

Batch-Normalization

February 10, 2023

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.

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

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

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

```
%reload_ext autoreload
```

```
[ ]: # 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.

```
[ ]: # 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: [-33.3140234  5.05791416 27.44645111]
```

```

stds: [27.28819317 24.23404506 37.06806543]
After batch normalization (gamma=1, beta=0)
mean: [-4.76285678e-16  4.21884749e-17 -7.21644966e-17]
std:  [0.99999999 0.99999999 1.          ]
After batch normalization (nontrivial gamma, beta)
means: [11. 12. 13.]
stds:  [0.99999999 1.99999998 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.

```

[ ]: # 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: [-0.06106739 -0.07932008  0.06920021]
stds:  [0.96350622 1.0371408  0.96590731]

```

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.

```
[ ]: # 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 gamma: batchnorm_forward(x, gamma, beta, bn_param)[0]
fb = lambda beta: 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:  2.283087106464751e-09
dgamma error:  2.2470551124754713e-12
dbeta error:  3.2755980183511937e-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 $1e-4$.

```
[ ]: 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: 2.260019462554289
W1 relative error: 2.2939796059403806e-05
W2 relative error: 2.5208679187998513e-05
W3 relative error: 3.7421749411570136e-10
b1 relative error: 2.7755575615628914e-09
b2 relative error: 2.220446049250313e-08
b3 relative error: 1.1267290312950767e-10
beta1 relative error: 1.673833264554499e-08
beta2 relative error: 7.642706478241727e-08
gamma1 relative error: 1.1554041624835892e-08
gamma2 relative error: 7.790797076454733e-08

```

```

Running check with reg = 3.14
Initial loss: 7.39844875087809
W1 relative error: 0.000652963884682915
W2 relative error: 8.078787959531888e-06
W3 relative error: 2.1674688235598615e-08
b1 relative error: 5.551115123125783e-09
b2 relative error: 5.551115123125783e-09
b3 relative error: 1.809701771811536e-10
beta1 relative error: 1.2958120554539437e-08
beta2 relative error: 1.4006846971724922e-06
gamma1 relative error: 1.4831508388965654e-08
gamma2 relative error: 2.743147438054341e-06

```

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.

```

[ ]: # 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: 2.296010
(Epoch 0 / 10) train acc: 0.168000; val_acc: 0.184000
(Epoch 1 / 10) train acc: 0.227000; val_acc: 0.244000
(Epoch 2 / 10) train acc: 0.323000; val_acc: 0.258000
(Epoch 3 / 10) train acc: 0.390000; val_acc: 0.277000
(Epoch 4 / 10) train acc: 0.436000; val_acc: 0.314000
(Epoch 5 / 10) train acc: 0.456000; val_acc: 0.304000
(Epoch 6 / 10) train acc: 0.522000; val_acc: 0.340000
(Epoch 7 / 10) train acc: 0.569000; val_acc: 0.332000
(Epoch 8 / 10) train acc: 0.585000; val_acc: 0.321000
(Epoch 9 / 10) train acc: 0.637000; val_acc: 0.310000
(Epoch 10 / 10) train acc: 0.686000; val_acc: 0.337000

```

```

(Iteration 1 / 200) loss: 2.302708
(Epoch 0 / 10) train acc: 0.104000; val_acc: 0.117000
(Epoch 1 / 10) train acc: 0.171000; val_acc: 0.142000
(Epoch 2 / 10) train acc: 0.248000; val_acc: 0.217000
(Epoch 3 / 10) train acc: 0.211000; val_acc: 0.199000
(Epoch 4 / 10) train acc: 0.287000; val_acc: 0.244000
(Epoch 5 / 10) train acc: 0.299000; val_acc: 0.244000
(Epoch 6 / 10) train acc: 0.337000; val_acc: 0.261000
(Epoch 7 / 10) train acc: 0.345000; val_acc: 0.231000
(Epoch 8 / 10) train acc: 0.412000; val_acc: 0.288000
(Epoch 9 / 10) train acc: 0.455000; val_acc: 0.292000
(Epoch 10 / 10) train acc: 0.450000; val_acc: 0.295000

```

```

[ ]: 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')

