

# Batch-Normalization

February 12, 2020

## 1 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.

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

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

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

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

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

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

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

```
data = get_CIFAR10_data()
```

```

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

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

```

## 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.

```

In [3]: # Check the training-time forward pass by checking means and variances
        # of features both before and after batch normalization

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

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

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

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

```

Before batch normalization:

```

means: [-15.68145296 -40.58541258  18.50642948]
stds:  [30.33999891 35.02984249 36.45512118]

```

After batch normalization (gamma=1, beta=0)

```

mean:  [ 2.89213098e-16  4.43881043e-16 -8.93729535e-17]

```

```

std: [0.99999999 1.          1.          ]
After batch normalization (nontrivial gamma, beta)
means: [11. 12. 13.]
stds: [0.99999999 1.99999999 2.99999999]

```

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

```

In [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: [ 0.03117551 -0.01980777  0.07824601]
stds:  [1.02294919  0.97886324  1.12352346]

```

## 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.

```

In [5]: # Gradient check batchnorm backward pass

N, D = 4, 5

```

```

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

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

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

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

dx error:  8.401738438975145e-10
dgamma error:  1.3160578264246991e-11
dbeta error:  3.2755691451867106e-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  $1e-4$ .

```

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 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: 2.305524453293663
W1 relative error: 0.0004251511174331715
W2 relative error: 0.0001600432781126721
W3 relative error: 4.0176784209594947e-10
b1 relative error: 5.551115123125783e-09
b2 relative error: 2.220446049250313e-08
b3 relative error: 1.0232207151980328e-10
beta1 relative error: 2.0934948479641335e-08
beta2 relative error: 1.1850738852962536e-07
gamma1 relative error: 2.692208373357594e-08
gamma2 relative error: 4.423932093186216e-09

```

```

Running check with reg = 3.14
Initial loss: 7.2375671353161914
W1 relative error: 2.7470559689743688e-05
W2 relative error: 1.3011673445356122e-06
W3 relative error: 4.981845768150892e-06
b1 relative error: 1.7763568394002505e-07
b2 relative error: 2.220446049250313e-08
b3 relative error: 1.6227850630042402e-10
beta1 relative error: 3.729162936195644e-09
beta2 relative error: 3.156742490997375e-08
gamma1 relative error: 3.7534402710616995e-09
gamma2 relative error: 2.0124868359752387e-08

```

## 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.

```

In [7]: # 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: 2.302887
(Epoch 0 / 10) train acc: 0.137000; val_acc: 0.135000
(Epoch 1 / 10) train acc: 0.365000; val_acc: 0.277000
(Epoch 2 / 10) train acc: 0.441000; val_acc: 0.334000
(Epoch 3 / 10) train acc: 0.499000; val_acc: 0.329000
(Epoch 4 / 10) train acc: 0.553000; val_acc: 0.319000
(Epoch 5 / 10) train acc: 0.617000; val_acc: 0.314000
(Epoch 6 / 10) train acc: 0.660000; val_acc: 0.337000
(Epoch 7 / 10) train acc: 0.714000; val_acc: 0.345000
(Epoch 8 / 10) train acc: 0.777000; val_acc: 0.339000
(Epoch 9 / 10) train acc: 0.768000; val_acc: 0.350000
(Epoch 10 / 10) train acc: 0.804000; val_acc: 0.338000
(Iteration 1 / 200) loss: 2.301780
(Epoch 0 / 10) train acc: 0.161000; val_acc: 0.143000
(Epoch 1 / 10) train acc: 0.226000; val_acc: 0.196000
(Epoch 2 / 10) train acc: 0.303000; val_acc: 0.266000
(Epoch 3 / 10) train acc: 0.327000; val_acc: 0.264000

```

```
(Epoch 4 / 10) train acc: 0.368000; val_acc: 0.295000
(Epoch 5 / 10) train acc: 0.422000; val_acc: 0.311000
(Epoch 6 / 10) train acc: 0.466000; val_acc: 0.308000
(Epoch 7 / 10) train acc: 0.454000; val_acc: 0.295000
(Epoch 8 / 10) train acc: 0.542000; val_acc: 0.315000
(Epoch 9 / 10) train acc: 0.559000; val_acc: 0.320000
(Epoch 10 / 10) train acc: 0.607000; val_acc: 0.318000
```

```
In [8]: fig, axes = plt.subplots(3, 1)

ax = axes[0]
ax.set_title('Training loss')
ax.set_xlabel('Iteration')

ax = axes[1]
ax.set_title('Training accuracy')
ax.set_xlabel('Epoch')

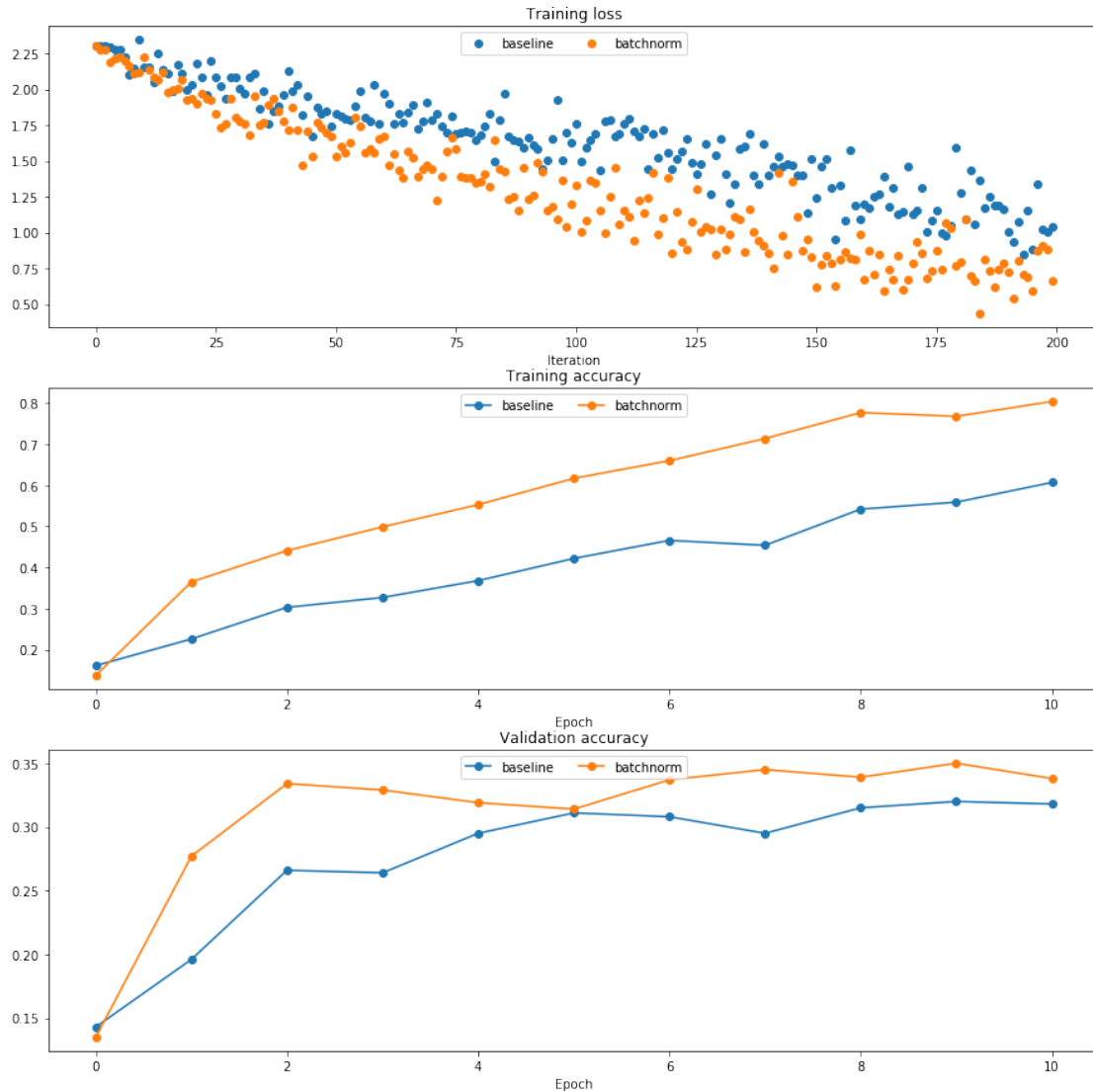
ax = axes[2]
ax.set_title('Validation accuracy')
ax.set_xlabel('Epoch')

ax = axes[0]
ax.plot(solver.loss_history, 'o', label='baseline')
ax.plot(bn_solver.loss_history, 'o', label='batchnorm')

ax = axes[1]
ax.plot(solver.train_acc_history, '-o', label='baseline')
ax.plot(bn_solver.train_acc_history, '-o', label='batchnorm')

ax = axes[2]
ax.plot(solver.val_acc_history, '-o', label='baseline')
ax.plot(bn_solver.val_acc_history, '-o', label='batchnorm')

for i in [1, 2, 3]:
    ax = axes[i - 1]
    ax.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.

```
In [9]: # 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

