== #
298 # END YOUR CODE HERE
299 # ===== #
300
301 elif mode == 'test':
302
303 # ===== #
304 # YOUR CODE HERE:
305 #   Implement the inverted dropout forward pass during test time.
306 # ===== #
307
308 out = x
309
310 # ===== #
311 # END YOUR CODE HERE
312 # ===== #
313
314 cache = (dropout_param, mask)
315 out = out.astype(x.dtype, copy=False)
316
317 return out, cache
318
319 def dropout_backward(dout, cache):
320 """
321 Perform the backward pass for (inverted) dropout.
322
323 Inputs:
324 - dout: Upstream derivatives, of any shape
325 - cache: (dropout_param, mask) from dropout_forward.
326 """
327 dropout_param, mask = cache
328 mode = dropout_param['mode']
329
330 dx = None
331 if mode == 'train':
332 # ===== #
333 # YOUR CODE HERE:
334 #   Implement the inverted dropout backward pass during training time.
335 # ===== #
336
337 (dropout_param, mask) = cache
338 dx = dout * mask
339
340 # ===== #
341 # END YOUR CODE HERE
342 # ===== #
343 elif mode == 'test':
344 # ===== #
345 # YOUR CODE HERE:
346 #   Implement the inverted dropout backward pass during test time.
347 # ===== #
348
349 dx = dout
350
351 # ===== #
352 # END YOUR CODE HERE
353 # ===== #
354 return dx
355
356 def svm_loss(x, y):
357 """
358 Computes the loss and gradient using for multiclass SVM classification.

```

```

359
360     Inputs:
361     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
362         for the ith input.
363     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
364         0 <= y[i] < C
365
366     Returns a tuple of:
367     - loss: Scalar giving the loss
368     - dx: Gradient of the loss with respect to x
369     """
370     N = x.shape[0]
371     correct_class_scores = x[np.arange(N), y]
372     margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
373     margins[np.arange(N), y] = 0
374     loss = np.sum(margins) / N
375     num_pos = np.sum(margins > 0, axis=1)
376     dx = np.zeros_like(x)
377     dx[margins > 0] = 1
378     dx[np.arange(N), y] -= num_pos
379     dx /= N
380     return loss, dx
381
382
383 def softmax_loss(x, y):
384     """
385     Computes the loss and gradient for softmax classification.
386
387     Inputs:
388     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
389         for the ith input.
390     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
391         0 <= y[i] < C
392
393     Returns a tuple of:
394     - loss: Scalar giving the loss
395     - dx: Gradient of the loss with respect to x
396     """
397
398     probs = np.exp(x - np.max(x, axis=1, keepdims=True))
399     probs /= np.sum(probs, axis=1, keepdims=True)
400     N = x.shape[0]
401     loss = -np.sum(np.log(probs[np.arange(N), y])) / N
402     dx = probs.copy()
403     dx[np.arange(N), y] -= 1
404     dx /= N
405     return loss, dx
406

```

```

1  from .layers import *
2
3
4  def affine_relu_forward(x, w, b):
5      """
6          Convenience layer that performs an affine transform followed by a ReLU
7
8          Inputs:
9          - x: Input to the affine layer
10         - w, b: Weights for the affine layer
11
12        Returns a tuple of:
13        - out: Output from the ReLU
14        - cache: Object to give to the backward pass
15    """
16    a, fc_cache = affine_forward(x, w, b)
17    out, relu_cache = relu_forward(a)
18    cache = (fc_cache, relu_cache)
19    return out, cache
20
21  def affine_relu_backward(dout, cache):
22      """
23          Backward pass for the affine-relu convenience layer
24      """
25    fc_cache, relu_cache = cache
26    da = relu_backward(dout, relu_cache)
27    dx, dw, db = affine_backward(da, fc_cache)
28    return dx, dw, db
29
30  def affine_batchnorm_relu_forward(x, w, b, gamma, beta, bn_params):
31      """
32          Convenience layer that performs an affine transform followed by batchnorm and a
33          ReLU
34
35          Inputs:
36          - x: Input to the affine layer
37          - w, b: Weights for the affine layer
38          - gamma, beta, bn_params: params for batchnorm
39
40          Returns a tuple of:
41          - out: Output from the ReLU
42          - cache: Object to give to the backward pass
43      """
44    a, fc_cache = affine_forward(x, w, b)
45    b, bn_cache = batchnorm_forward(a, gamma, beta, bn_params)
46    out, relu_cache = relu_forward(b)
47    cache = (fc_cache, bn_cache, relu_cache)
48    return out, cache
49
50  def affine_batchnorm_relu_backward(dout, cache):
51      """
52          Backward pass for the affine-batchnorm-relu convenience layer
53      """
54    fc_cache, bn_cache, relu_cache = cache
55    da = relu_backward(dout, relu_cache)
56    dbn, dgamma, dbeta = batchnorm_backward(da, bn_cache)
57    dx, dw, db = affine_backward(dbn, fc_cache)
58    return dx, dw, db, dgamma, dbeta
59
```