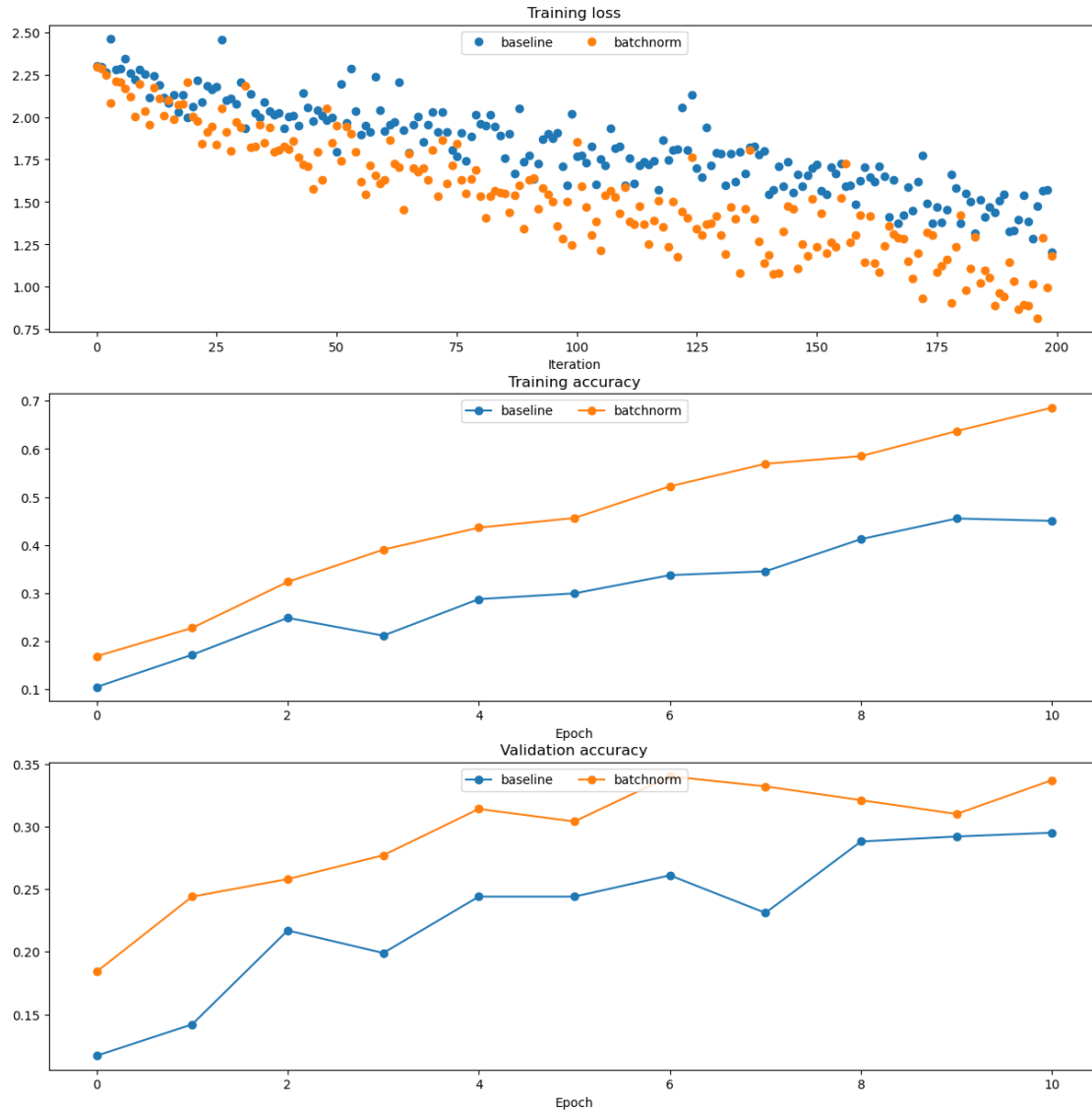
ax = axes[0]
ax.plot(solver.loss_history, 'o', label='baseline')
ax.plot(bn_solver.loss_history, 'o', label='batchnorm')

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.

```
[ ]: # 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