```

```
/Users/edwardzhang/Desktop/ece247/HW4/HW4-code/nndl/layers.py:426: RuntimeWarning: divide by z
    loss = -np.sum(np.log(probs[np.arange(N), y])) / N
```

```
Running weight scale 17 / 20
Running weight scale 18 / 20
Running weight scale 19 / 20
Running weight scale 20 / 20
```

```
In [10]: # 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))

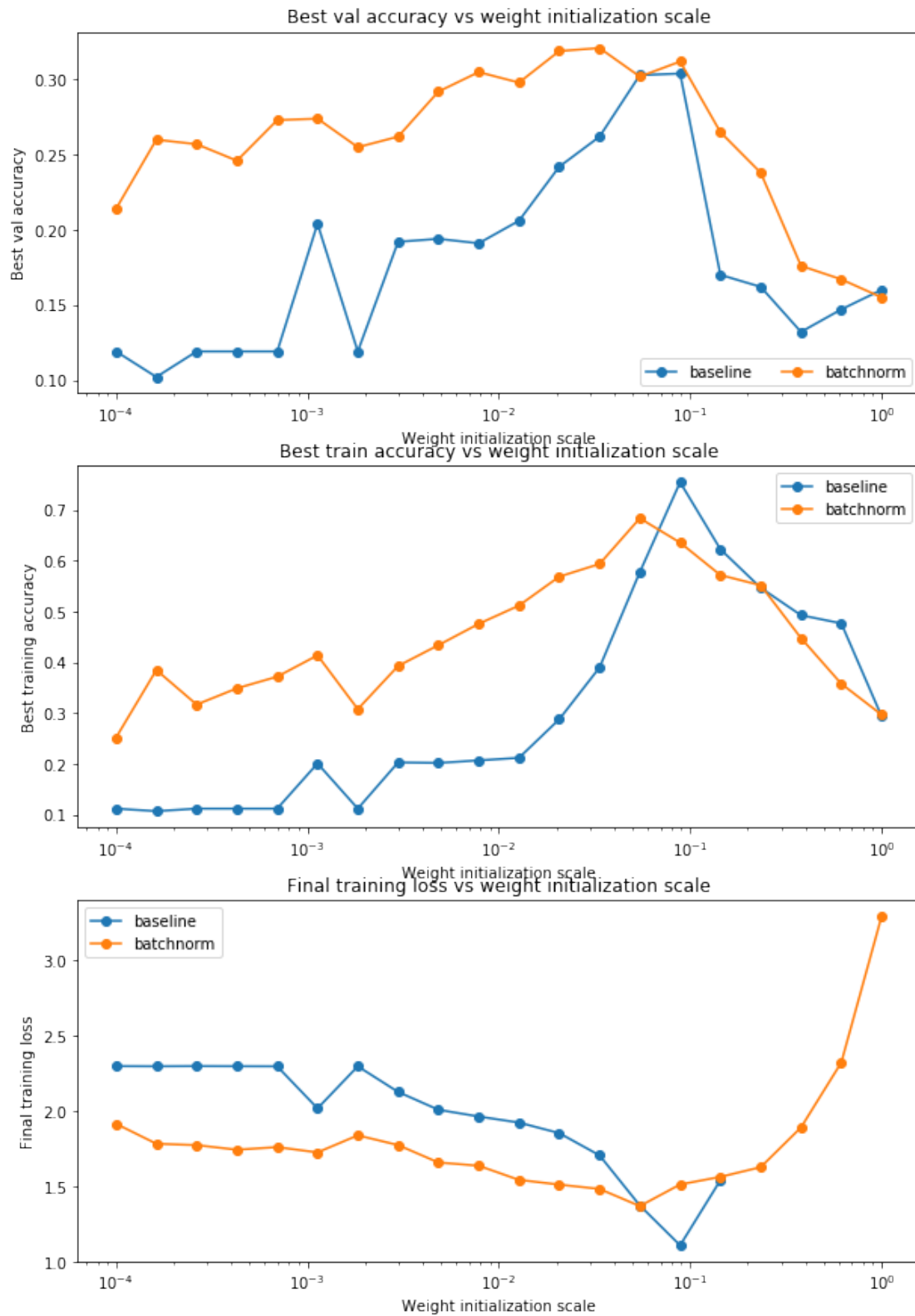
    final_train_loss.append(np.mean(solvers[ws].loss_history[-100:]))
    bn_final_train_loss.append(np.mean(bn_solvers[ws].loss_history[-100:]))

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:**

This makes sense because with batchnorm we have greater tolerance for poor weight initializations.

# Dropout

February 12, 2020

## 1 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.

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

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

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

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

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

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

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

```

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

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

```

## 1.1 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 [3]: x = np.random.randn(500, 500) + 10

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

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

```

```

Running tests with p = 0.3
Mean of input: 9.997243690342161
Mean of train-time output: 10.051213555038881
Mean of test-time output: 9.997243690342161
Fraction of train-time output set to zero: 0.698476
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.6
Mean of input: 9.997243690342161
Mean of train-time output: 9.997800418539825
Mean of test-time output: 9.997243690342161
Fraction of train-time output set to zero: 0.400136
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.75
Mean of input: 9.997243690342161
Mean of train-time output: 9.990889323112652
Mean of test-time output: 9.997243690342161
Fraction of train-time output set to zero: 0.250436
Fraction of test-time output set to zero: 0.0

```

## 1.2 Dropout backward pass

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

```
In [4]: x = np.random.randn(10, 10) + 10
        dout = np.random.randn(*x.shape)

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

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

dx relative error: 5.4456117023484937e-11
```

## 1.3 Implement a fully connected neural network with dropout layers

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

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

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

```
In [5]: N, D, H1, H2, C = 2, 15, 20, 30, 10
        X = np.random.randn(N, D)
        y = np.random.randint(C, size=(N,))

        for dropout in [0.5, 0.75, 1.0]:
            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 = 0.5
Initial loss: 2.309771209610118
W1 relative error: 2.694274363733021e-07
W2 relative error: 7.439246147919978e-08
W3 relative error: 1.910371122296728e-08
b1 relative error: 4.112891126518e-09
b2 relative error: 5.756217724722137e-10
b3 relative error: 1.3204470857080166e-10
```

```
Running check with dropout = 0.75
Initial loss: 2.306133548427975
W1 relative error: 8.72986097970181e-08
W2 relative error: 2.9777307885797295e-07
W3 relative error: 1.8832780806174298e-08
b1 relative error: 5.379486003985169e-08
b2 relative error: 3.6529949080385546e-09
b3 relative error: 9.987242764516995e-11
```

```
Running check with dropout = 1.0
Initial loss: 2.3053332250963194
W1 relative error: 1.2744095365229032e-06
W2 relative error: 4.678743300473988e-07
W3 relative error: 4.331673892536035e-08
b1 relative error: 4.0853539035931665e-08
b2 relative error: 1.951342257912746e-09
b3 relative error: 9.387142701440351e-11
```

## 1.4 Dropout as a regularizer

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

```
In [6]: # Train two identical nets, one with dropout and one without
```

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

solvers = {}
```



```

dropout_choices = [0.6, 1.0]
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: 2.300199
(Epoch 0 / 25) train acc: 0.158000; val_acc: 0.127000
(Epoch 1 / 25) train acc: 0.132000; val_acc: 0.121000
(Epoch 2 / 25) train acc: 0.204000; val_acc: 0.170000
(Epoch 3 / 25) train acc: 0.240000; val_acc: 0.192000
(Epoch 4 / 25) train acc: 0.312000; val_acc: 0.274000
(Epoch 5 / 25) train acc: 0.314000; val_acc: 0.269000
(Epoch 6 / 25) train acc: 0.364000; val_acc: 0.252000
(Epoch 7 / 25) train acc: 0.390000; val_acc: 0.281000
(Epoch 8 / 25) train acc: 0.386000; val_acc: 0.290000
(Epoch 9 / 25) train acc: 0.372000; val_acc: 0.267000
(Epoch 10 / 25) train acc: 0.424000; val_acc: 0.286000
(Epoch 11 / 25) train acc: 0.396000; val_acc: 0.275000
(Epoch 12 / 25) train acc: 0.458000; val_acc: 0.299000
(Epoch 13 / 25) train acc: 0.496000; val_acc: 0.305000
(Epoch 14 / 25) train acc: 0.492000; val_acc: 0.299000
(Epoch 15 / 25) train acc: 0.550000; val_acc: 0.296000
(Epoch 16 / 25) train acc: 0.584000; val_acc: 0.297000
(Epoch 17 / 25) train acc: 0.582000; val_acc: 0.309000
(Epoch 18 / 25) train acc: 0.612000; val_acc: 0.306000
(Epoch 19 / 25) train acc: 0.628000; val_acc: 0.323000
(Epoch 20 / 25) train acc: 0.608000; val_acc: 0.324000
(Iteration 101 / 125) loss: 1.369535
(Epoch 21 / 25) train acc: 0.644000; val_acc: 0.330000
(Epoch 22 / 25) train acc: 0.708000; val_acc: 0.341000
(Epoch 23 / 25) train acc: 0.690000; val_acc: 0.297000
(Epoch 24 / 25) train acc: 0.740000; val_acc: 0.300000
(Epoch 25 / 25) train acc: 0.750000; val_acc: 0.329000
(Iteration 1 / 125) loss: 2.300607
(Epoch 0 / 25) train acc: 0.172000; val_acc: 0.167000
(Epoch 1 / 25) train acc: 0.210000; val_acc: 0.197000
(Epoch 2 / 25) train acc: 0.284000; val_acc: 0.240000
(Epoch 3 / 25) train acc: 0.302000; val_acc: 0.246000
(Epoch 4 / 25) train acc: 0.392000; val_acc: 0.289000

```