```

1 import numpy as np
2 from .layers import *
3 from .layer_utils import *
4
5
6 class TwoLayerNet(object):
7     """
8         A two-layer fully-connected neural network with ReLU nonlinearity and
9         softmax loss that uses a modular layer design. We assume an input dimension
10        of D, a hidden dimension of H, and perform classification over C classes.
11
12        The architecture should be affine - relu - affine - softmax.
13
14        Note that this class does not implement gradient descent; instead, it
15        will interact with a separate Solver object that is responsible for running
16        optimization.
17
18        The learnable parameters of the model are stored in the dictionary
19        self.params that maps parameter names to numpy arrays.
20        """
21
22    def __init__(self, input_dim=3*32*32, hidden_dims=100, num_classes=10,
23                 dropout=0, weight_scale=1e-3, reg=0.0):
24        """
25            Initialize a new network.
26
27            Inputs:
28            - input_dim: An integer giving the size of the input
29            - hidden_dims: An integer giving the size of the hidden layer
30            - num_classes: An integer giving the number of classes to classify
31            - dropout: Scalar between 0 and 1 giving dropout strength.
32            - weight_scale: Scalar giving the standard deviation for random
33                initialization of the weights.
34            - reg: Scalar giving L2 regularization strength.
35        """
36        self.params = {}
37        self.reg = reg
38
39        # ===== #
40        # YOUR CODE HERE:
41        #   Initialize W1, W2, b1, and b2. Store these as self.params['W1'],
42        #   self.params['W2'], self.params['b1'] and self.params['b2']. The
43        #   biases are initialized to zero and the weights are initialized
44        #   so that each parameter has mean 0 and standard deviation weight_scale.
45        #   The dimensions of W1 should be (input_dim, hidden_dim) and the
46        #   dimensions of W2 should be (hidden_dims, num_classes)
47        # ===== #
48
49
50        # ===== #
51        # END YOUR CODE HERE
52        # ===== #
53
54    def loss(self, X, y=None):
55        """
56            Compute loss and gradient for a minibatch of data.
57
58            Inputs:
59            - X: Array of input data of shape (N, d_1, ..., d_k)
60            - y: Array of labels, of shape (N,). y[i] gives the label for X[i].
61
62            Returns:
63            If y is None, then run a test-time forward pass of the model and return:
64            - scores: Array of shape (N, C) giving classification scores, where
65                scores[i, c] is the classification score for X[i] and class c.
66
67            If y is not None, then run a training-time forward and backward pass and
68            return a tuple of:
69            - loss: Scalar value giving the loss
70            - grads: Dictionary with the same keys as self.params, mapping parameter
71                names to gradients of the loss with respect to those parameters.
72        """

```