```

```

[ ]: # 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')

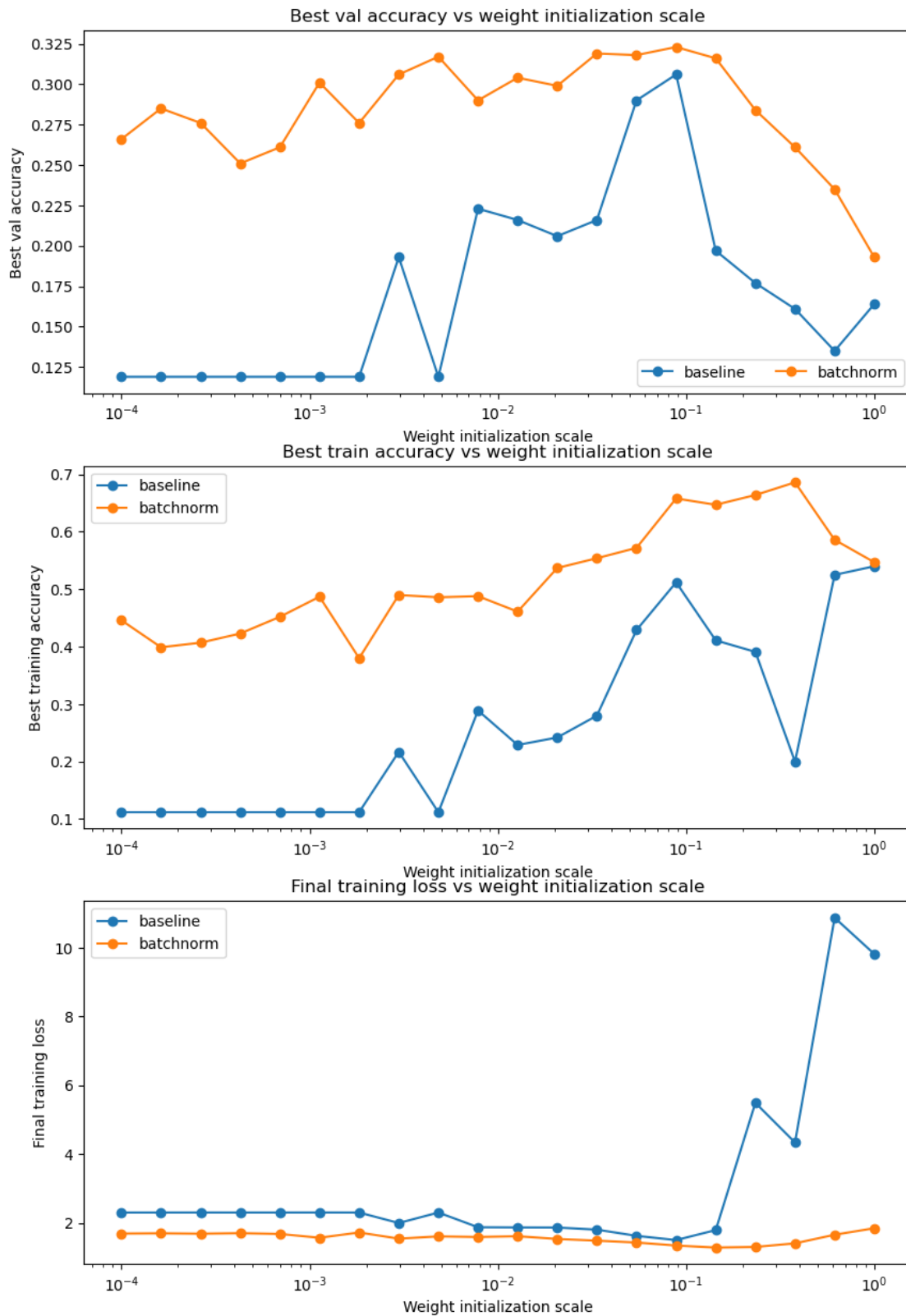
plt.subplot(3, 1, 2)
plt.title('Best train accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Best training accuracy')
plt.semilogx(weight_scales, best_train_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_train_accs, '-o', label='batchnorm')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

For the baseline, when the weight initialization is very small, we could observe from the first 2 plots that the accuracy is very low. It implies that the gradient vanishing happens. And when the weight initialization is too large, we could also observe from the 3rd plots that the training loss increases very fast which implies that the gradient exploding happens. The results tell us that baseline is very sensitive to the weight initialization.

For the batchnorm, we could find that the impact of weight initialization is much decreased because the normalization helps to strict the value of too large/small gradients. But from the first 2 plots, when weight initialization is too large, the overfitting occurs.

Generally, I think BN could help to restrict the impact of gradient vanishing/exploding and also help to prevent the overfitting at some extent.

Dropout

February 10, 2023

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.

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

```
[ ]: ## 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, \u
    ↪ 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))))
```

```
[ ]: # 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.

```
[ ]: 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())
    print( )
```

```
Running tests with p = 0.3
Mean of input: 10.002898472632062
Mean of train-time output: 10.025619592528209
Mean of test-time output: 10.002898472632062
Fraction of train-time output set to zero: 0.699396
Fraction of test-time output set to zero: 0.0
```

```
Running tests with p = 0.6
Mean of input: 10.002898472632062
Mean of train-time output: 10.015838393782959
Mean of test-time output: 10.002898472632062
Fraction of train-time output set to zero: 0.399256
Fraction of test-time output set to zero: 0.0
```

```
Running tests with p = 0.75
Mean of input: 10.002898472632062
```

```
Mean of train-time output: 10.015224946990443
Mean of test-time output: 10.002898472632062
Fraction of train-time output set to zero: 0.248928
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:

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

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

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

```
dx relative error: 5.4456104932567176e-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 W1 gradient relative error is on the order of $1e-6$ (the largest of all the relative errors).

```
[ ]: 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: 2.3050396176612384
W1 relative error: 7.839091367974606e-07
W2 relative error: 2.5656696352835975e-07
W3 relative error: 7.66140281048889e-08
b1 relative error: 4.461314594443551e-09
b2 relative error: 7.630256229692641e-10
b3 relative error: 1.7671730210721276e-10

```

```

Running check with dropout = 0.75
Initial loss: 2.2924165725936616
W1 relative error: 1.913764152741491e-07
W2 relative error: 1.3174644422188826e-07
W3 relative error: 3.246749313511893e-07
b1 relative error: 7.519165775969613e-09
b2 relative error: 1.9613765149833703e-09
b3 relative error: 1.4849637974979418e-10

```

```

Running check with dropout = 1.0
Initial loss: 2.3055208155971316
W1 relative error: 5.474025246189813e-07
W2 relative error: 7.600892283230569e-08
W3 relative error: 4.657460911645141e-08
b1 relative error: 6.88212284937315e-09
b2 relative error: 1.5631103560839323e-09
b3 relative error: 1.6519207248783954e-10

