

# Optimization

February 9, 2022

## 1 Sunay Bhat W2022 HW4 C247

## 2 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`.

```
[2]: ## 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))))
```

```
[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,)
```

## 2.1 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

### 2.1.1 Test all functions you copy and pasted

```
[4]: 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: 2.675168729800476e-10  
dw error: 4.5751263918702534e-11  
db error: 1.0489400292828403e-11

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

difference: 4.999999798022158e-08

If relu\_forward function is working, error should be less than 1e-9:  
dx error: 3.2756077382282424e-12

If affine\_relu\_forward and affine\_relu\_backward are working, error should be less than 1e-9::  
dx error: 2.1636207865197628e-10  
dw error: 4.0634236850757724e-10  
db error: 3.2756095922011834e-12

Running check with reg = 0  
Initial loss: 2.316413617603716  
W1 relative error: 2.2472746956718977e-06  
W2 relative error: 1.81258129736145e-07  
W3 relative error: 7.210768694035575e-08  
b1 relative error: 1.1357459070280958e-08  
b2 relative error: 1.361352057975109e-08  
b3 relative error: 1.0847308582773077e-10  
Running check with reg = 3.14  
Initial loss: 6.902940007165628  
W1 relative error: 1.0949097055726692e-08  
W2 relative error: 1.970241029245354e-08  
W3 relative error: 1.1413874871688892e-08  
b1 relative error: 3.013348285170618e-08  
b2 relative error: 3.4054435276750802e-09  
b3 relative error: 2.4670371432898317e-10

### 3 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.

#### 3.1 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.

```
[7]: 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

```

### 3.2 SGD + Nesterov momentum

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

```

[8]: 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

```

### 3.3 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.

```

[10]: 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)

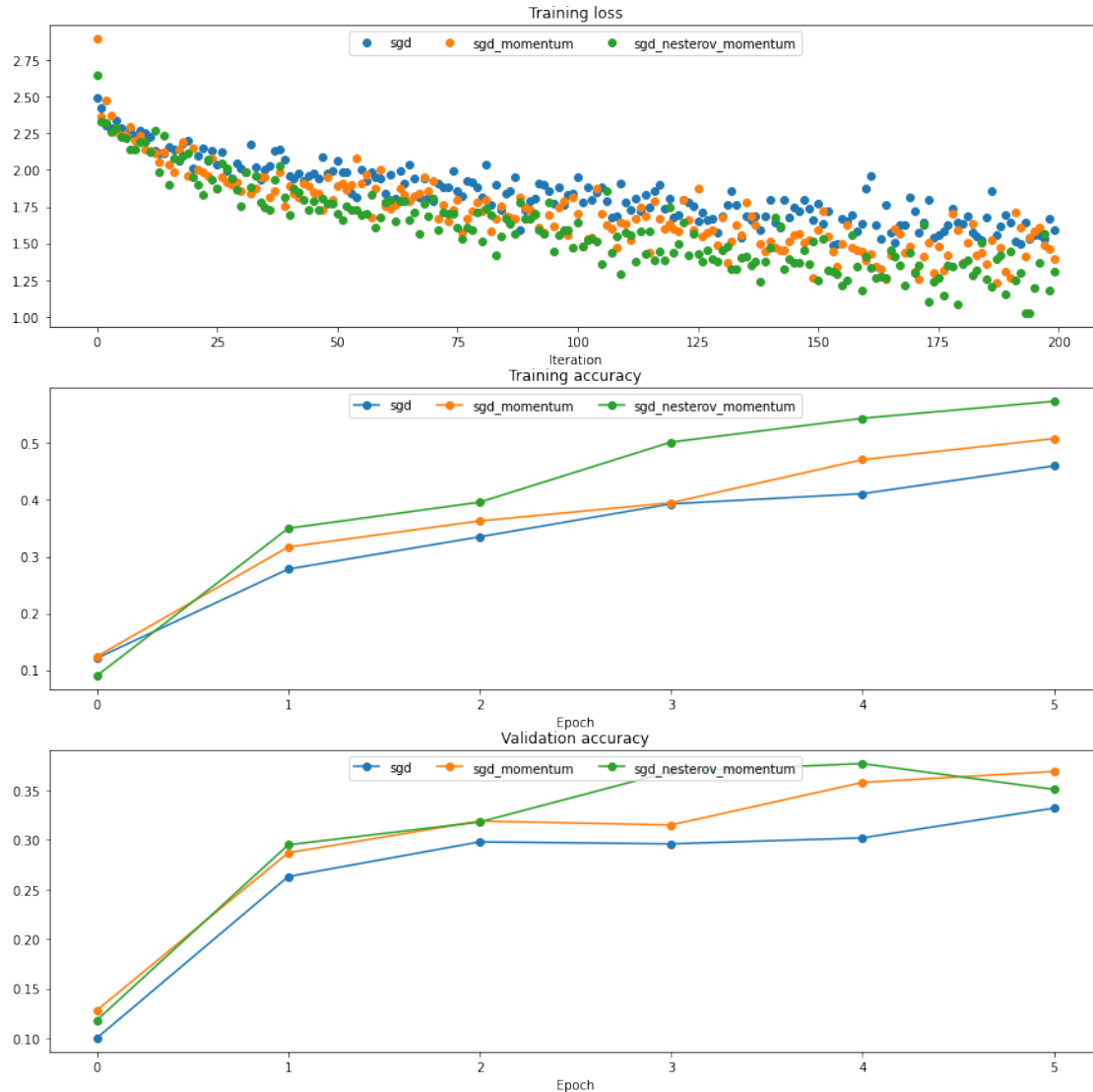
    plt.subplot(3, 1, 2)
    plt.plot(solver.train_acc_history, '-o', label=update_rule)

    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`



### 3.4 RMSProp

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

```
[11]: 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

```

### 3.5 Adaptive moments

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

```

[16]: # 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