```

(Epoch 5 / 25) train acc: 0.420000; val_acc: 0.274000
(Epoch 6 / 25) train acc: 0.420000; val_acc: 0.304000
(Epoch 7 / 25) train acc: 0.474000; val_acc: 0.293000
(Epoch 8 / 25) train acc: 0.516000; val_acc: 0.330000
(Epoch 9 / 25) train acc: 0.566000; val_acc: 0.322000
(Epoch 10 / 25) train acc: 0.620000; val_acc: 0.321000
(Epoch 11 / 25) train acc: 0.656000; val_acc: 0.317000
(Epoch 12 / 25) train acc: 0.676000; val_acc: 0.319000
(Epoch 13 / 25) train acc: 0.680000; val_acc: 0.304000
(Epoch 14 / 25) train acc: 0.752000; val_acc: 0.323000
(Epoch 15 / 25) train acc: 0.802000; val_acc: 0.321000
(Epoch 16 / 25) train acc: 0.804000; val_acc: 0.300000
(Epoch 17 / 25) train acc: 0.868000; val_acc: 0.303000
(Epoch 18 / 25) train acc: 0.894000; val_acc: 0.298000
(Epoch 19 / 25) train acc: 0.910000; val_acc: 0.282000
(Epoch 20 / 25) train acc: 0.926000; val_acc: 0.316000
(Iteration 101 / 125) loss: 0.245816
(Epoch 21 / 25) train acc: 0.950000; val_acc: 0.282000
(Epoch 22 / 25) train acc: 0.958000; val_acc: 0.292000
(Epoch 23 / 25) train acc: 0.966000; val_acc: 0.307000
(Epoch 24 / 25) train acc: 0.966000; val_acc: 0.285000
(Epoch 25 / 25) train acc: 0.970000; val_acc: 0.284000

```

In [7]: *# Plot train and validation accuracies of the two models*

```

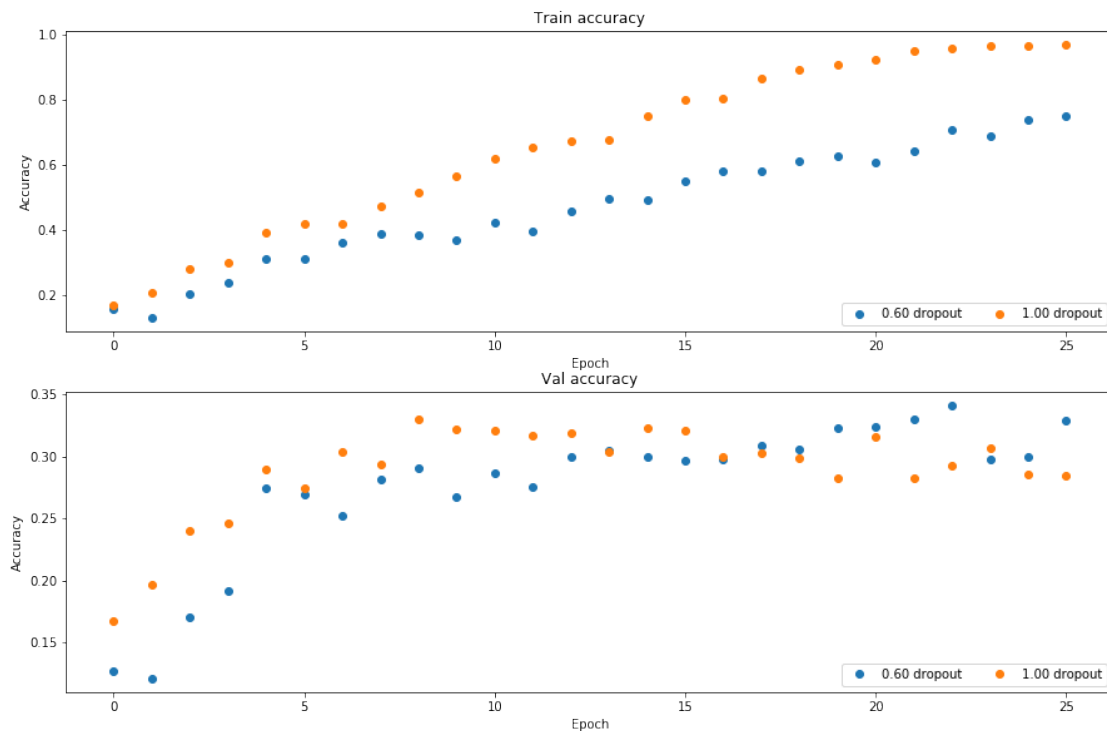
train_accs = []
val_accs = []
for dropout in dropout_choices:
    solver = solvers[dropout]
    train_accs.append(solver.train_acc_history[-1])
    val_accs.append(solver.val_acc_history[-1])

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

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

```

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



## 1.5 Question

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

## 1.6 Answer:

Yes, dropout is performing regularization because with dropout, our model has a training accuracy much closer to validation accuracy whereas our model without dropout has a much higher training accuracy than validation accuracy.

**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:

$\min(\text{floor}((X - 32\%)) / 23\%, 1)$  where if you get 55% or higher validation accuracy, you get full points.

```
In [10]: # ===== #
# YOUR CODE HERE:
# Implement a FC-net that achieves at least 55% validation accuracy
# on CIFAR-10.
```

```

# ===== #

optimizer = 'adam'
best_model = None

layer_dims = [500, 500, 500]
weight_scale = 0.01
learning_rate = 5e-3
lr_decay = 0.9
use_batchnorm = True
dropout = 0.5

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

solver = Solver(model, data,
                 num_epochs=10, batch_size=256,
                 update_rule=optimizer,
                 optim_config={
                     'learning_rate': learning_rate,
                 },
                 lr_decay=lr_decay,
                 verbose=True, print_every=50)

solver.train()

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

```

```

(Iteration 1 / 1910) loss: 2.313253
(Epoch 0 / 10) train acc: 0.156000; val_acc: 0.140000
(Iteration 51 / 1910) loss: 1.900724
(Iteration 101 / 1910) loss: 1.752910
(Iteration 151 / 1910) loss: 1.790791
(Epoch 1 / 10) train acc: 0.434000; val_acc: 0.448000
(Iteration 201 / 1910) loss: 1.627303
(Iteration 251 / 1910) loss: 1.690706
(Iteration 301 / 1910) loss: 1.710692
(Iteration 351 / 1910) loss: 1.608809
(Epoch 2 / 10) train acc: 0.492000; val_acc: 0.486000
(Iteration 401 / 1910) loss: 1.618772
(Iteration 451 / 1910) loss: 1.501616
(Iteration 501 / 1910) loss: 1.434100
(Iteration 551 / 1910) loss: 1.591375
(Epoch 3 / 10) train acc: 0.533000; val_acc: 0.509000
(Iteration 601 / 1910) loss: 1.412635
(Iteration 651 / 1910) loss: 1.277013
(Iteration 701 / 1910) loss: 1.451348
(Iteration 751 / 1910) loss: 1.477529

```

(Epoch 4 / 10) train acc: 0.544000; val\_acc: 0.516000  
(Iteration 801 / 1910) loss: 1.440390  
(Iteration 851 / 1910) loss: 1.421336  
(Iteration 901 / 1910) loss: 1.403937  
(Iteration 951 / 1910) loss: 1.349680  
(Epoch 5 / 10) train acc: 0.536000; val\_acc: 0.549000  
(Iteration 1001 / 1910) loss: 1.371560  
(Iteration 1051 / 1910) loss: 1.296324  
(Iteration 1101 / 1910) loss: 1.444857  
(Epoch 6 / 10) train acc: 0.566000; val\_acc: 0.525000  
(Iteration 1151 / 1910) loss: 1.390857  
(Iteration 1201 / 1910) loss: 1.342500  
(Iteration 1251 / 1910) loss: 1.409043  
(Iteration 1301 / 1910) loss: 1.310190  
(Epoch 7 / 10) train acc: 0.590000; val\_acc: 0.529000  
(Iteration 1351 / 1910) loss: 1.275100  
(Iteration 1401 / 1910) loss: 1.400575  
(Iteration 1451 / 1910) loss: 1.372440  
(Iteration 1501 / 1910) loss: 1.381053  
(Epoch 8 / 10) train acc: 0.603000; val\_acc: 0.546000  
(Iteration 1551 / 1910) loss: 1.342825  
(Iteration 1601 / 1910) loss: 1.404309  
(Iteration 1651 / 1910) loss: 1.274069  
(Iteration 1701 / 1910) loss: 1.350118  
(Epoch 9 / 10) train acc: 0.619000; val\_acc: 0.540000  
(Iteration 1751 / 1910) loss: 1.212386  
(Iteration 1801 / 1910) loss: 1.276342  
(Iteration 1851 / 1910) loss: 1.240554  
(Iteration 1901 / 1910) loss: 1.283522  
(Epoch 10 / 10) train acc: 0.618000; val\_acc: 0.560000

```

1 import numpy as np
2 import pdb
3 import copy
4
5 from .layers import *
6 from .layer_utils import *
7
8 """
9 This code was originally written for CS 231n at Stanford University
10 (cs231n.stanford.edu). It has been modified in various areas for use in the
11 ECE 239AS class at UCLA. This includes the descriptions of what code to
12 implement as well as some slight potential changes in variable names to be
13 consistent with class nomenclature. We thank Justin Johnson & Serena Yeung
14 for
15 permission to use this code. To see the original version, please visit
16 cs231n.stanford.edu.
17 """
18 class TwoLayerNet(object):
19     """
20     A two-layer fully-connected neural network with ReLU nonlinearity and
21     softmax loss that uses a modular layer design. We assume an input
22     dimension
23     of D, a hidden dimension of H, and perform classification over C classes.
24
25     The architecture should be affine - relu - affine - softmax.
26
27     Note that this class does not implement gradient descent; instead, it
28     will interact with a separate Solver object that is responsible for
29     running
30     optimization.
31
32     The learnable parameters of the model are stored in the dictionary
33     self.params that maps parameter names to numpy arrays.
34     """
35     def __init__(self, input_dim=3*32*32, hidden_dims=100, num_classes=10,
36                 dropout=0, weight_scale=1e-3, reg=0.0):
37         """
38         Initialize a new network.
39
40         Inputs:
41         - input_dim: An integer giving the size of the input
42         - hidden_dims: An integer giving the size of the hidden layer
43         - num_classes: An integer giving the number of classes to classify
44         - dropout: Scalar between 0 and 1 giving dropout strength.
45         - weight_scale: Scalar giving the standard deviation for random
46           initialization of the weights.
47         - reg: Scalar giving L2 regularization strength.
48         """
49         self.params = {}
50         self.reg = reg
51
52         # ===== #
53         # YOUR CODE HERE:
54         # Initialize W1, W2, b1, and b2. Store these as self.params['W1'],
55         # self.params['W2'], self.params['b1'] and self.params['b2']. The
56         # biases are initialized to zero and the weights are initialized
57         # so that each parameter has mean 0 and standard deviation
58         weight_scale.

```