```

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.

```
[ ]: # 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: 2.299614
(Epoch 0 / 25) train acc: 0.158000; val_acc: 0.155000
(Epoch 1 / 25) train acc: 0.162000; val_acc: 0.162000
(Epoch 2 / 25) train acc: 0.254000; val_acc: 0.231000
(Epoch 3 / 25) train acc: 0.298000; val_acc: 0.256000
(Epoch 4 / 25) train acc: 0.344000; val_acc: 0.284000
(Epoch 5 / 25) train acc: 0.342000; val_acc: 0.264000
(Epoch 6 / 25) train acc: 0.416000; val_acc: 0.297000
(Epoch 7 / 25) train acc: 0.458000; val_acc: 0.299000
(Epoch 8 / 25) train acc: 0.472000; val_acc: 0.309000
(Epoch 9 / 25) train acc: 0.558000; val_acc: 0.310000
(Epoch 10 / 25) train acc: 0.610000; val_acc: 0.323000
(Epoch 11 / 25) train acc: 0.610000; val_acc: 0.292000
(Epoch 12 / 25) train acc: 0.658000; val_acc: 0.322000
(Epoch 13 / 25) train acc: 0.732000; val_acc: 0.313000
(Epoch 14 / 25) train acc: 0.752000; val_acc: 0.290000
(Epoch 15 / 25) train acc: 0.792000; val_acc: 0.324000
(Epoch 16 / 25) train acc: 0.818000; val_acc: 0.303000
(Epoch 17 / 25) train acc: 0.848000; val_acc: 0.289000
(Epoch 18 / 25) train acc: 0.882000; val_acc: 0.296000
(Epoch 19 / 25) train acc: 0.924000; val_acc: 0.289000
(Epoch 20 / 25) train acc: 0.934000; val_acc: 0.284000
(Iteration 101 / 125) loss: 0.231835
(Epoch 21 / 25) train acc: 0.952000; val_acc: 0.285000

```

```

(Epoch 22 / 25) train acc: 0.970000; val_acc: 0.303000
(Epoch 23 / 25) train acc: 0.984000; val_acc: 0.275000
(Epoch 24 / 25) train acc: 0.988000; val_acc: 0.285000
(Epoch 25 / 25) train acc: 0.984000; val_acc: 0.279000
(Iteration 1 / 125) loss: 2.304071
(Epoch 0 / 25) train acc: 0.220000; val_acc: 0.170000
(Epoch 1 / 25) train acc: 0.290000; val_acc: 0.206000
(Epoch 2 / 25) train acc: 0.270000; val_acc: 0.213000
(Epoch 3 / 25) train acc: 0.328000; val_acc: 0.255000
(Epoch 4 / 25) train acc: 0.392000; val_acc: 0.298000
(Epoch 5 / 25) train acc: 0.402000; val_acc: 0.293000
(Epoch 6 / 25) train acc: 0.436000; val_acc: 0.299000
(Epoch 7 / 25) train acc: 0.508000; val_acc: 0.322000
(Epoch 8 / 25) train acc: 0.560000; val_acc: 0.323000
(Epoch 9 / 25) train acc: 0.592000; val_acc: 0.324000
(Epoch 10 / 25) train acc: 0.634000; val_acc: 0.302000
(Epoch 11 / 25) train acc: 0.670000; val_acc: 0.298000
(Epoch 12 / 25) train acc: 0.760000; val_acc: 0.317000
(Epoch 13 / 25) train acc: 0.736000; val_acc: 0.310000
(Epoch 14 / 25) train acc: 0.788000; val_acc: 0.311000
(Epoch 15 / 25) train acc: 0.782000; val_acc: 0.300000
(Epoch 16 / 25) train acc: 0.854000; val_acc: 0.304000
(Epoch 17 / 25) train acc: 0.882000; val_acc: 0.300000
(Epoch 18 / 25) train acc: 0.908000; val_acc: 0.267000
(Epoch 19 / 25) train acc: 0.924000; val_acc: 0.300000
(Epoch 20 / 25) train acc: 0.946000; val_acc: 0.281000
(Iteration 101 / 125) loss: 0.234402
(Epoch 21 / 25) train acc: 0.966000; val_acc: 0.295000
(Epoch 22 / 25) train acc: 0.980000; val_acc: 0.281000
(Epoch 23 / 25) train acc: 0.978000; val_acc: 0.283000
(Epoch 24 / 25) train acc: 0.982000; val_acc: 0.305000
(Epoch 25 / 25) train acc: 0.992000; val_acc: 0.296000

```

```

[ ]: # 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')