```

### 3.6 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.

```

[17]: 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)

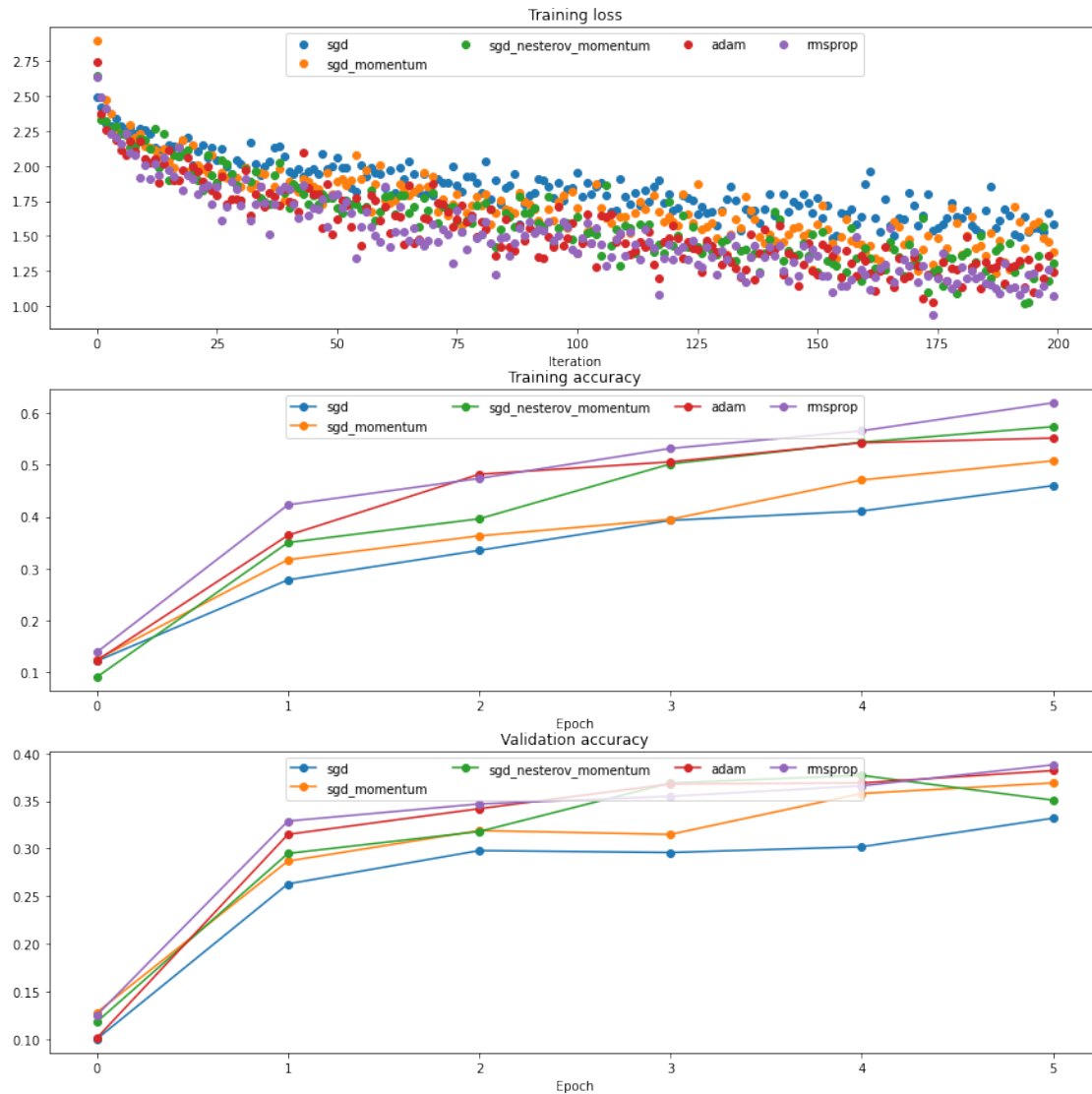
plt.subplot(3, 1, 3)
plt.plot(solver.val_acc_history, '-o', label=update_rule)

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



### 3.7 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+% on CIFAR-10.

```
[18]: 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: 2.328181
(Epoch 0 / 10) train acc: 0.111000; val_acc: 0.133000
(Iteration 51 / 4900) loss: 1.797477
(Iteration 101 / 4900) loss: 1.675898
(Iteration 151 / 4900) loss: 1.727979
(Iteration 201 / 4900) loss: 1.663434
(Iteration 251 / 4900) loss: 1.600791
(Iteration 301 / 4900) loss: 1.877411
(Iteration 351 / 4900) loss: 1.550363
(Iteration 401 / 4900) loss: 1.666523
(Iteration 451 / 4900) loss: 1.634747
(Epoch 1 / 10) train acc: 0.432000; val_acc: 0.416000
(Iteration 501 / 4900) loss: 1.745033
(Iteration 551 / 4900) loss: 1.468143
(Iteration 601 / 4900) loss: 1.722282
(Iteration 651 / 4900) loss: 1.660672
(Iteration 701 / 4900) loss: 1.701815
(Iteration 751 / 4900) loss: 1.529541
```

(Iteration 801 / 4900) loss: 1.438237  
(Iteration 851 / 4900) loss: 1.741340  
(Iteration 901 / 4900) loss: 1.583205  
(Iteration 951 / 4900) loss: 1.363590  
(Epoch 2 / 10) train acc: 0.485000; val\_acc: 0.468000  
(Iteration 1001 / 4900) loss: 1.431547  
(Iteration 1051 / 4900) loss: 1.718218  
(Iteration 1101 / 4900) loss: 1.316846  
(Iteration 1151 / 4900) loss: 1.379579  
(Iteration 1201 / 4900) loss: 1.423866  
(Iteration 1251 / 4900) loss: 1.591148  
(Iteration 1301 / 4900) loss: 1.342587  
(Iteration 1351 / 4900) loss: 1.503359  
(Iteration 1401 / 4900) loss: 1.142978  
(Iteration 1451 / 4900) loss: 1.567900  
(Epoch 3 / 10) train acc: 0.506000; val\_acc: 0.517000  
(Iteration 1501 / 4900) loss: 1.248408  
(Iteration 1551 / 4900) loss: 1.344610  
(Iteration 1601 / 4900) loss: 1.365671  
(Iteration 1651 / 4900) loss: 1.105743  
(Iteration 1701 / 4900) loss: 1.458173  
(Iteration 1751 / 4900) loss: 1.290630  
(Iteration 1801 / 4900) loss: 1.407816  
(Iteration 1851 / 4900) loss: 1.232504  
(Iteration 1901 / 4900) loss: 1.153526  
(Iteration 1951 / 4900) loss: 1.234272  
(Epoch 4 / 10) train acc: 0.542000; val\_acc: 0.501000  
(Iteration 2001 / 4900) loss: 1.300892  
(Iteration 2051 / 4900) loss: 1.210190  
(Iteration 2101 / 4900) loss: 1.410650  
(Iteration 2151 / 4900) loss: 1.131898  
(Iteration 2201 / 4900) loss: 1.463518  
(Iteration 2251 / 4900) loss: 1.187098  
(Iteration 2301 / 4900) loss: 0.996405  
(Iteration 2351 / 4900) loss: 1.310067  
(Iteration 2401 / 4900) loss: 1.176746  
(Epoch 5 / 10) train acc: 0.568000; val\_acc: 0.510000  
(Iteration 2451 / 4900) loss: 1.322918  
(Iteration 2501 / 4900) loss: 1.195075  
(Iteration 2551 / 4900) loss: 1.012166  
(Iteration 2601 / 4900) loss: 1.093007  
(Iteration 2651 / 4900) loss: 1.414172  
(Iteration 2701 / 4900) loss: 1.218426  
(Iteration 2751 / 4900) loss: 0.984611  
(Iteration 2801 / 4900) loss: 1.178021  
(Iteration 2851 / 4900) loss: 1.183064  
(Iteration 2901 / 4900) loss: 1.375793  
(Epoch 6 / 10) train acc: 0.590000; val\_acc: 0.517000