```

57     # The dimensions of W1 should be (input_dim, hidden_dim) and the
58     # dimensions of W2 should be (hidden_dims, num_classes)
59     # ===== #
60
61
62     # ===== #
63     # END YOUR CODE HERE
64     # ===== #
65
66     def loss(self, X, y=None):
67         """
68         Compute loss and gradient for a minibatch of data.
69
70         Inputs:
71         - X: Array of input data of shape (N, d_1, ..., d_k)
72         - y: Array of labels, of shape (N,). y[i] gives the label for X[i].
73
74         Returns:
75         If y is None, then run a test-time forward pass of the model and
return:
76         - scores: Array of shape (N, C) giving classification scores, where
77           scores[i, c] is the classification score for X[i] and class c.
78
79         If y is not None, then run a training-time forward and backward pass
and
80         return a tuple of:
81         - loss: Scalar value giving the loss
82         - grads: Dictionary with the same keys as self.params, mapping
parameter
83           names to gradients of the loss with respect to those parameters.
84         """
85         scores = None
86
87         # ===== #
88         # YOUR CODE HERE:
89         # Implement the forward pass of the two-layer neural network. Store
90         # the class scores as the variable 'scores'. Be sure to use the
layers
91         # you prior implemented.
92         # ===== #
93
94         # ===== #
95         # END YOUR CODE HERE
96         # ===== #
97
98         # If y is None then we are in test mode so just return scores
99         if y is None:
100             return scores
101
102         loss, grads = 0, {}
103         # ===== #
104         # YOUR CODE HERE:
105         # Implement the backward pass of the two-layer neural net. Store
106         # the loss as the variable 'loss' and store the gradients in the
107         # 'grads' dictionary. For the grads dictionary, grads['W1'] holds
108         # the gradient for W1, grads['b1'] holds the gradient for b1, etc.
109         # i.e., grads[k] holds the gradient for self.params[k].
110         #

```

```

111         # Add L2 regularization, where there is an added cost
112         0.5*self.reg*W^2
113         # for each W. Be sure to include the 0.5 multiplying factor to
114         # match our implementation.
115         # And be sure to use the layers you prior implemented.
116         # ===== #
117
118         # ===== #
119         # END YOUR CODE HERE
120         # ===== #
121
122         return loss, grads
123
124
125 class FullyConnectedNet(object):
126     """
127     A fully-connected neural network with an arbitrary number of hidden
128     layers,
129     ReLU nonlinearities, and a softmax loss function. This will also
130     implement
131     dropout and batch normalization as options. For a network with L layers,
132     the architecture will be
133     {affine - [batch norm] - relu - [dropout]} x (L - 1) - affine - softmax
134     where batch normalization and dropout are optional, and the {...} block
135     is
136     repeated L - 1 times.
137
138     Similar to the TwoLayerNet above, learnable parameters are stored in the
139     self.params dictionary and will be learned using the Solver class.
140     """
141     def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,
142                 dropout=0, use_batchnorm=False, reg=0.0,
143                 weight_scale=1e-2, dtype=np.float32, seed=None):
144         """
145         Initialize a new FullyConnectedNet.
146
147         Inputs:
148         - hidden_dims: A list of integers giving the size of each hidden
149         layer.
150         - input_dim: An integer giving the size of the input.
151         - num_classes: An integer giving the number of classes to classify.
152         - dropout: Scalar between 0 and 1 giving dropout strength. If
153         dropout=1 then
154         the network should not use dropout at all.
155         - use_batchnorm: Whether or not the network should use batch
156         normalization.
157         - reg: Scalar giving L2 regularization strength.
158         - weight_scale: Scalar giving the standard deviation for random
159         initialization of the weights.
160         - dtype: A numpy datatype object; all computations will be performed
161         using
162         this datatype. float32 is faster but less accurate, so you should
163         use
164         float64 for numeric gradient checking.

```



```

160         - seed: If not None, then pass this random seed to the dropout
layers. This
161         will make the dropout layers deterministic so we can gradient check
the
162         model.
163         """
164         self.use_batchnorm = use_batchnorm
165         self.use_dropout = dropout < 1
166         self.reg = reg
167         self.num_layers = 1 + len(hidden_dims)
168         self.dtype = dtype
169         self.params = {}
170
171         # ===== #
172         # YOUR CODE HERE:
173         # Initialize all parameters of the network in the self.params
dictionary.
174         # The weights and biases of layer 1 are W1 and b1; and in general
the
175         # weights and biases of layer i are Wi and bi. The
176         # biases are initialized to zero and the weights are initialized
177         # so that each parameter has mean 0 and standard deviation
weight_scale.
178         #
179         # BATCHNORM: Initialize the gammas of each layer to 1 and the beta
180         # parameters to zero. The gamma and beta parameters for layer 1
should
181         # be self.params['gamma1'] and self.params['beta1']. For layer 2,
they
182         # should be gamma2 and beta2, etc. Only use batchnorm if
self.use_batchnorm
183         # is true and DO NOT do batch normalize the output scores.
184         # ===== #
185
186         dims = hidden_dims[:]
187         dims.insert(0, input_dim)
188         dims.append(num_classes)
189
190         for i in range(len(dims) - 1):
191             self.params['W{}'.format(i+1)] = weight_scale *
np.random.randn(dims[i], dims[i+1])
192             self.params['b{}'.format(i+1)] = np.zeros(dims[i+1])
193
194         if use_batchnorm:
195             for i in range(self.num_layers - 1):
196                 self.params['gamma{}'.format(i+1)] = np.ones(hidden_dims[i])
197                 self.params['beta{}'.format(i+1)] = np.zeros(hidden_dims[i])
198
199
200         # ===== #
201         # END YOUR CODE HERE
202         # ===== #
203
204         # When using dropout we need to pass a dropout_param dictionary to
each
205         # dropout layer so that the layer knows the dropout probability and
the mode
206         # (train / test). You can pass the same dropout_param to each dropout
layer.
207         self.dropout_param = {}

```

```

208         if self.use_dropout:
209             self.dropout_param = {'mode': 'train', 'p': dropout}
210         if seed is not None:
211             self.dropout_param['seed'] = seed
212
213         # With batch normalization we need to keep track of running means and
214         # variances, so we need to pass a special bn_param object to each
batch
215         # normalization layer. You should pass self.bn_params[0] to the
forward pass
216         # of the first batch normalization layer, self.bn_params[1] to the
forward
217         # pass of the second batch normalization layer, etc.
218         self.bn_params = []
219         if self.use_batchnorm:
220             self.bn_params = [{'mode': 'train'} for i in
np.arange(self.num_layers - 1)]
221
222         # Cast all parameters to the correct datatype
223         for k, v in self.params.items():
224             self.params[k] = v.astype(dtype)
225
226
227     def loss(self, X, y=None):
228         """
229         Compute loss and gradient for the fully-connected net.
230
231         Input / output: Same as TwoLayerNet above.
232         """
233         X = X.astype(self.dtype)
234         mode = 'test' if y is None else 'train'
235
236         # Set train/test mode for batchnorm params and dropout param since
they
237         # behave differently during training and testing.
238         if self.dropout_param is not None:
239             self.dropout_param['mode'] = mode
240         if self.use_batchnorm:
241             for bn_param in self.bn_params:
242                 bn_param[mode] = mode
243
244         scores = None
245
246         # ===== #
247         # YOUR CODE HERE:
248         # Implement the forward pass of the FC net and store the output
249         # scores as the variable "scores".
250         #
251         # BATCHNORM: If self.use_batchnorm is true, insert a bathnorm layer
252         # between the affine_forward and relu_forward layers. You may
253         # also write an affine_batchnorm_relu() function in layer_utils.py.
254         #
255         # DROPOUT: If dropout is non-zero, insert a dropout layer after
256         # every ReLU layer.
257         # ===== #
258
259         caches = []
260         dropout_caches = []
261         x = X
262         for i in range(1, self.num_layers + 1):

```

```

263         w, b = self.params['W{}'.format(i)], self.params['b{}'.format(i)]
264         if (i == self.num_layers):
265             scores, cache = affine_forward(x, w, b)
266         else:
267             if self.use_batchnorm:
268                 x, cache = affine_batchnorm_relu_forward(x, w, b,
self.params['gamma{}'.format(i)], self.params['beta{}'.format(i)],
self.bn_params[i-1])
269             else:
270                 x, cache = affine_relu_forward(x, w, b)
271
272             if self.use_dropout:
273                 x, dropout_cache = dropout_forward(x, self.dropout_param)
274                 dropout_caches.append(dropout_cache)
275             caches.append(cache)
276
277             # ===== #
278             # END YOUR CODE HERE
279             # ===== #
280
281             # If test mode return early
282             if mode == 'test':
283                 return scores
284
285             loss, grads = 0.0, {}
286             # ===== #
287             # YOUR CODE HERE:
288             # Implement the backwards pass of the FC net and store the
gradients
289             # in the grads dict, so that grads[k] is the gradient of
self.params[k]
290             # Be sure your L2 regularization includes a 0.5 factor.
291             #
292             # BATCHNORM: Incorporate the backward pass of the batchnorm.
293             #
294             # DROPOUT: Incorporate the backward pass of dropout.
295             # ===== #
296
297
298             sm_loss, dout = softmax_loss(scores, y)
299             reg_loss = 0.0
300             for i in range(1, self.num_layers + 1):
301                 reg_loss += 0.5 * self.reg *
(np.sum(self.params['W{}'.format(i)]**2))
302             loss = sm_loss + reg_loss
303
304             for i in range(self.num_layers, 0, -1):
305                 cur_w = 'W{}'.format(i)
306                 cur_b = 'b{}'.format(i)
307                 if (i == self.num_layers):
308                     dout, grads[cur_w], grads[cur_b] = affine_backward(dout,
caches.pop())
309                 else:
310                     if self.use_dropout:
311                         dout = dropout_backward(dout, dropout_caches.pop())
312                     if self.use_batchnorm:
313                         dout, grads[cur_w], grads[cur_b], grads['gamma{}'.format(i)],
grads['beta{}'.format(i)] = affine_batchnorm_relu_backward(dout,
caches.pop())
314                 else:

```