```

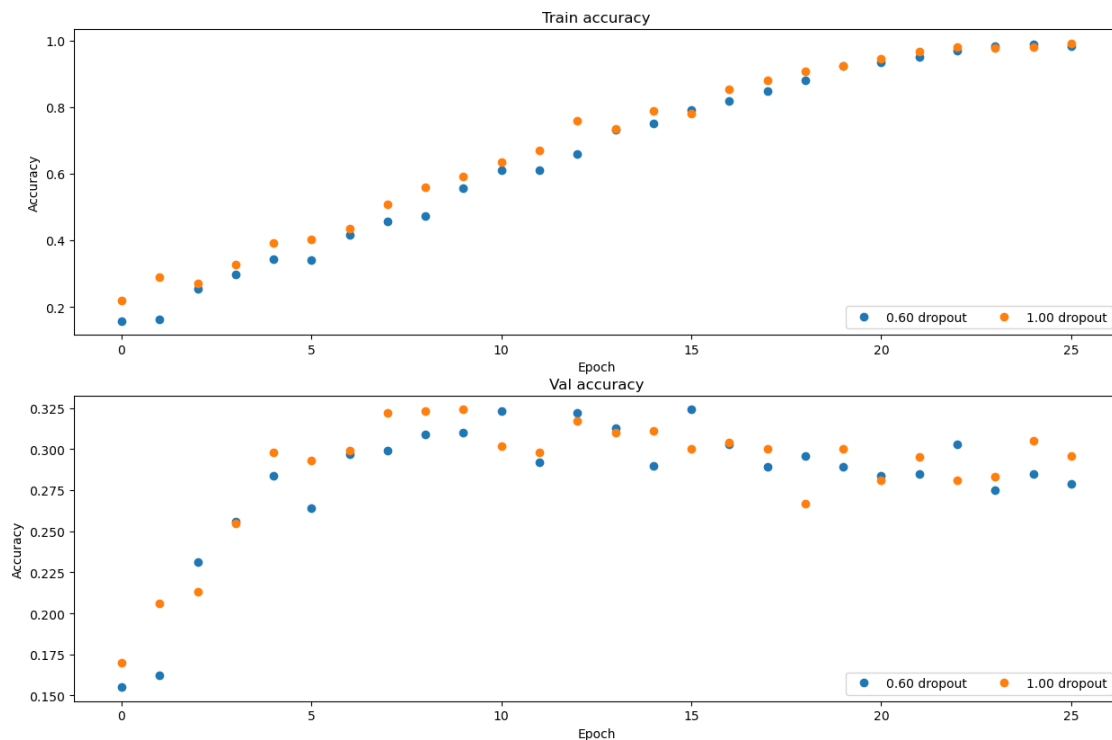
```

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:

I think here dropout's impact is not obvious considering we only trained on 500 data. The training accuracy is almost 100% but the validation accuracy is only about 30%. It implies that our model

is seriously overfitting and we do need dropout.

I changed the dropout from [0.6, 1.0] to [0.2, 1.0]. The results showed that model still achieved nearly 100% on training accuracy, but on validation accuracy with dropout the network performed better obviously compared to network without dropout.

Then I changed the num_train to 5000 and tested it again. The result showed that with/without dropout the training accuracy still nearly to 90%. But on validation data, with dropout the network performed better (not too much though).

Therefore we could say that dropout's capability: it could help to increase our model's generalization and prevent the overfitting at some extent.

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.

```
[ ]: # ===== #
# YOUR CODE HERE:
#   Implement a FC-net that achieves at least 55% validation accuracy
#   on CIFAR-10.
# ===== #
optimizer = 'adam'
dropout=0.6
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, dropout = dropout)

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()

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

(Iteration 1 / 4900) loss: 2.335502

(Epoch 0 / 10) train acc: 0.155000; val_acc: 0.166000

(Iteration 51 / 4900) loss: 1.886194
(Iteration 101 / 4900) loss: 1.833522
(Iteration 151 / 4900) loss: 1.899841
(Iteration 201 / 4900) loss: 1.610528
(Iteration 251 / 4900) loss: 1.404449
(Iteration 301 / 4900) loss: 1.631404
(Iteration 351 / 4900) loss: 1.855381
(Iteration 401 / 4900) loss: 1.708536
(Iteration 451 / 4900) loss: 1.649088
(Epoch 1 / 10) train acc: 0.447000; val_acc: 0.449000
(Iteration 501 / 4900) loss: 1.621702
(Iteration 551 / 4900) loss: 1.495131
(Iteration 601 / 4900) loss: 1.505964
(Iteration 651 / 4900) loss: 1.464100
(Iteration 701 / 4900) loss: 1.504684
(Iteration 751 / 4900) loss: 1.487932
(Iteration 801 / 4900) loss: 1.551561
(Iteration 851 / 4900) loss: 1.412798
(Iteration 901 / 4900) loss: 1.461804
(Iteration 951 / 4900) loss: 1.455468
(Epoch 2 / 10) train acc: 0.496000; val_acc: 0.482000
(Iteration 1001 / 4900) loss: 1.543609
(Iteration 1051 / 4900) loss: 1.407979
(Iteration 1101 / 4900) loss: 1.411821
(Iteration 1151 / 4900) loss: 1.285107
(Iteration 1201 / 4900) loss: 1.381421
(Iteration 1251 / 4900) loss: 1.612857
(Iteration 1301 / 4900) loss: 1.434685
(Iteration 1351 / 4900) loss: 1.440963
(Iteration 1401 / 4900) loss: 1.556632
(Iteration 1451 / 4900) loss: 1.562366
(Epoch 3 / 10) train acc: 0.542000; val_acc: 0.507000
(Iteration 1501 / 4900) loss: 1.597746
(Iteration 1551 / 4900) loss: 1.219000
(Iteration 1601 / 4900) loss: 1.442981
(Iteration 1651 / 4900) loss: 1.468734
(Iteration 1701 / 4900) loss: 1.297483
(Iteration 1751 / 4900) loss: 1.543932
(Iteration 1801 / 4900) loss: 1.529823
(Iteration 1851 / 4900) loss: 1.478847
(Iteration 1901 / 4900) loss: 1.441085
(Iteration 1951 / 4900) loss: 1.617001
(Epoch 4 / 10) train acc: 0.547000; val_acc: 0.549000
(Iteration 2001 / 4900) loss: 1.346974
(Iteration 2051 / 4900) loss: 1.407445
(Iteration 2101 / 4900) loss: 1.271642
(Iteration 2151 / 4900) loss: 1.608927
(Iteration 2201 / 4900) loss: 1.432519