```

(Iteration 2951 / 4900) loss: 1.241499
(Iteration 3001 / 4900) loss: 1.120312
(Iteration 3051 / 4900) loss: 1.253601
(Iteration 3101 / 4900) loss: 1.033850
(Iteration 3151 / 4900) loss: 1.190524
(Iteration 3201 / 4900) loss: 1.256515
(Iteration 3251 / 4900) loss: 1.346379
(Iteration 3301 / 4900) loss: 1.221121
(Iteration 3351 / 4900) loss: 0.930110
(Iteration 3401 / 4900) loss: 0.902628
(Epoch 7 / 10) train acc: 0.609000; val_acc: 0.512000
(Iteration 3451 / 4900) loss: 1.125307
(Iteration 3501 / 4900) loss: 0.982823
(Iteration 3551 / 4900) loss: 1.093506
(Iteration 3601 / 4900) loss: 0.961197
(Iteration 3651 / 4900) loss: 1.136582
(Iteration 3701 / 4900) loss: 0.921714
(Iteration 3751 / 4900) loss: 1.174514
(Iteration 3801 / 4900) loss: 0.932937
(Iteration 3851 / 4900) loss: 1.150155
(Iteration 3901 / 4900) loss: 0.833050
(Epoch 8 / 10) train acc: 0.636000; val_acc: 0.529000
(Iteration 3951 / 4900) loss: 0.956006
(Iteration 4001 / 4900) loss: 0.977857
(Iteration 4051 / 4900) loss: 0.992840
(Iteration 4101 / 4900) loss: 1.016208
(Iteration 4151 / 4900) loss: 0.886169
(Iteration 4201 / 4900) loss: 1.030166
(Iteration 4251 / 4900) loss: 0.976255
(Iteration 4301 / 4900) loss: 0.995391
(Iteration 4351 / 4900) loss: 0.812811
(Iteration 4401 / 4900) loss: 1.104995
(Epoch 9 / 10) train acc: 0.685000; val_acc: 0.533000
(Iteration 4451 / 4900) loss: 1.174991
(Iteration 4501 / 4900) loss: 0.851488
(Iteration 4551 / 4900) loss: 0.877790
(Iteration 4601 / 4900) loss: 1.069405
(Iteration 4651 / 4900) loss: 0.914333
(Iteration 4701 / 4900) loss: 0.774547
(Iteration 4751 / 4900) loss: 0.965120
(Iteration 4801 / 4900) loss: 0.937213
(Iteration 4851 / 4900) loss: 0.781859
(Epoch 10 / 10) train acc: 0.709000; val_acc: 0.536000

```

```

[19]: 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: 0.536

Test set accuracy: 0.523

[ ]:

# Batch-Normalization

February 9, 2022

## 1 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`.

```
[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))))
```

```
[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,)
```

## 1.1 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.

```
[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: [18.08245532 14.30449559 33.41358387]
  stds:  [29.32085642 38.28340501 30.98650269]
```



```

After batch normalization (gamma=1, beta=0)
mean:  [ 1.30173650e-16 -1.88182803e-16 -4.57134330e-16]
std:   [0.99999999 1.          0.99999999]
After batch normalization (nontrivial gamma, beta)
means:  [11. 12. 13.]
stds:   [0.99999999 1.99999999 2.99999998]

```

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

```

[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:  [0.03389563 0.04485845 0.10096212]
stds:   [1.1150012  0.93355312 0.94601145]

```

## 1.2 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.

```

[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:  4.102799594114112e-10
dgamma error:  4.8208809343424334e-11
dbeta error:  3.276959277354581e-12

```

### 1.3 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  $1e-4$ .

```

[11]: 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: 2.2244131561672003
W1 relative error: 1.910603882187157e-05
W2 relative error: 1.5530850127155482e-05
W3 relative error: 4.0381176170850605e-10
b1 relative error: 0.0022204460492503126
b2 relative error: 4.440892098500626e-07
b3 relative error: 1.2766396166406513e-10
beta1 relative error: 4.080835971806538e-09
beta2 relative error: 1.090036130965299e-08
gamma1 relative error: 9.107694431370659e-09
gamma2 relative error: 3.662332876020335e-08