```
315         dout, grads[cur_w], grads[cur_b] = affine_relu_backward(dout,
316 caches.pop())
317         grads[cur_w] += self.reg * self.params[cur_w]
318         # ===== #
319         # END YOUR CODE HERE
320         # ===== #
321
322     return loss, grads
323
```

```

1 from .layers import *
2
3 """
4 This code was originally written for CS 231n at Stanford University
5 (cs231n.stanford.edu). It has been modified in various areas for use in the
6 ECE 239AS class at UCLA. This includes the descriptions of what code to
7 implement as well as some slight potential changes in variable names to be
8 consistent with class nomenclature. We thank Justin Johnson & Serena Yeung
9 for
10 permission to use this code. To see the original version, please visit
11 cs231n.stanford.edu.
12 """
13 def affine_relu_forward(x, w, b):
14     """
15     Convenience layer that performs an affine transform followed by a ReLU
16
17     Inputs:
18     - x: Input to the affine layer
19     - w, b: Weights for the affine layer
20
21     Returns a tuple of:
22     - out: Output from the ReLU
23     - cache: Object to give to the backward pass
24     """
25     a, fc_cache = affine_forward(x, w, b)
26     out, relu_cache = relu_forward(a)
27     cache = (fc_cache, relu_cache)
28     return out, cache
29
30
31 def affine_relu_backward(dout, cache):
32     """
33     Backward pass for the affine-relu convenience layer
34     """
35     fc_cache, relu_cache = cache
36     da = relu_backward(dout, relu_cache)
37     dx, dw, db = affine_backward(da, fc_cache)
38     return dx, dw, db
39
40 def affine_batchnorm_relu_forward(x, w, b, gamma, beta, bn_params):
41     a, fc_cache = affine_forward(x, w, b)
42     b, batch_cache = batchnorm_forward(a, gamma, beta, bn_params)
43     out, relu_cache = relu_forward(b)
44     cache = (fc_cache, batch_cache, relu_cache)
45     return out, cache
46
47 def affine_batchnorm_relu_backward(dout, cache):
48     fc_cache, batch_cache, relu_cache = cache
49     da = relu_backward(dout, relu_cache)
50     db, dgamma, dbeta = batchnorm_backward(da, batch_cache)
51     dx, dw, db = affine_backward(db, fc_cache)
52     return dx, dw, db, dgamma, dbeta

```

```

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

```

```

55
56
57 def sgd_momentum(w, dw, config=None):
58     """
59     Performs stochastic gradient descent with momentum.
60
61     config format:
62     - learning_rate: Scalar learning rate.
63     - momentum: Scalar between 0 and 1 giving the momentum value.
64       Setting momentum = 0 reduces to sgd.
65     - velocity: A numpy array of the same shape as w and dw used to store a
moving
66     average of the gradients.
67     """
68     if config is None: config = {}
69     config.setdefault('learning_rate', 1e-2)
70     config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't
there
71     v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets
it to zero.
72
73     # ===== #
74     # YOUR CODE HERE:
75     #   Implement the momentum update formula. Return the updated weights
76     #   as next_w, and the updated velocity as v.
77     # ===== #
78
79     alpha = config.get('momentum')
80     eps = config.get('learning_rate')
81     v = alpha*v - eps*dw
82     next_w = w+v
83
84     # ===== #
85     # END YOUR CODE HERE
86     # ===== #
87
88     config['velocity'] = v
89
90     return next_w, config
91
92 def sgd_nesterov_momentum(w, dw, config=None):
93     """
94     Performs stochastic gradient descent with Nesterov momentum.
95
96     config format:
97     - learning_rate: Scalar learning rate.
98     - momentum: Scalar between 0 and 1 giving the momentum value.
99       Setting momentum = 0 reduces to sgd.
100    - velocity: A numpy array of the same shape as w and dw used to store a
moving
101    average of the gradients.
102    """
103    if config is None: config = {}
104    config.setdefault('learning_rate', 1e-2)
105    config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't
there
106    v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets
it to zero.
107
108    # ===== #

```

```

109 # YOUR CODE HERE:
110 # Implement the momentum update formula. Return the updated weights
111 # as next_w, and the updated velocity as v.
112 # ===== #
113 alpha = config.get('momentum')
114 eps = config.get('learning_rate')
115
116 v_old = v
117 v = alpha*v_old - eps*dw
118 next_w = w + v + alpha*(v - v_old)
119
120 # ===== #
121 # END YOUR CODE HERE
122 # ===== #
123
124 config['velocity'] = v
125
126 return next_w, config
127
128 def rmsprop(w, dw, config=None):
129     """
130     Uses the RMSProp update rule, which uses a moving average of squared
131     gradient
132     values to set adaptive per-parameter learning rates.
133
134     config format:
135     - learning_rate: Scalar learning rate.
136     - decay_rate: Scalar between 0 and 1 giving the decay rate for the
137     squared
138     gradient cache.
139     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
140     - beta: Moving average of second moments of gradients.
141     """
142     if config is None: config = {}
143     config.setdefault('learning_rate', 1e-2)
144     config.setdefault('decay_rate', 0.99)
145     config.setdefault('epsilon', 1e-8)
146     config.setdefault('a', np.zeros_like(w))
147
148     next_w = None
149
150     # ===== #
151     # YOUR CODE HERE:
152     # Implement RMSProp. Store the next value of w as next_w. You need
153     # to also store in config['a'] the moving average of the second
154     # moment gradients, so they can be used for future gradients.
155
156     Concretely,
157     # config['a'] corresponds to "a" in the lecture notes.
158     # ===== #
159
160     v = config.get('epsilon')
161     a = config.get('a')
162     lr = config.get('learning_rate')
163     beta = config.get('decay_rate')
164
165     a = beta*a + (1-beta)*(dw**2)
166     next_w = w - np.multiply(lr / (np.sqrt(a) + v), dw)
167
168     config['a'] = a

```



```

166 # ===== #
167 # END YOUR CODE HERE
168 # ===== #
169
170 return next_w, config
171
172
173 def adam(w, dw, config=None):
174     """
175     Uses the Adam update rule, which incorporates moving averages of both the
176     gradient and its square and a bias correction term.
177
178     config format:
179     - learning_rate: Scalar learning rate.
180     - beta1: Decay rate for moving average of first moment of gradient.
181     - beta2: Decay rate for moving average of second moment of gradient.
182     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
183     - m: Moving average of gradient.
184     - v: Moving average of squared gradient.
185     - t: Iteration number.
186     """
187     if config is None: config = {}
188     config.setdefault('learning_rate', 1e-3)
189     config.setdefault('beta1', 0.9)
190     config.setdefault('beta2', 0.999)
191     config.setdefault('epsilon', 1e-8)
192     config.setdefault('v', np.zeros_like(w))
193     config.setdefault('a', np.zeros_like(w))
194     config.setdefault('t', 0)
195
196     next_w = None
197
198     # ===== #
199     # YOUR CODE HERE:
200     # Implement Adam. Store the next value of w as next_w. You need
201     # to also store in config['a'] the moving average of the second
202     # moment gradients, and in config['v'] the moving average of the
203     # first moments. Finally, store in config['t'] the increasing time.
204     # ===== #
205
206     t = config.get('t') + 1
207     v = config.get('v')
208     beta1 = config.get('beta1')
209     v = beta1*v + (1-beta1)*dw
210
211     a = config.get('a')
212     beta2 = config.get('beta2')
213     a = beta2*a + (1-beta2)*(dw**2)
214
215     v_u = v / (1 - beta1**t)
216     a_u = a / (1 - beta2**t)
217
218     lr = config.get('learning_rate')
219     eps = config.get('epsilon')
220     next_w = w - (lr / (np.sqrt(a_u) + eps) * v_u)
221
222     config['a'] = a
223     config['v'] = v
224     config['t'] = t
225

```

```
226 # ===== #
227 # END YOUR CODE HERE
228 # ===== #
229
230 return next_w, config
231
```