(Iteration 2251 / 4900) loss: 1.403091
(Iteration 2301 / 4900) loss: 1.231440
(Iteration 2351 / 4900) loss: 1.260361
(Iteration 2401 / 4900) loss: 1.280835
(Epoch 5 / 10) train acc: 0.577000; val_acc: 0.553000
(Iteration 2451 / 4900) loss: 1.469316
(Iteration 2501 / 4900) loss: 1.210299
(Iteration 2551 / 4900) loss: 1.254554
(Iteration 2601 / 4900) loss: 1.238156
(Iteration 2651 / 4900) loss: 1.262420
(Iteration 2701 / 4900) loss: 1.289816
(Iteration 2751 / 4900) loss: 1.298516
(Iteration 2801 / 4900) loss: 1.324259
(Iteration 2851 / 4900) loss: 1.237017
(Iteration 2901 / 4900) loss: 1.258834
(Epoch 6 / 10) train acc: 0.585000; val_acc: 0.574000
(Iteration 2951 / 4900) loss: 1.273547
(Iteration 3001 / 4900) loss: 1.275819
(Iteration 3051 / 4900) loss: 1.148007
(Iteration 3101 / 4900) loss: 1.340733
(Iteration 3151 / 4900) loss: 1.362481
(Iteration 3201 / 4900) loss: 1.279735
(Iteration 3251 / 4900) loss: 1.152038
(Iteration 3301 / 4900) loss: 1.138591
(Iteration 3351 / 4900) loss: 1.284542
(Iteration 3401 / 4900) loss: 1.145100
(Epoch 7 / 10) train acc: 0.614000; val_acc: 0.556000
(Iteration 3451 / 4900) loss: 1.152454
(Iteration 3501 / 4900) loss: 1.057248
(Iteration 3551 / 4900) loss: 1.189292
(Iteration 3601 / 4900) loss: 1.263029
(Iteration 3651 / 4900) loss: 1.143761
(Iteration 3701 / 4900) loss: 1.268868
(Iteration 3751 / 4900) loss: 1.089459
(Iteration 3801 / 4900) loss: 1.033787
(Iteration 3851 / 4900) loss: 1.426417
(Iteration 3901 / 4900) loss: 1.216828
(Epoch 8 / 10) train acc: 0.593000; val_acc: 0.569000
(Iteration 3951 / 4900) loss: 1.308395
(Iteration 4001 / 4900) loss: 1.214057
(Iteration 4051 / 4900) loss: 1.152115
(Iteration 4101 / 4900) loss: 1.283493
(Iteration 4151 / 4900) loss: 1.297368
(Iteration 4201 / 4900) loss: 1.207147
(Iteration 4251 / 4900) loss: 1.205163
(Iteration 4301 / 4900) loss: 1.183897
(Iteration 4351 / 4900) loss: 1.144597
(Iteration 4401 / 4900) loss: 1.354404

```
(Epoch 9 / 10) train acc: 0.643000; val_acc: 0.575000
(Iteration 4451 / 4900) loss: 1.123538
(Iteration 4501 / 4900) loss: 1.196001
(Iteration 4551 / 4900) loss: 1.001292
(Iteration 4601 / 4900) loss: 1.160576
(Iteration 4651 / 4900) loss: 1.212783
(Iteration 4701 / 4900) loss: 1.230949
(Iteration 4751 / 4900) loss: 1.254236
(Iteration 4801 / 4900) loss: 1.232851
(Iteration 4851 / 4900) loss: 1.364168
(Epoch 10 / 10) train acc: 0.652000; val_acc: 0.580000
```

```
[ ]: 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: 0.58
Test set accuracy: 0.564
```