```

```

Running check with reg = 3.14
Initial loss: 6.916972135377012
W1 relative error: 3.4878789210583476e-06
W2 relative error: 3.6523344204693075e-07
W3 relative error: 2.3927497037876514e-08
b1 relative error: 1.3183898417423734e-08
b2 relative error: 4.440892098500626e-08
b3 relative error: 1.9828315572957545e-10
beta1 relative error: 1.6796183376744786e-07
beta2 relative error: 4.351905510834597e-09
gamma1 relative error: 3.066521454196165e-07
gamma2 relative error: 4.6196135505950856e-09

```

## 1.4 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.

```

[13]: # 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: 2.336318
(Epoch 0 / 10) train acc: 0.113000; val_acc: 0.091000
(Epoch 1 / 10) train acc: 0.342000; val_acc: 0.276000
(Epoch 2 / 10) train acc: 0.427000; val_acc: 0.317000
(Epoch 3 / 10) train acc: 0.453000; val_acc: 0.321000
(Epoch 4 / 10) train acc: 0.540000; val_acc: 0.339000
(Epoch 5 / 10) train acc: 0.559000; val_acc: 0.326000
(Epoch 6 / 10) train acc: 0.621000; val_acc: 0.351000
(Epoch 7 / 10) train acc: 0.667000; val_acc: 0.324000
(Epoch 8 / 10) train acc: 0.712000; val_acc: 0.342000
(Epoch 9 / 10) train acc: 0.757000; val_acc: 0.328000
(Epoch 10 / 10) train acc: 0.814000; val_acc: 0.343000
(Iteration 1 / 200) loss: 2.303048
(Epoch 0 / 10) train acc: 0.143000; val_acc: 0.140000
(Epoch 1 / 10) train acc: 0.269000; val_acc: 0.239000

```

```
(Epoch 2 / 10) train acc: 0.340000; val_acc: 0.254000
(Epoch 3 / 10) train acc: 0.364000; val_acc: 0.283000
(Epoch 4 / 10) train acc: 0.416000; val_acc: 0.302000
(Epoch 5 / 10) train acc: 0.477000; val_acc: 0.328000
(Epoch 6 / 10) train acc: 0.484000; val_acc: 0.329000
(Epoch 7 / 10) train acc: 0.563000; val_acc: 0.316000
(Epoch 8 / 10) train acc: 0.597000; val_acc: 0.326000
(Epoch 9 / 10) train acc: 0.619000; val_acc: 0.293000
(Epoch 10 / 10) train acc: 0.667000; val_acc: 0.324000
```

```
[14]: plt.subplot(3, 1, 1)
plt.title('Training loss')
plt.xlabel('Iteration')

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

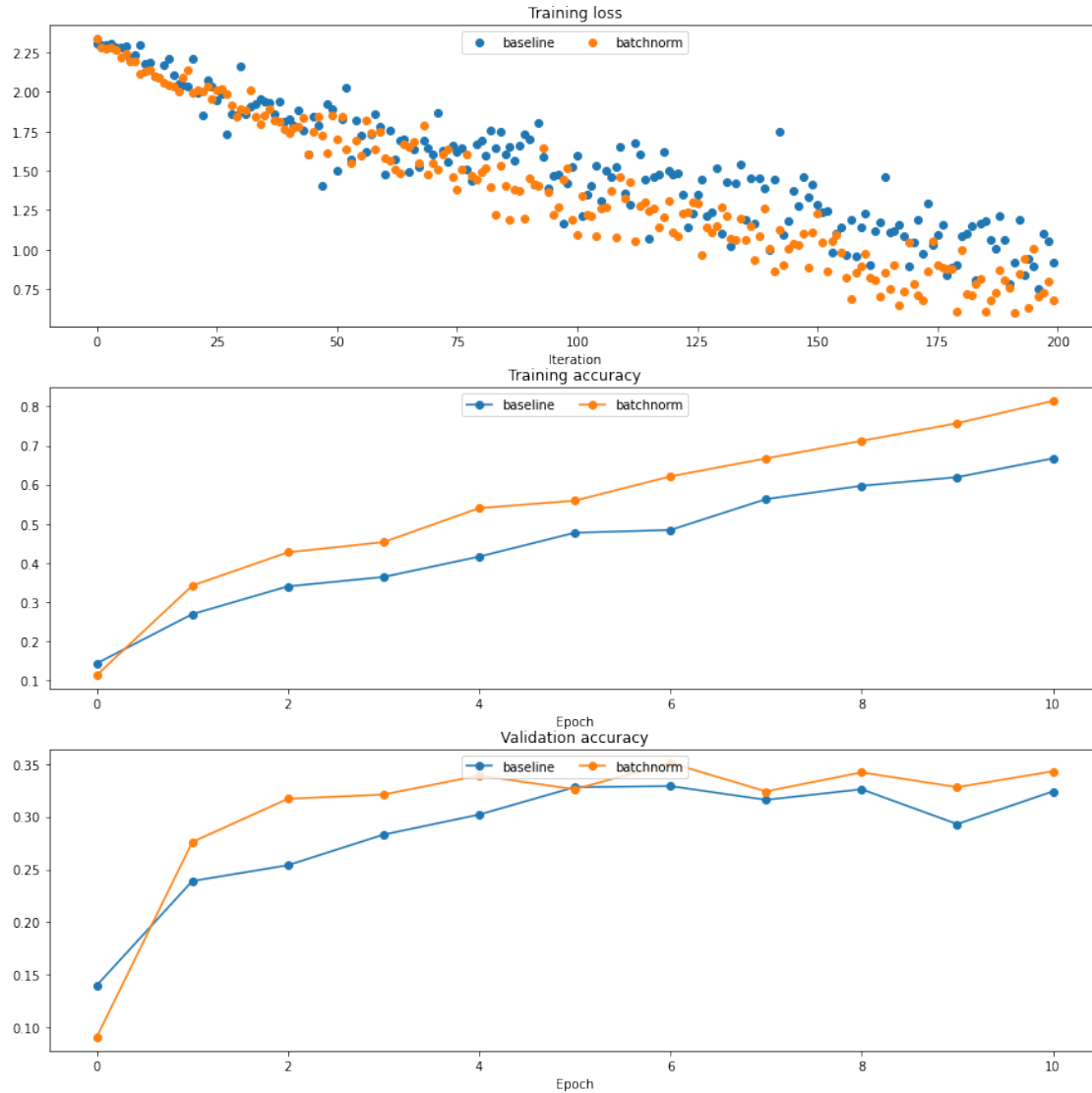
plt.subplot(3, 1, 3)
plt.title('Validation accuracy')
plt.xlabel('Epoch')

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()
```