```

73     scores = None
74
75     # ===== #
76     # YOUR CODE HERE:
77     #   Implement the forward pass of the two-layer neural network. Store
78     #   the class scores as the variable 'scores'. Be sure to use the layers
79     #   you prior implemented.
80     # ===== #
81
82     # ===== #
83     # END YOUR CODE HERE
84     # ===== #
85
86     # If y is None then we are in test mode so just return scores
87     if y is None:
88         return scores
89
90     loss, grads = 0, {}
91     # ===== #
92     # YOUR CODE HERE:
93     #   Implement the backward pass of the two-layer neural net. Store
94     #   the loss as the variable 'loss' and store the gradients in the
95     #   'grads' dictionary. For the grads dictionary, grads['W1'] holds
96     #   the gradient for W1, grads['b1'] holds the gradient for b1, etc.
97     #   i.e., grads[k] holds the gradient for self.params[k].
98     #
99     # Add L2 regularization, where there is an added cost 0.5*self.reg*W^2
100    # for each W. Be sure to include the 0.5 multiplying factor to
101    # match our implementation.
102    #
103    # And be sure to use the layers you prior implemented.
104    # ===== #
105
106    # ===== #
107    # END YOUR CODE HERE
108    # ===== #
109
110    return loss, grads
111
112
113 class FullyConnectedNet(object):
114     """
115     A fully-connected neural network with an arbitrary number of hidden layers,
116     ReLU nonlinearities, and a softmax loss function. This will also implement
117     dropout and batch normalization as options. For a network with L layers,
118     the architecture will be
119
120     {affine - [batch norm] - relu - [dropout]} x (L - 1) - affine - softmax
121
122     where batch normalization and dropout are optional, and the {...} block is
123     repeated L - 1 times.
124
125     Similar to the TwoLayerNet above, learnable parameters are stored in the
126     self.params dictionary and will be learned using the Solver class.
127     """
128
129     def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,
130                  dropout=0, use_batchnorm=False, reg=0.0,
131                  weight_scale=1e-2, dtype=np.float32, seed=None):
132         """
133             Initialize a new FullyConnectedNet.
134
135             Inputs:
136             - hidden_dims: A list of integers giving the size of each hidden layer.
137             - input_dim: An integer giving the size of the input.
138             - num_classes: An integer giving the number of classes to classify.
139             - dropout: Scalar between 0 and 1 giving dropout strength. If dropout=0 then
140                 the network should not use dropout at all.
141             - use_batchnorm: Whether or not the network should use batch normalization.
142             - reg: Scalar giving L2 regularization strength.
143             - weight_scale: Scalar giving the standard deviation for random
144                 initialization of the weights.

```

```

145     - dtype: A numpy datatype object; all computations will be performed using
146       this datatype. float32 is faster but less accurate, so you should use
147       float64 for numeric gradient checking.
148     - seed: If not None, then pass this random seed to the dropout layers. This
149       will make the dropout layers deterministic so we can gradient check the
150       model.
151     """
152     self.use_batchnorm = use_batchnorm
153     self.use_dropout = dropout > 0
154     self.reg = reg
155     self.num_layers = 1 + len(hidden_dims)
156     self.dtype = dtype
157     self.params = {}
158
159     # ===== #
160     # YOUR CODE HERE:
161     #   Initialize all parameters of the network in the self.params dictionary.
162     #   The weights and biases of layer 1 are  $W_1$  and  $b_1$ ; and in general the
163     #   weights and biases of layer  $i$  are  $W_i$  and  $b_i$ . The
164     #   biases are initialized to zero and the weights are initialized
165     #   so that each parameter has mean 0 and standard deviation weight_scale.
166     #
167     #   BATCHNORM: Initialize the gammas of each layer to 1 and the beta
168     #   parameters to zero. The gamma and beta parameters for layer 1 should
169     #   be self.params['gamma1'] and self.params['beta1']. For layer 2, they
170     #   should be gamma2 and beta2, etc. Only use batchnorm if self.use_batchnorm
171     #   is true and DO NOT do batch normalize the output scores.
172     # ===== #
173     prev_dim = input_dim
174     for i in range(self.num_layers):
175         w = 'W'+str(i+1)
176         b = 'b'+str(i+1)
177         if i == len(hidden_dims):
178             self.params[w] = np.random.normal(loc=0.0, scale=weight_scale, size =
179                 (prev_dim, num_classes))
180             self.params[b] = np.zeros((num_classes))
181         else:
182             self.params[w] = np.random.normal(loc=0.0, scale=weight_scale, size =
183                 (prev_dim, hidden_dims[i]))
184             self.params[b] = np.zeros((hidden_dims[i]))
185             prev_dim = hidden_dims[i]
186             if self.use_batchnorm:
187                 self.params['gamma'+str(i+1)] = np.ones((hidden_dims[i]))
188                 self.params['beta'+str(i+1)] = np.zeros((hidden_dims[i]))
189
190     # ===== #
191     # END YOUR CODE HERE
192     # ===== #
193
194     # When using dropout we need to pass a dropout_param dictionary to each
195     # dropout layer so that the layer knows the dropout probability and the mode
196     # (train / test). You can pass the same dropout_param to each dropout layer.
197     self.dropout_param = {}
198     if self.use_dropout:
199         self.dropout_param = {'mode': 'train', 'p': dropout}
200         if seed is not None:
201             self.dropout_param['seed'] = seed
202
203         # With batch normalization we need to keep track of running means and
204         # variances, so we need to pass a special bn_param object to each batch
205         # normalization layer. You should pass self.bn_params[0] to the forward pass
206         # of the first batch normalization layer, self.bn_params[1] to the forward
207         # pass of the second batch normalization layer, etc.
208         self.bn_params = []
209         if self.use_batchnorm:
210             self.bn_params = [{ 'mode': 'train' } for i in np.arange(self.num_layers - 1)]
211
212         # Cast all parameters to the correct datatype
213         for k, v in self.params.items():
214             self.params[k] = v.astype(dtype)

```