fc_net

February 10, 2023

```
[ ]: # %load fc_net.py
import numpy as np
import pdb

from .layers import *
from .layer_utils import *

class TwoLayerNet(object):
    """
    A two-layer fully-connected neural network with ReLU nonlinearity and
    softmax loss that uses a modular layer design. We assume an input dimension
    of  $D$ , a hidden dimension of  $H$ , and perform classification over  $C$  classes.

    The architecture should be affine - relu - affine - softmax.

    Note that this class does not implement gradient descent; instead, it
    will interact with a separate Solver object that is responsible for running
    optimization.

    The learnable parameters of the model are stored in the dictionary
    self.params that maps parameter names to numpy arrays.
    """

    def __init__(self, input_dim=3*32*32, hidden_dims=100, num_classes=10,
                  dropout=1, weight_scale=1e-3, reg=0.0):
        """
        Initialize a new network.

        Inputs:
        - input_dim: An integer giving the size of the input
        - hidden_dims: An integer giving the size of the hidden layer
        - num_classes: An integer giving the number of classes to classify
        - dropout: Scalar between 0 and 1 giving dropout strength.
        - weight_scale: Scalar giving the standard deviation for random
          initialization of the weights.
        - reg: Scalar giving L2 regularization strength.
        """
```



```

"""
self.params = {}
self.reg = reg

# ===== #
# YOUR CODE HERE:
# Initialize W1, W2, b1, and b2. Store these as self.params['W1'],
# self.params['W2'], self.params['b1'] and self.params['b2']. The
# biases are initialized to zero and the weights are initialized
# so that each parameter has mean 0 and standard deviation 1/weight_scale.
# The dimensions of W1 should be (input_dim, hidden_dim) and the
# dimensions of W2 should be (hidden_dims, num_classes)
# ===== #
W1, W2 = self.params['W1'], self.params['W2']
b1, b2 = self.params['b1'], self.params['b2']
output_first_layer, cache_first = affine_relu_forward(X, W1, b1)
scores, cache_sec = affine_forward(output_first_layer, W2, b2)

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

def loss(self, X, y=None):
    """
    Compute loss and gradient for a minibatch of data.

    Inputs:
    - X: Array of input data of shape (N, d_1, ..., d_k)
    - y: Array of labels, of shape (N,). y[i] gives the label for X[i].

    Returns:
    If y is None, then run a test-time forward pass of the model and return:
    - scores: Array of shape (N, C) giving classification scores, where
      scores[i, c] is the classification score for X[i] and class c.

    If y is not None, then run a training-time forward and backward pass and
    return a tuple of:
    - loss: Scalar value giving the loss
    - grads: Dictionary with the same keys as self.params, mapping parameter
      names to gradients of the loss with respect to those parameters.
    """
    scores = None

    # ===== #
    # YOUR CODE HERE:
    # Implement the forward pass of the two-layer neural network. Store

```

```

# the class scores as the variable 'scores'. Be sure to use the
↪ layers
# you prior implemented.
# ===== #

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

# If y is None then we are in test mode so just return scores
if y is None:
    return scores

loss, grads = 0, {}
# ===== #
# YOUR CODE HERE:
# Implement the backward pass of the two-layer neural net. Store
# the loss as the variable 'loss' and store the gradients in the
# 'grads' dictionary. For the grads dictionary, grads['W1'] holds
# the gradient for W1, grads['b1'] holds the gradient for b1, etc.
# i.e., grads[k] holds the gradient for self.params[k].
#
# Add L2 regularization, where there is an added cost  $0.5 * \text{self.reg} * W^2$ 
# for each W. Be sure to include the 0.5 multiplying factor to
# match our implementation.
#
# And be sure to use the layers you prior implemented.
# ===== #

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

return loss, grads

```

```

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=1, 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=1
    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 < 1
    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.
# =====
hidd_len = len(hidden_dims)
shape_one = input_dim
for i in range(hidd_len):
    self.params['W'+str(i+1)] = np.random.randn(shape_one,
↪hidden_dims[i]) * weight_scale
    self.params['b'+str(i+1)] = np.zeros(hidden_dims[i])
    shape_one = hidden_dims[i]
    if self.use_batchnorm is True:
        self.params['gamma'+str(i+1)] = np.ones(hidden_dims[i])
        self.params['beta'+str(i+1)] = np.zeros(hidden_dims[i])
    self.params['W'+str(hidd_len+1)] = np.random.randn(shape_one,
↪num_classes) * weight_scale
    self.params['b'+str(hidd_len+1)] = np.zeros(num_classes)
pass

# ===== #
# 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.
    # ===== #
    output = {}