## 1.5 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.

```
[15]: # 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

```

```

[16]: # 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()

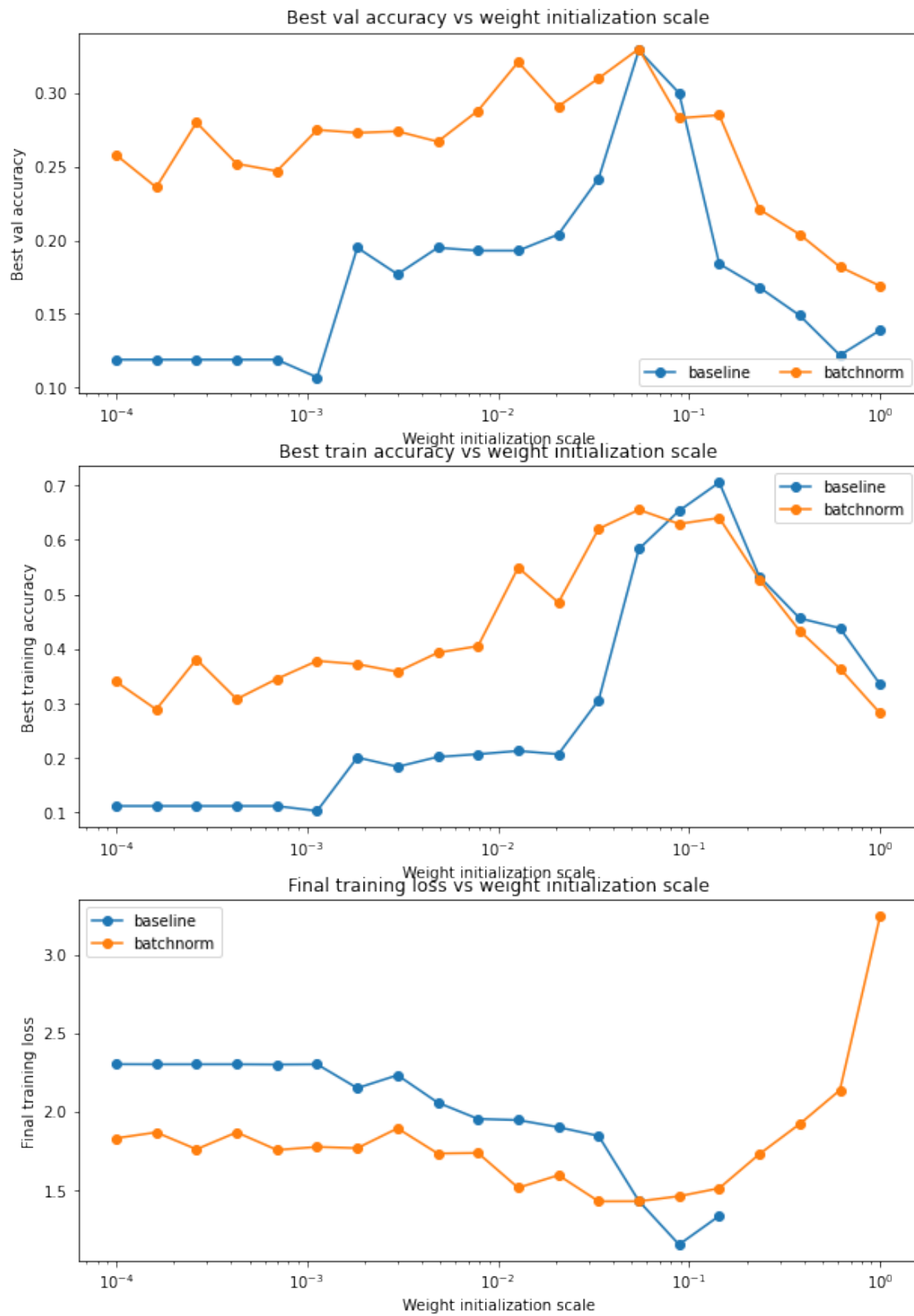
plt.subplot(3, 1, 3)
plt.title('Final training loss vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Final training loss')
plt.semilogx(weight_scales, final_train_loss, '-o', label='baseline')
plt.semilogx(weight_scales, bn_final_train_loss, '-o', label='batchnorm')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

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

```



```
plt.show()
```



### 1.6 Question:

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

### 1.7 Answer:

**Across weight initializations, batchnorm generally performs better and is more stable to changes in the initialization.**

When our weight initializations were low ( $< 10^{-1}$ ), the baseline network likely has a **vanishing gradient problem** with weights going to zero. When the weights are initialized larger, we still get decent training accuracy, but **poorer generalization to val accuracy probably due to exploding gradients/weights**. Thus, the baseline is only comparable - and slightly better - than batchnorm in a narrow window of weight initializations around  $10^{-1}$ . Overall batchnorm is a far more robust technique.

[ ]:

import numpy as np

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

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

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

Returns:  
- next\_w: The next point after the update.  
- config: The config dictionary to be passed to the next iteration of the  
  update rule.

NOTE: For most update rules, the default learning rate will probably not perform  
well; however the default values of the other hyperparameters should work well  
for a variety of different problems.