```

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

```

```

60     - b: A numpy array of biases, of shape (M,)
61
62     Returns a tuple of:
63     - dx: Gradient with respect to x, of shape (N, d1, ..., d_k)
64     - dw: Gradient with respect to w, of shape (D, M)
65     - db: Gradient with respect to b, of shape (M,)
66     """
67     x, w, b = cache
68     dx, dw, db = None, None, None
69
70     # ===== #
71     # YOUR CODE HERE:
72     # Calculate the gradients for the backward pass.
73     # Notice:
74     # dout is N x M
75     # dx should be N x d1 x ... x dk; it relates to dout through
multiplication with w, which is D x M
76     # dw should be D x M; it relates to dout through multiplication with x,
which is N x D after reshaping
77     # db should be M; it is just the sum over dout examples
78     # ===== #
79
80     dx = dout.dot(w.T).reshape(x.shape)
81     dw = x.reshape(x.shape[0], -1).T.dot(dout)
82     db = dout.sum(axis=0)
83
84     # ===== #
85     # END YOUR CODE HERE
86     # ===== #
87
88     return dx, dw, db
89
90 def relu_forward(x):
91     """
92     Computes the forward pass for a layer of rectified linear units (ReLUs).
93
94     Input:
95     - x: Inputs, of any shape
96
97     Returns a tuple of:
98     - out: Output, of the same shape as x
99     - cache: x
100    """
101    # ===== #
102    # YOUR CODE HERE:
103    # Implement the ReLU forward pass.
104    # ===== #
105
106    out = x * (x > 0)
107
108    # ===== #
109    # END YOUR CODE HERE
110    # ===== #
111
112    cache = x
113    return out, cache
114
115 def relu_backward(dout, cache):
116     """
117

```

```

118     Computes the backward pass for a layer of rectified linear units (ReLUs).
119
120     Input:
121     - dout: Upstream derivatives, of any shape
122     - cache: Input x, of same shape as dout
123
124     Returns:
125     - dx: Gradient with respect to x
126     """
127     x = cache
128
129     # ===== #
130     # YOUR CODE HERE:
131     #     Implement the ReLU backward pass
132     # ===== #
133
134     dx = dout * (cache > 0)
135
136     # ===== #
137     # END YOUR CODE HERE
138     # ===== #
139
140     return dx
141
142 def batchnorm_forward(x, gamma, beta, bn_param):
143     """
144     Forward pass for batch normalization.
145
146     During training the sample mean and (uncorrected) sample variance are
147     computed from minibatch statistics and used to normalize the incoming
148     data.
149     During training we also keep an exponentially decaying running mean of
150     the mean
151     and variance of each feature, and these averages are used to normalize
152     data
153     at test-time.
154
155     At each timestep we update the running averages for mean and variance
156     using
157     an exponential decay based on the momentum parameter:
158
159     running_mean = momentum * running_mean + (1 - momentum) * sample_mean
160     running_var = momentum * running_var + (1 - momentum) * sample_var
161
162     Note that the batch normalization paper suggests a different test-time
163     behavior: they compute sample mean and variance for each feature using a
164     large number of training images rather than using a running average. For
165     this implementation we have chosen to use running averages instead since
166     they do not require an additional estimation step; the torch7
167     implementation
168     of batch normalization also uses running averages.
169
170     Input:
171     - x: Data of shape (N, D)
172     - gamma: Scale parameter of shape (D,)
173     - beta: Shift parameter of shape (D,)
174     - bn_param: Dictionary with the following keys:
175         - mode: 'train' or 'test'; required
176         - eps: Constant for numeric stability
177         - momentum: Constant for running mean / variance.

```

```

173     - running_mean: Array of shape (D,) giving running mean of features
174     - running_var: Array of shape (D,) giving running variance of features
175
176 Returns a tuple of:
177 - out: of shape (N, D)
178 - cache: A tuple of values needed in the backward pass
179 """
180 mode = bn_param['mode']
181 eps = bn_param.get('eps', 1e-5)
182 momentum = bn_param.get('momentum', 0.9)
183
184 N, D = x.shape
185 running_mean = bn_param.get('running_mean', np.zeros(D, dtype=x.dtype))
186 running_var = bn_param.get('running_var', np.zeros(D, dtype=x.dtype))
187
188 out, cache = None, None
189 if mode == 'train':
190
191     # ===== #
192     # YOUR CODE HERE:
193     #   A few steps here:
194     #   (1) Calculate the running mean and variance of the minibatch.
195     #   (2) Normalize the activations with the running mean and
196 variance.
197     #   (3) Scale and shift the normalized activations. Store this
198     #       as the variable 'out'
199     #   (4) Store any variables you may need for the backward pass in
200     #       the 'cache' variable.
201     # ===== #
202
203     mean = np.mean(x, axis=0)
204     var = np.sum((x - mean)**2, axis=0) / N
205
206     running_mean = momentum * running_mean + (1 - momentum) * mean
207     running_var = momentum * running_var + (1 - momentum) * var
208
209     var_eps_sum_inv = 1 / np.sqrt(var + eps)
210     x_mean_diff = (x - mean)
211     x_n = var_eps_sum_inv * x_mean_diff
212     out = gamma * x_n + beta
213
214     cache = (x_n, x, gamma, var_eps_sum_inv, x_mean_diff)
215
216     # ===== #
217     # END YOUR CODE HERE
218     # ===== #
219 elif mode == 'test':
220     # ===== #
221     # YOUR CODE HERE:
222     #   Calculate the testing time normalized activation. Normalize
223 using
224     #   the running mean and variance, and then scale and shift
225 appropriately.
226     #   Store the output as 'out'.
227     # ===== #
228
229     x_n = (x - running_mean) / np.sqrt(running_var + eps)
230     out = gamma * x_n + beta
231
232     # ===== #

```

```

230         # END YOUR CODE HERE
231         # ===== #
232     else:
233         raise ValueError('Invalid forward batchnorm mode "%s"' % mode)
234
235     # Store the updated running means back into bn_param
236     bn_param['running_mean'] = running_mean
237     bn_param['running_var'] = running_var
238
239     return out, cache
240
241 def batchnorm_backward(dout, cache):
242     """
243     Backward pass for batch normalization.
244
245     For this implementation, you should write out a computation graph for
246     batch normalization on paper and propagate gradients backward through
247     intermediate nodes.
248
249     Inputs:
250     - dout: Upstream derivatives, of shape (N, D)
251     - cache: Variable of intermediates from batchnorm_forward.
252
253     Returns a tuple of:
254     - dx: Gradient with respect to inputs x, of shape (N, D)
255     - dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
256     - dbeta: Gradient with respect to shift parameter beta, of shape (D,)
257     """
258     dx, dgamma, dbeta = None, None, None
259
260     # ===== #
261     # YOUR CODE HERE:
262     # Implement the batchnorm backward pass, calculating dx, dgamma, and
263     # ===== #
264     dbeta.
265
266     x_n, x, gamma, var_eps_sum_inv, x_mean_diff = cache
267     N = x.shape[0]
268     dgamma = np.sum(dout * x_n, axis=0)
269     dbeta = np.sum(dout, axis=0)
270
271     dxhat = dout * gamma
272     da = var_eps_sum_inv * dxhat
273     db = x_mean_diff * dxhat
274     dc = -(var_eps_sum_inv**2) * db
275     de = 0.5 * var_eps_sum_inv * dc
276     dsig = np.sum(de, axis=0)
277     dmu = -var_eps_sum_inv * np.sum(dxhat, axis=0) - dsig * (2/N) *
278     np.sum(x_mean_diff, axis=0)
279
280     dx = da + (2*x_mean_diff/N) * dsig + (dmu/N)
281
282     # ===== #
283     # END YOUR CODE HERE
284     # ===== #
285
286     return dx, dgamma, dbeta
287
288 def dropout_forward(x, dropout_param):
289     """

```

```

288     Performs the forward pass for (inverted) dropout.
289
290     Inputs:
291     - x: Input data, of any shape
292     - dropout_param: A dictionary with the following keys:
293         - p: Dropout parameter. We keep each neuron output with probability p.
294         - mode: 'test' or 'train'. If the mode is train, then perform dropout;
295           if the mode is test, then just return the input.
296         - seed: Seed for the random number generator. Passing seed makes this
297           function deterministic, which is needed for gradient checking but not
in
298         real networks.
299
300     Outputs:
301     - out: Array of the same shape as x.
302     - cache: A tuple (dropout_param, mask). In training mode, mask is the
dropout
303     mask that was used to multiply the input; in test mode, mask is None.
304     """
305     p, mode = dropout_param['p'], dropout_param['mode']
306     if 'seed' in dropout_param:
307         np.random.seed(dropout_param['seed'])
308
309     mask = None
310     out = None
311
312     if mode == 'train':
313         # ===== #
314         # YOUR CODE HERE:
315         #     Implement the inverted dropout forward pass during training time.
316
317         #     Store the masked and scaled activations in out, and store the
318         #     dropout mask as the variable mask.
319         # ===== #
320
321         mask = (np.random.rand(*x.shape) < p) / p
322         out = x * mask
323
324         # ===== #
325         # END YOUR CODE HERE
326         # ===== #
327     elif mode == 'test':
328
329         # ===== #
330         # YOUR CODE HERE:
331         #     Implement the inverted dropout forward pass during test time.
332         # ===== #
333
334         out = x
335
336         # ===== #
337         # END YOUR CODE HERE
338         # ===== #
339
340     cache = (dropout_param, mask)
341     out = out.astype(x.dtype, copy=False)
342
343     return out, cache
344

```



```

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

```

```

403     dx[np.arange(N), y] -= num_pos
404     dx /= N
405     return loss, dx
406
407
408 def softmax_loss(x, y):
409     """
410     Computes the loss and gradient for softmax classification.
411
412     Inputs:
413     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth
class
414       for the ith input.
415     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
416       0 <= y[i] < C
417
418     Returns a tuple of:
419     - loss: Scalar giving the loss
420     - dx: Gradient of the loss with respect to x
421     """
422
423     probs = np.exp(x - np.max(x, axis=1, keepdims=True))
424     probs /= np.sum(probs, axis=1, keepdims=True)
425     N = x.shape[0]
426     loss = -np.sum(np.log(probs[np.arange(N), y])) / N
427     dx = probs.copy()
428     dx[np.arange(N), y] -= 1
429     dx /= N
430     return loss, dx
431

```

# Optimization

February 12, 2020

## 0.1 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.

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

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

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

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

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

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

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