```

215
216     def loss(self, X, y=None):
217         """
218             Compute loss and gradient for the fully-connected net.
219
220             Input / output: Same as TwoLayerNet above.
221         """
222         X = X.astype(self.dtype)
223         mode = 'test' if y is None else 'train'
224
225         # Set train/test mode for batchnorm params and dropout param since they
226         # behave differently during training and testing.
227         if self.dropout_param is not None:
228             self.dropout_param['mode'] = mode
229         if self.use_batchnorm:
230             for bn_param in self.bn_params:
231                 bn_param['mode'] = mode
232
233         scores = None
234
235         # ===== #
236         # YOUR CODE HERE:
237         #   Implement the forward pass of the FC net and store the output
238         #   scores as the variable "scores".
239         #
240         # BATCHNORM: If self.use_batchnorm is true, insert a batchnorm layer
241         # between the affine_forward and relu_forward layers. You may
242         # also write an affine_batchnorm_relu() function in layer_utils.py.
243         #
244         # DROPOUT: If dropout is non-zero, insert a dropout layer after
245         # every ReLU layer.
246         # ===== #
247         act = []
248         cache = []
249         masks = []
250         for i in range(self.num_layers):
251             w = 'W'+str(i+1)
252             b = 'b'+str(i+1)
253             gamma = 'gamma'+str(i+1)
254             beta = 'beta'+str(i+1)
255             temp_c = []
256             if i==0:
257                 if self.use_batchnorm:
258                     a,c =
259                         affine_batchnorm_relu_forward(X,self.params[w],self.params[b],self.params[gamma],self.params[beta],self.bn_params[i])
260                 else:
261                     a,c = affine_relu_forward(X,self.params[w],self.params[b])
262                     temp_c.append(c)
263                     if self.use_dropout:
264                         a,dc = dropout_forward(a,self.dropout_param)
265                         temp_c.append(dc)
266                     act.append(a)
267             elif i==self.num_layers-1:
268                 a,c = affine_forward(act[i-1],self.params[w],self.params[b])
269                 temp_c.append(c)
270                 scores=a
271             else:
272                 if self.use_batchnorm:
273                     a,c =
274                         affine_batchnorm_relu_forward(act[i-1],self.params[w],self.params[b],self.params[gamma],self.params[beta],self.bn_params[i])
275                 else:
276                     a,c = affine_relu_forward(act[i-1],self.params[w],self.params[b])
277                     temp_c.append(c)
278                     if self.use_dropout:
279                         a,dc = dropout_forward(a,self.dropout_param)
280                         temp_c.append(dc)
281                     act.append(a)
282                     temp_c.append(c)
283                     cache.append(temp_c)

```

```

283 # ===== #
284 # END YOUR CODE HERE
285 # ===== #
286
287 # If test mode return early
288 if mode == 'test':
289     return scores
290
291 loss, grads = 0.0, {}
292 # ===== #
293 # YOUR CODE HERE:
294 # Implement the backwards pass of the FC net and store the gradients
295 # in the grads dict, so that grads[k] is the gradient of self.params[k]
296 # Be sure your L2 regularization includes a 0.5 factor.
297 #
298 # BATCHNORM: Incorporate the backward pass of the batchnorm.
299 #
300 # DROPOUT: Incorporate the backward pass of dropout.
301 # ===== #
302 loss=0
303 loss, ds = softmax_loss(scores,y)
304 for i in np.arange(self.num_layers-1,-1,-1):
305     w = 'W'+str(i+1)
306     b = 'b'+str(i+1)
307     gamma = 'gamma'+str(i+1)
308     beta = 'beta'+str(i+1)
309     loss += self.reg*0.5*np.sum(np.square(self.params[w]))
310     if i==self.num_layers-1:
311         da,grads[w], grads[b] = affine_backward(ds,cache[i][0])
312     else:
313         if self.use_dropout:
314             da = dropout_backward(da,cache[i][1])
315         if self.use_batchnorm:
316             da,grads[w], grads[b], grads[gamma], grads[beta] =
317                 affine_batchnorm_relu_backward(da,cache[i][0])
318         else:
319             da,grads[w], grads[b] = affine_relu_backward(da,cache[i][0])
320     grads[w] += self.reg*self.params[w]
321
322 # ===== #
323 # END YOUR CODE HERE
324 # ===== #
325
326 return loss, grads
327

```