For efficiency, update rules may perform in-place updates, mutating w and  
setting next\_w equal to w.

def sgd(w, dw, config=None):  
 """  
 Performs vanilla stochastic gradient descent.  
  
 config format:  
 - learning\_rate: Scalar learning rate.  
 """  
 if config is None: config = {}  
 config.setdefault('learning\_rate', 1e-2)  
  
 w -= config['learning\_rate'] \* dw  
 return w, config

def sgd\_momentum(w, dw, config=None):  
 """  
 Performs stochastic gradient descent with momentum.  
  
 config format:  
 - learning\_rate: Scalar learning rate.  
 - momentum: Scalar between 0 and 1 giving the momentum value.  
 Setting momentum = 0 reduces to sgd.  
 - velocity: A numpy array of the same shape as w and dw used to store a moving  
 average of the gradients.  
 """  
 if config is None: config = {}  
 config.setdefault('learning\_rate', 1e-2)  
 config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there  
 v = config.get('velocity', np.zeros\_like(w)) # gets velocity, else sets it to zero.  
  
 # ===== #  
 # YOUR CODE HERE:  
 # Implement the momentum update formula. Return the updated weights  
 # as next\_w, and the updated velocity as v.  
 # ===== #  
  
 v = config['momentum'] \* v - config['learning\_rate'] \* dw  
 next\_w = w + v  
  
 # ===== #  
 # END YOUR CODE HERE  
 # ===== #  
  
 config['velocity'] = v  
  
 return next\_w, config

def sgd\_nesterov\_momentum(w, dw, config=None):  
 """  
 Performs stochastic gradient descent with Nesterov momentum.  
  
 config format:  
 - learning\_rate: Scalar learning rate.  
 - momentum: Scalar between 0 and 1 giving the momentum value.  
 Setting momentum = 0 reduces to sgd.  
 - velocity: A numpy array of the same shape as w and dw used to store a moving  
 average of the gradients.  
 """  
 if config is None: config = {}  
 config.setdefault('learning\_rate', 1e-2)  
 config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there  
 v = config.get('velocity', np.zeros\_like(w)) # gets velocity, else sets it to zero.  
  
 # ===== #  
 # YOUR CODE HERE:  
 # Implement the momentum update formula. Return the updated weights  
 # as next\_w, and the updated velocity as v.  
 # ===== #  
  
 # Use change of vars implementation  
 v\_old = v  
 v = config['momentum'] \* v - config['learning\_rate'] \* dw  
 next\_w = w + v + config['momentum'] \* (v - v\_old)  
  
 # ===== #  
 # END YOUR CODE HERE  
 # ===== #  
  
 config['velocity'] = v  
  
 return next\_w, config

def rmsprop(w, dw, config=None):  
 """  
 Uses the RMSProp update rule, which uses a moving average of squared gradient  
 values to set adaptive per-parameter learning rates.  
  
 config format:  
 - learning\_rate: Scalar learning rate.  
 - decay\_rate: Scalar between 0 and 1 giving the decay rate for the squared  
 gradient cache.  
 - epsilon: Small scalar used for smoothing to avoid dividing by zero.  
 - beta: Moving average of second moments of gradients.  
 """  
 if config is None: config = {}  
 config.setdefault('learning\_rate', 1e-2)  
 config.setdefault('decay\_rate', 0.99)  
 config.setdefault('epsilon', 1e-8)  
 config.setdefault('a', np.zeros\_like(w))  
  
 next\_w = None  
  
 # ===== #  
 # YOUR CODE HERE:  
 # Implement RMSProp. Store the next value of w as next\_w. You need  
 # to also store in config['a'] the moving average of the second  
 # moment gradients, so they can be used for future gradients. Concretely,  
 # config['a'] corresponds to "a" in the lecture notes.  
 # ===== #  
  
 config['a'] = config['a'] \* config['decay\_rate'] + (1-config['decay\_rate']) \* dw\*\*2  
 next\_w = w - config['learning\_rate']/(np.sqrt(config['a'] + config['epsilon']) \* dw  
  
 # ===== #  
 # END YOUR CODE HERE  
 # ===== #  
  
 return next\_w, config

def adam(w, dw, config=None):  
 """  
 Uses the Adam update rule, which incorporates moving averages of both the  
 gradient and its square and a bias correction term.  
  
 config format:  
 - learning\_rate: Scalar learning rate.  
 - beta1: Decay rate for moving average of first moment of gradient.  
 - beta2: Decay rate for moving average of second moment of gradient.  
 - epsilon: Small scalar used for smoothing to avoid dividing by zero.  
 - m: Moving average of gradient.  
 - v: Moving average of squared gradient.  
 - t: Iteration number.  
 """  
 if config is None: config = {}  
 config.setdefault('learning\_rate', 1e-3)  
 config.setdefault('beta1', 0.9)  
 config.setdefault('beta2', 0.999)  
 config.setdefault('epsilon', 1e-8)  
 config.setdefault('v', np.zeros\_like(w))  
 config.setdefault('a', np.zeros\_like(w))  
 config.setdefault('t', 0)  
  
 next\_w = None  
  
 # ===== #  
 # YOUR CODE HERE:  
 # Implement Adam. Store the next value of w as next\_w. You need  
 # to also store in config['a'] the moving average of the second  
 # moment gradients, and in config['v'] the moving average of the  
 # first moments. Finally, store in config['t'] the increasing time.  
 # ===== #  
  
 config['t'] += 1  
  
 # Moment Updates  
 config['v'] = config['beta1'] \* config['v'] + (1-config['beta1']) \* dw  
 config['a'] = config['beta2'] \* config['a'] + (1-config['beta2']) \* dw\*\*2  
  
 # Bias Corection  
 v\_corr = config['v'] / (1-config['beta1]\*\*config['t'])  
 a\_corr = config['a'] / (1-config['beta2']\*\*config['t'])  
  
 # Param Update  
 next\_w = w - config['learning\_rate']/(np.sqrt(a\_corr) + config['epsilon']) \* v\_corr  
 # ===== #  
 # END YOUR CODE HERE  
 # ===== #  
  
 return next\_w, config



```
import numpy as np
```

```
def affine_forward(x, w, b):
```

```
    """
    Computes the forward pass for an affine (fully-connected) layer.
```

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

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

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

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

```
    out = x.reshape(x.shape[0], np.prod(x.shape[1:])) @ w + b
```

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

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

```
def affine_backward(dout, cache):
```

```
    """
    Computes the backward pass for an affine layer.
```