```
data = get_CIFAR10_data()
```

```

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

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

```

## 0.2 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`

### 0.2.1 Test all functions you copy and pasted

```

In [3]: from nndl.layer_tests import *

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

```

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

If `affine_backward` is working, error should be less than  $1e-9$ ::  
 dx error:  $3.856575683705302e-09$   
 dw error:  $1.1399491310769201e-10$   
 db error:  $5.219731727479396e-11$

If `relu_forward` function is working, difference should be around  $1e-8$ :  
 difference:  $4.999999798022158e-08$

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

```
dx error: 3.275596545811805e-12
```

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

```
dx error: 2.383177924909149e-09
```

```
dw error: 1.83548767388022e-10
```

```
db error: 4.082423466114065e-12
```

Running check with `reg = 0`

```
Initial loss: 2.304275608380092
```

```
W1 relative error: 8.417754181333801e-06
```

```
W2 relative error: 5.398832855664135e-06
```

```
W3 relative error: 1.5441950905045603e-06
```

```
b1 relative error: 1.633739334783617e-08
```

```
b2 relative error: 1.891295003590916e-09
```

```
b3 relative error: 6.558446622618634e-11
```

Running check with `reg = 3.14`

```
Initial loss: 7.137305601344814
```

```
W1 relative error: 2.8720064677982203e-08
```

```
W2 relative error: 1.0381035536940991e-08
```

```
W3 relative error: 4.5518133076821233e-08
```

```
b1 relative error: 6.266803412277905e-08
```

```
b2 relative error: 1.5191260046727443e-08
```

```
b3 relative error: 1.3955883581600704e-10
```

## 1 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.

### 1.1 SGD + momentum

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

```
In [4]: from nndl.optim import sgd_momentum
```

```
N, D = 4, 5
```

```
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
```

```
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
```

```
v = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)
```

```
config = {'learning_rate': 1e-3, 'velocity': v}
```

```
next_w, _ = sgd_momentum(w, dw, config=config)
```

```

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

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

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

```

## 1.2 SGD + Nesterov momentum

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

In [5]: `from ndl.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

```

### 1.3 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.

```
In [6]: num_train = 4000
        small_data = {
            'X_train': data['X_train'][:num_train],
            'y_train': data['y_train'][:num_train],
            'X_val': data['X_val'],
            'y_val': data['y_val'],
        }

        solvers = {}

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

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

        fig, axes = plt.subplots(3, 1)

        ax = axes[0]
        ax.set_title('Training loss')
        ax.set_xlabel('Iteration')

        ax = axes[1]
        ax.set_title('Training accuracy')
        ax.set_xlabel('Epoch')

        ax = axes[2]
        ax.set_title('Validation accuracy')
        ax.set_xlabel('Epoch')
```

```

for update_rule, solver in solvers.items():
    ax = axes[0]
    ax.plot(solver.loss_history, 'o', label=update_rule)

    ax = axes[1]
    ax.plot(solver.train_acc_history, '-o', label=update_rule)

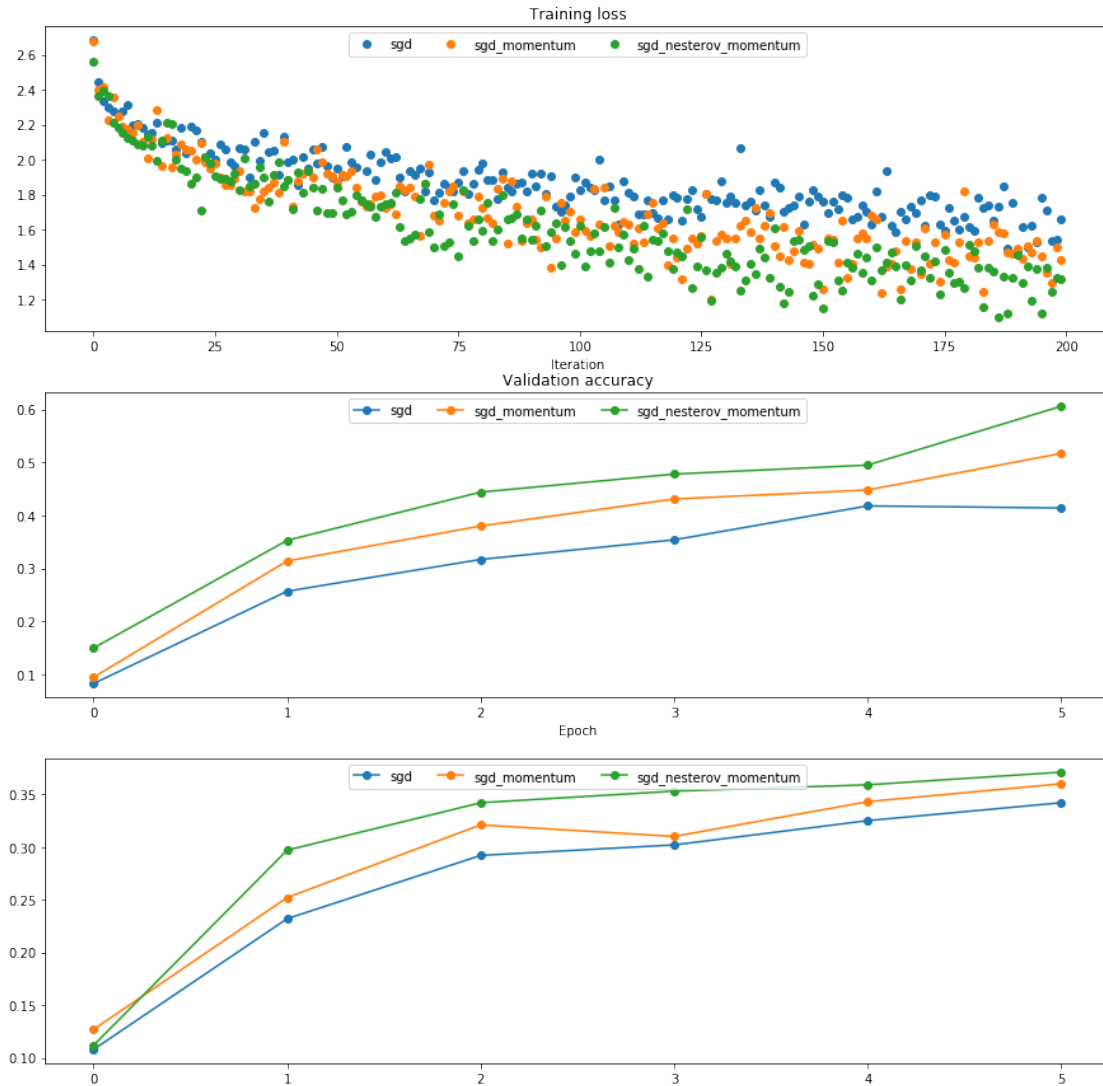
    ax = axes[2]
    ax.plot(solver.val_acc_history, '-o', label=update_rule)

for i in [1, 2, 3]:
    ax = axes[i - 1]
    ax.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`





## 1.4 RMSProp

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

In [7]: `from nndl.optim import rmsprop`

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

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

```

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

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

```

```

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

```

## 1.5 Adaptive moments

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

```

In [8]: # Test Adam implementation; you should see errors around 1e-7 or less
from nndl.optim import adam

```

```

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

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

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

```

```
[ 0.67473684,  0.69421053,  0.71368421,  0.73315789,  0.75263158],
[ 0.77210526,  0.79157895,  0.81105263,  0.83052632,  0.85      ]])
```

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

```
next_w error: 1.1395691798535431e-07
a error: 4.208314038113071e-09
v error: 4.214963193114416e-09
```

## 1.6 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.