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

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

```
    """
    x, w, b = cache
```

```
    dx, dw, db = None, None, None
```

```
    # ===== #
    # YOUR CODE HERE:
    #   Calculate the gradients for the backward pass.
    # Notice:
    #   dout is N x M
    #   dx should be N x d_1 x ... x d_k; it relates to dout through multiplication with w, which is D x M
```

```
    #   dw should be D x M; it relates to dout through multiplication with x, which is N x D after reshaping
    #   db should be M; it is just the sum over dout examples
    # ===== #
```

```
    dx = (dout @ w.T).reshape(x.shape)
    dw = x.reshape(x.shape[0], np.prod(x.shape[1:])).T @ dout
    db = np.sum(dout, axis=0)
```

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

```
    return dx, dw, db
```

```
def relu_forward(x):
```

```
    """
    Computes the forward pass for a layer of rectified linear units (ReLUs).
```

```
    Input:
    - x: Inputs, of any shape
```

```
    Returns a tuple of:
    - out: Output, of the same shape as x
    - cache: x
    """
```

```
    # ===== #
    # YOUR CODE HERE:
    #   Implement the ReLU forward pass.
    # ===== #
```

```
    out = np.maximum(0,x)
```

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

```
    cache = x
    return out, cache
```

```
def relu_backward(dout, cache):
```

```
    """
    Computes the backward pass for a layer of rectified linear units (ReLUs).
```

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

```
    Returns:
    - dx: Gradient with respect to x
    """
```

```
    x = cache
```

```
    # ===== #
    # YOUR CODE HERE:
    #   Implement the ReLU backward pass
    # ===== #
```

```
    dx = (x > 0) * dout
```

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

```
    return dx
```

```
def batchnorm_forward(x, gamma, beta, bn_param):
```

```
    """
    Forward pass for batch normalization.
```

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

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

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

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

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

```
    Returns a tuple of:
    - out: of shape (N, D)
    - cache: A tuple of values needed in the backward pass
    """
```

```
    mode = bn_param['mode']
    eps = bn_param.get('eps', 1e-5)
    momentum = bn_param.get('momentum', 0.9)
```

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

```
    out, cache = None, None
```

```
    if mode == 'train':
```

```
        # ===== #
        # YOUR CODE HERE:
        #   A few steps here:
        #   (1) Calculate the running mean and variance of the minibatch.
        #   (2) Normalize the activations with the running mean and variance.
        #   (3) Scale and shift the normalized activations. Store this
        #       as the variable 'out'
        #   (4) Store any variables you may need for the backward pass in
        #       the 'cache' variable.
        # ===== #

        batch_mean = x.mean(axis=0)
        batch_var = x.var(axis=0)
        running_mean = momentum * running_mean + (1 - momentum) * batch_mean
        running_var = momentum * running_var + (1 - momentum) * batch_var

        x_norm = (x - batch_mean) / np.sqrt(batch_var + eps)
        out = gamma * x_norm + beta

        cache = {
            'x_norm': x_norm,
            'gamma': gamma,
            'batch_var': batch_var,
            'eps': eps,
            'a': (x - batch_mean),
            'mean': batch_mean,
        }
```

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

```
    elif mode == 'test':
```

```
        # ===== #
        # YOUR CODE HERE:
        #   Calculate the testing time normalized activation. Normalize using
        #   the running mean and variance, and then scale and shift appropriately.
        #   Store the output as 'out'.
        # ===== #

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

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

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

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

```
    return out, cache
```

```
def batchnorm_backward(dout, cache):
```

```
    """
    Backward pass for batch normalization.
```

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

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

```
    Returns a tuple of:
    - dx: Gradient with respect to inputs x, of shape (N, D)
    - dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
    - dbeta: Gradient with respect to shift parameter beta, of shape (D,)
    """
```

```
    dx, dgamma, dbeta = None, None, None
```

```
    # ===== #
    # YOUR CODE HERE:
    #   Implement the batchnorm backward pass, calculating dx, dgamma, and dbeta.
    # ===== #
```

```
    M = dout.shape[0]
```

```
    dbeta = dout.sum(axis=0)
```

```
    dgamma = (cache['x_norm'] * dout).sum(axis=0)
```

```
    # dx_norm
    dx_norm = dout * cache['gamma']
    da = dx_norm / np.sqrt(cache['batch_var'] + cache['eps'])
    dmu = -1 * dx_norm.sum(axis=0) / np.sqrt(cache['batch_var'] + cache['eps'])
    db = cache['a'] * dx_norm
    dc = -1 * db / (cache['batch_var'] + cache['eps'])
    de = dc * 1/2 * 1/np.sqrt(cache['batch_var'] + cache['eps'])
    dvar = de.sum(axis=0)
```

```
    dx = da + 2/M * cache['a'] * dvar + 1/M * dmu
```

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

```
    return dx, dgamma, dbeta
```

```
def dropout_forward(x, dropout_param):
```

```
    """
    Performs the forward pass for (inverted) dropout.
```

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

```
    Outputs:
    - out: Array of the same shape as x.
    - cache: A tuple (dropout_param, mask). In training mode, mask is the dropout
      mask that was used to multiply the input; in test mode, mask is None.
    """
```

```
    p, mode = dropout_param['p'], dropout_param['mode']
```

```
    if 'seed' in dropout_param:
        np.random.seed(dropout_param['seed'])
```

```
    mask = None
```

```
    out = None
```

```
    if mode == 'train':
```

```
        # ===== #
        # YOUR CODE HERE:
        #   Implement the inverted dropout forward pass during training time.
        #   Store the masked and scaled activations in out, and store the
        #   dropout mask as the variable mask.
        # ===== #

        mask = (np.random.random_sample(x.shape) >= p) / (1 - p)
        out = x * mask

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

```
    elif mode == 'test':
```

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

```
        out = x
```

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

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

```
    return out, cache
```

```
def dropout_backward(dout, cache):
```

```
    """
    Perform the backward pass for (inverted) dropout.
```

```
    Inputs:
    - dout: Upstream derivatives, of any shape
    - cache: (dropout_param, mask) from dropout_forward.
```

```
    """
    dropout_param, mask = cache
    mode = dropout_param['mode']
```

```
    dx = None
```

```
    if mode == 'train':
```

```
        # ===== #
        # YOUR CODE HERE:
        #   Implement the inverted dropout backward pass during training time.
        # ===== #
```

```
        dx = dout * mask
```

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

```
    elif mode == 'test':
```

```
        # ===== #
        # YOUR CODE HERE:
        #   Implement the inverted dropout backward pass during test time.
        # ===== #
```

```
        dx = dout
```

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

```
    return dx
```

```
def svm_loss(x, y):
```

```
    """
    Computes the loss and gradient using for multiclass SVM classification.
```

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

```
    Returns a tuple of:
    - loss: Scalar giving the loss
    - dx: Gradient of the loss with respect to x
    """
```

```
    N = x.shape[0]
    correct_class_scores = x[np.arange(N), y]
    margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
    margins[np.arange(N), y] = 0
    loss = np.sum(margins) / N
```

```
    num_pos = np.sum(margins > 0, axis=1)
```

```
    dx = np.zeros_like(x)
```

```
    dx[margins > 0] = 1
```

```
    dx[np.arange(N), y] -= num_pos
```

```
    dx /= N
```

```
    return loss, dx
```

```
def softmax_loss(x, y):
```

```
    """
    Computes the loss and gradient for softmax classification.
```

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

```
    Returns a tuple of:
    - loss: Scalar giving the loss
    - dx: Gradient of the loss with respect to x
    """
```

```
    probs = np.exp(x - np.max(x, axis=1, keepdims=True))
    probs /= np.sum(probs, axis=1, keepdims=True)
```

```
    N = x.shape[0]
```

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

```
    dx = probs.copy()
```

```
    dx[np.arange(N), y] -= 1
```

```
    dx /= N
```

```
    return loss, dx
```



```
import numpy as np
from .layers import *
from .layer_utils import *

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

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

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

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

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

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

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

        # Concat dims for full NN
        dims = [input_dim] + hidden_dims + [num_classes]
        for layer in range(self.num_layers):
            self.params['W' + str(layer + 1)] = np.random.normal(0, weight_scale, (dims[layer], dims[layer + 1]))
            self.params['b' + str(layer + 1)] = np.zeros(dims[layer + 1])

            if self.use_batchnorm and (layer != (self.num_layers-1)):
                self.params['gamma' + str(layer + 1)] = np.ones(dims[layer + 1])
                self.params['beta' + str(layer + 1)] = np.zeros(dims[layer + 1])

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

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

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

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

def loss(self, X, y=None):
    """
    Compute loss and gradient for the fully-connected net.

    Input / output: Same as TwoLayerNet above.
    """
    X = X.astype(self.dtype)
    mode = 'test' if y is None else 'train'

    # Set train/test mode for batchnorm params and dropout param since they
    # behave differently during training and testing.
    if self.dropout_param is not None:
        self.dropout_param['mode'] = mode
    if self.use_batchnorm:
        for bn_param in self.bn_params:
            bn_param[mode] = mode

    scores = None

    # ===== #
    # YOUR CODE HERE:
    #   Implement the forward pass of the FC net and store the output
    #   scores as the variable "scores".
    #
    #   BATCHNORM: If self.use_batchnorm is true, insert a bathnorm layer
    #   between the affine_forward and relu_forward layers. You may
    #   also write an affine_batchnorm_relu() function in layer_utils.py.
    #
    #   DROPOUT: If dropout is non-zero, insert a dropout layer after
    #   every ReLU layer.
    # ===== #

    a = {}
    norm = {}
    h = {}
    drop = {}
    drop[0] = [X]

    for layer in range(self.num_layers):
        #Affine
        a[layer + 1] = affine_forward(drop[layer][0], self.params['W' + str(layer + 1)], self.params['b' + str(layer + 1)])

        if layer < (self.num_layers-1):
            # BatchNorm
            if self.use_batchnorm: norm[layer + 1] = batchnorm_forward(a[layer + 1][0], self.params['gamma' + str(layer + 1)],
                self.params['beta' + str(layer + 1)], self.bn_params[layer])
            else: norm[layer + 1] = a[layer + 1]
            # ReLU
            h[layer + 1] = relu_forward(norm[layer + 1][0])
            # Dropout
            if self.use_dropout: drop[layer + 1] = dropout_forward(h[layer + 1][0], self.dropout_param)
            else: drop[layer + 1] = h[layer + 1]

    scores = a[self.num_layers][0]

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

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

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

    loss, dout = softmax_loss(scores, y)
    Ws = [self.params['W' + str(i + 1)] for i in range(self.num_layers)]

    loss += 0.5 * self.reg * sum([np.linalg.norm(weight, 'fro')**2 for weight in Ws])
    das = {}
    dhs = {}
    ddrops = {}
    dnorms = {}
    dgammas = {}
    dbetas = {}
    dws = {}
    dbs = {}
    das[self.num_layers] = dout

    for layer in reversed(range(self.num_layers)):
        ddrops[layer], dws[layer + 1], dbs[layer + 1] = affine_backward(das[layer + 1], a[layer + 1][1])
        if layer != 0:
            if self.use_dropout: dhs[layer] = dropout_backward(ddrops[layer], drop[layer][1])
            else: dhs[layer] = ddrops[layer]
            dnorms[layer] = relu_backward(dhs[layer], h[layer][1])
            if self.use_batchnorm: das[layer], dgammas[layer], dbetas[layer] = batchnorm_backward(dnorms[layer], norm[layer][1])
            else: das[layer] = dnorms[layer]

    for layer in range(self.num_layers):
        grads['W' + str(layer + 1)] = dws[layer + 1] + self.reg * self.params['W' + str(layer + 1)]
        grads['b' + str(layer + 1)] = dbs[layer + 1].T
        if layer != (self.num_layers-1) and self.use_batchnorm:
            grads['gamma' + str(layer + 1)] = dgammas[layer + 1]
            grads['beta' + str(layer + 1)] = dbetas[layer + 1].T

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

    return loss, grads
```