```
In [9]: learning_rates = {'rmsprop': 2e-4, 'adam': 1e-3}
```

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

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

fig, axes = plt.subplots(3, 1)

ax = axes[0]
ax.set_title('Training loss')
ax.set_xlabel('Iteration')

ax = axes[1]
ax.set_title('Training accuracy')
ax.set_xlabel('Epoch')

ax = axes[2]
ax.set_title('Validation accuracy')
ax.set_xlabel('Epoch')
```

```

for update_rule, solver in solvers.items():
    ax = axes[0]
    ax.plot(solver.loss_history, 'o', label=update_rule)

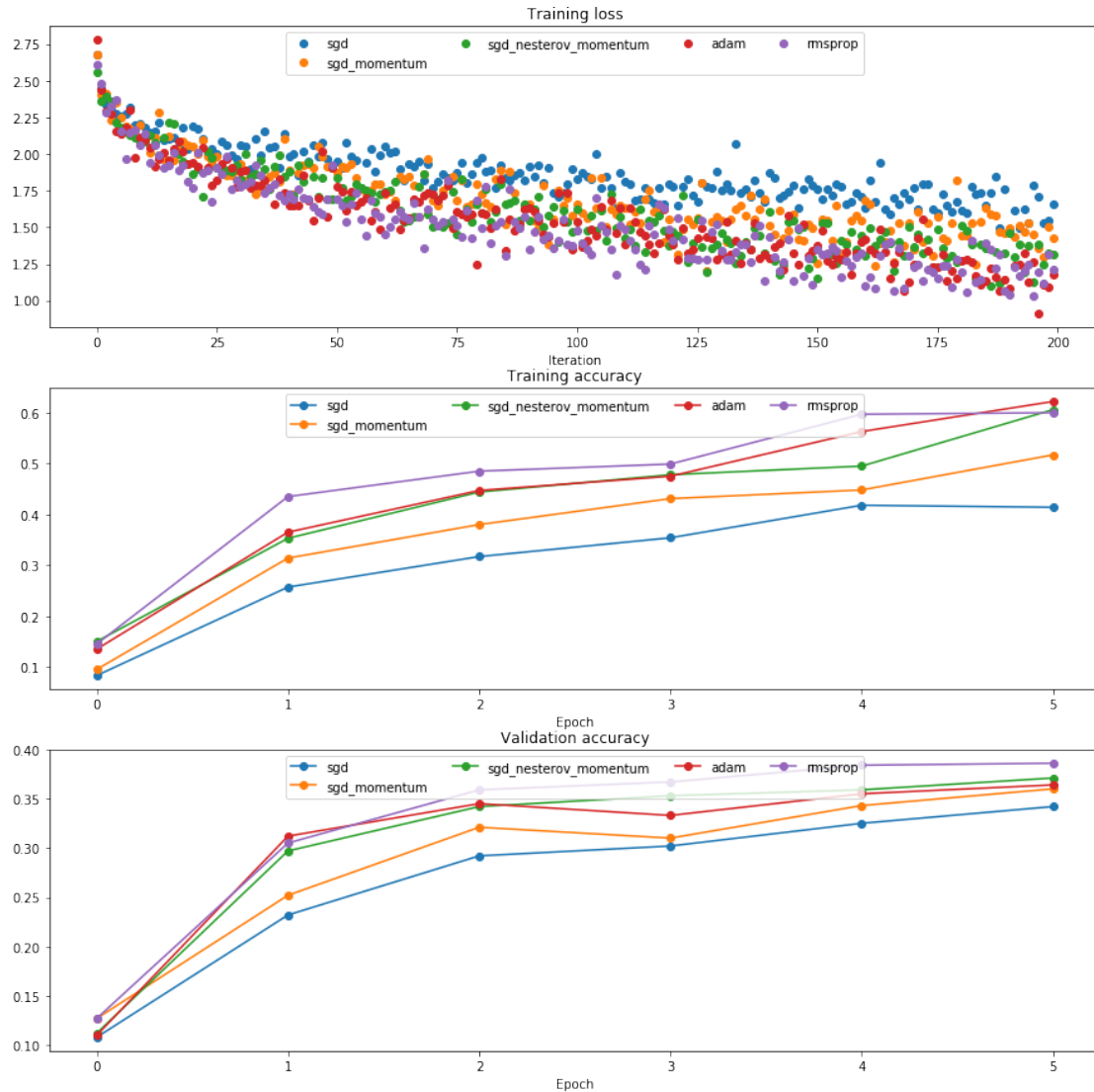
    ax = axes[1]
    ax.plot(solver.train_acc_history, '-o', label=update_rule)

    ax = axes[2]
    ax.plot(solver.val_acc_history, '-o', label=update_rule)

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

```

Optimizing with adam  
 Optimizing with rmsprop



## 1.7 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.

```
In [10]: optimizer = 'adam'
         best_model = None

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

```

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

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

solver.train()

(Iteration 1 / 4900) loss: 2.313495
(Epoch 0 / 10) train acc: 0.156000; val_acc: 0.172000
(Iteration 51 / 4900) loss: 1.910098
(Iteration 101 / 4900) loss: 1.714113
(Iteration 151 / 4900) loss: 1.877789
(Iteration 201 / 4900) loss: 1.653570
(Iteration 251 / 4900) loss: 1.608789
(Iteration 301 / 4900) loss: 1.682090
(Iteration 351 / 4900) loss: 1.653318
(Iteration 401 / 4900) loss: 1.520741
(Iteration 451 / 4900) loss: 1.724113
(Epoch 1 / 10) train acc: 0.446000; val_acc: 0.426000
(Iteration 501 / 4900) loss: 1.544322
(Iteration 551 / 4900) loss: 1.527392
(Iteration 601 / 4900) loss: 1.489695
(Iteration 651 / 4900) loss: 1.560724
(Iteration 701 / 4900) loss: 1.644997
(Iteration 751 / 4900) loss: 1.679174
(Iteration 801 / 4900) loss: 1.531670
(Iteration 851 / 4900) loss: 1.607276
(Iteration 901 / 4900) loss: 1.585658
(Iteration 951 / 4900) loss: 1.537134
(Epoch 2 / 10) train acc: 0.478000; val_acc: 0.451000
(Iteration 1001 / 4900) loss: 1.662541
(Iteration 1051 / 4900) loss: 1.497583
(Iteration 1101 / 4900) loss: 1.347863
(Iteration 1151 / 4900) loss: 1.315546
(Iteration 1201 / 4900) loss: 1.391136
(Iteration 1251 / 4900) loss: 1.424261
(Iteration 1301 / 4900) loss: 1.429473
(Iteration 1351 / 4900) loss: 1.449509
(Iteration 1401 / 4900) loss: 1.517006
(Iteration 1451 / 4900) loss: 1.246213
(Epoch 3 / 10) train acc: 0.501000; val_acc: 0.489000

```

(Iteration 1501 / 4900) loss: 1.473849  
(Iteration 1551 / 4900) loss: 1.469772  
(Iteration 1601 / 4900) loss: 1.206832  
(Iteration 1651 / 4900) loss: 1.558015  
(Iteration 1701 / 4900) loss: 1.261236  
(Iteration 1751 / 4900) loss: 1.283223  
(Iteration 1801 / 4900) loss: 1.147611  
(Iteration 1851 / 4900) loss: 1.295039  
(Iteration 1901 / 4900) loss: 1.461217  
(Iteration 1951 / 4900) loss: 1.262842  
(Epoch 4 / 10) train acc: 0.566000; val\_acc: 0.501000  
(Iteration 2001 / 4900) loss: 1.274226  
(Iteration 2051 / 4900) loss: 1.259313  
(Iteration 2101 / 4900) loss: 1.390710  
(Iteration 2151 / 4900) loss: 1.129790  
(Iteration 2201 / 4900) loss: 1.166850  
(Iteration 2251 / 4900) loss: 1.097454  
(Iteration 2301 / 4900) loss: 1.355590  
(Iteration 2351 / 4900) loss: 1.390996  
(Iteration 2401 / 4900) loss: 1.044549  
(Epoch 5 / 10) train acc: 0.540000; val\_acc: 0.497000  
(Iteration 2451 / 4900) loss: 1.165498  
(Iteration 2501 / 4900) loss: 1.160091  
(Iteration 2551 / 4900) loss: 1.183314  
(Iteration 2601 / 4900) loss: 1.026732  
(Iteration 2651 / 4900) loss: 1.092496  
(Iteration 2701 / 4900) loss: 1.128480  
(Iteration 2751 / 4900) loss: 1.157050  
(Iteration 2801 / 4900) loss: 1.156004  
(Iteration 2851 / 4900) loss: 1.297502  
(Iteration 2901 / 4900) loss: 1.214705  
(Epoch 6 / 10) train acc: 0.600000; val\_acc: 0.508000  
(Iteration 2951 / 4900) loss: 1.200887  
(Iteration 3001 / 4900) loss: 0.932238  
(Iteration 3051 / 4900) loss: 1.087853  
(Iteration 3101 / 4900) loss: 1.194822  
(Iteration 3151 / 4900) loss: 1.073080  
(Iteration 3201 / 4900) loss: 1.199952  
(Iteration 3251 / 4900) loss: 1.170383  
(Iteration 3301 / 4900) loss: 1.162548  
(Iteration 3351 / 4900) loss: 1.339277  
(Iteration 3401 / 4900) loss: 1.163136  
(Epoch 7 / 10) train acc: 0.597000; val\_acc: 0.503000  
(Iteration 3451 / 4900) loss: 1.037553  
(Iteration 3501 / 4900) loss: 0.868173  
(Iteration 3551 / 4900) loss: 1.296146  
(Iteration 3601 / 4900) loss: 1.108326  
(Iteration 3651 / 4900) loss: 1.123196

```

(Iteration 3701 / 4900) loss: 1.132850
(Iteration 3751 / 4900) loss: 1.023332
(Iteration 3801 / 4900) loss: 1.197241
(Iteration 3851 / 4900) loss: 1.072664
(Iteration 3901 / 4900) loss: 1.036933
(Epoch 8 / 10) train acc: 0.650000; val_acc: 0.532000
(Iteration 3951 / 4900) loss: 0.911799
(Iteration 4001 / 4900) loss: 0.937844
(Iteration 4051 / 4900) loss: 1.229134
(Iteration 4101 / 4900) loss: 1.001406
(Iteration 4151 / 4900) loss: 1.190598
(Iteration 4201 / 4900) loss: 0.964132
(Iteration 4251 / 4900) loss: 1.061063
(Iteration 4301 / 4900) loss: 0.923960
(Iteration 4351 / 4900) loss: 0.900450
(Iteration 4401 / 4900) loss: 0.939629
(Epoch 9 / 10) train acc: 0.657000; val_acc: 0.497000
(Iteration 4451 / 4900) loss: 0.885418
(Iteration 4501 / 4900) loss: 0.812531
(Iteration 4551 / 4900) loss: 0.834904
(Iteration 4601 / 4900) loss: 1.116214
(Iteration 4651 / 4900) loss: 0.990060
(Iteration 4701 / 4900) loss: 0.930970
(Iteration 4751 / 4900) loss: 0.921465
(Iteration 4801 / 4900) loss: 0.923087
(Iteration 4851 / 4900) loss: 0.877881
(Epoch 10 / 10) train acc: 0.698000; val_acc: 0.519000

```

```

In [11]: y_test_pred = np.argmax(model.loss(data['X_test']), axis=1)
         y_val_pred = np.argmax(model.loss(data['X_val']), axis=1)
         print('Validation set accuracy: {}'.format(np.mean(y_val_pred == data['y_val'])))
         print('Test set accuracy: {}'.format(np.mean(y_test_pred == data['y_test'])))

```

Validation set accuracy: 0.532

Test set accuracy: 0.521