```

```

output[0] = X
cache = {}
for i in range(self.num_layers):
    j = i + 1
    weights = str('W') + str(j)
    bias = str('b') + str(j)

    W_j = self.params[weights]
    b_j = self.params[bias]

    if j == self.num_layers:
        output[j], cache[i] = affine_forward(output[i], W_j, b_j)
    else:
        if self.use_dropout and self.use_batchnorm:
            gammas = str('gamma') + str(j)
            betas = str('beta') + str(j)
            gamma_j = self.params[gammas]
            beta_j = self.params[betas]
            output[j], cache[i] =
↪affine_batchnorm_relu_dropout_forward(output[i], W_j, b_j, gamma_j, beta_j,
↪self.bn_params[i], self.dropout_param)

            elif self.use_batchnorm is True:
                gammas = str('gamma') + str(j)
                betas = str('beta') + str(j)
                gamma_j = self.params[gammas]
                beta_j = self.params[betas]
                output[j], cache[i] = affine_batchnorm_relu_forward(output[i],
↪W_j, b_j, gamma_j, beta_j, self.bn_params[i])

            else:
                output[j], cache[i] = affine_relu_forward(output[i], W_j, b_j)

scores = output[self.num_layers]
pass

# ===== #
# 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)
reg_loss = 0
for i in range(self.num_layers):
    j = i + 1
    weights = str('W') + str(j)
    W_j = self.params[weights]
    reg_loss = reg_loss + 0.5 * self.reg * np.sum(W_j * W_j)
loss = loss + reg_loss

#compute grad
hidd_len = self.num_layers-1
# weights and bias of the last layer
dx, last_weight, last_bias = affine_backward(dout, cache[self.
↪num_layers-1])
    grads['W'+str(self.num_layers)] = last_weight + self.reg * self.
↪params['W'+str(self.num_layers)]
    grads['b'+str(self.num_layers)] = last_bias
# weights, bias, gamma, beta of the rest layers
if self.use_dropout and self.use_batchnorm:
    for i in range(hidd_len, 0, -1):
        dx, dw, db, dgamma, dbeta = ↪
↪affine_batchnorm_relu_dropout_backward(dx, cache[i-1])
        grads['W'+str(i)] = dw + self.reg * self.params['W'+str(i)]
        grads['b'+str(i)] = db
        grads['gamma'+str(i)] = dgamma
        grads['beta'+str(i)] = dbeta
elif self.use_batchnorm is True:
    for i in range(hidd_len, 0, -1):
        dx, dw, db, dgamma, dbeta = affine_batchnorm_relu_backward(dx, ↪
↪cache[i-1])
        grads['W'+str(i)] = dw + self.reg * self.params['W'+str(i)]
        grads['b'+str(i)] = db
        grads['gamma'+str(i)] = dgamma
        grads['beta'+str(i)] = dbeta
else:
    for i in range(self.num_layers-1, 0, -1):
        dx, dw, db = affine_relu_backward(dx, cache[i-1])

```

```
        grads['W'+str(i)] = dw + self.reg * self.params['W'+str(i)]
        grads['b'+str(i)] = db

    pass

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

    return loss, grads
```


layer_utils

February 10, 2023

```
[ ]: # %load layer_utils.py
from .layers import *

def affine_relu_forward(x, w, b):
    """
    Convenience layer that performs an affine transform followed by a ReLU

    Inputs:
    - x: Input to the affine layer
    - w, b: Weights for the affine layer

    Returns a tuple of:
    - out: Output from the ReLU
    - cache: Object to give to the backward pass
    """
    a, fc_cache = affine_forward(x, w, b)
    out, relu_cache = relu_forward(a)
    cache = (fc_cache, relu_cache)
    return out, cache

def affine_relu_backward(dout, cache):
    """
    Backward pass for the affine-relu convenience layer
    """
    fc_cache, relu_cache = cache
    da = relu_backward(dout, relu_cache)
    dx, dw, db = affine_backward(da, fc_cache)
    return dx, dw, db

def affine_batchnorm_relu_forward(x, w, b, gamma, beta, bn_param):
    #forward of Affine - BN - Relu
    aff_out, aff_cache = affine_forward(x, w, b)
    bn_out, bn_cache = batchnorm_forward(aff_out, gamma, beta, bn_param)
    out, relu_cache = relu_forward(bn_out)
    cache = (aff_cache, bn_cache, relu_cache)
    return out, cache
```

```

def affine_batchnorm_relu_backward(dout, cache):
    aff_cache, bn_cache, relu_cache = cache
    d_relu = relu_backward(dout, relu_cache)
    db, dgamma, dbeta = batchnorm_backward(d_relu, bn_cache)
    dx, dw, db = affine_backward(db, aff_cache)
    return dx, dw, db, dgamma, dbeta

def affine_batchnorm_relu_dropout_forward(x, w, b, gamma, beta, bn_param, dropout_param):
    #forward of Affine - BN - Relu
    aff_out, aff_cache = affine_forward(x, w, b)
    bn_out, bn_cache = batchnorm_forward(aff_out, gamma, beta, bn_param)
    relu_out, relu_cache = relu_forward(bn_out)
    out, dropout_cache = dropout_forward(relu_out, dropout_param)
    cache = (aff_cache, bn_cache, relu_cache, dropout_cache)
    return out, cache

def affine_batchnorm_relu_dropout_backward(dout, cache):
    aff_cache, bn_cache, relu_cache, dropout_cache = cache
    d_dropout = dropout_backward(dout, dropout_cache)
    d_relu = relu_backward(d_dropout, relu_cache)
    db, dgamma, dbeta = batchnorm_backward(d_relu, bn_cache)
    dx, dw, db = affine_backward(db, aff_cache)
    return dx, dw, db, dgamma, dbeta

```

layers

February 10, 2023

```
[ ]: # %load layers.py
import numpy as np
import pdb

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)
    """
    out = None
    # ===== #
    # 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.
    # ===== #

    num_inputs = x.shape[0]
    input_shape = x.shape[1:]
    input_size = np.prod(input_shape)
    x_reshape = x.reshape(num_inputs, input_size)
    out = np.dot(x_reshape, 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
    # ===== #
    num_inputs = x.shape[0]
    input_shape = x.shape[1:]
    input_size = np.prod(input_shape)
    x_reshape = x.reshape(num_inputs, input_size)
    x_shape = x.shape
    db = np.sum(dout, axis = 0)
    dx = np.dot(dout, w.T).reshape(x_shape)
    dw = np.dot(x_reshape.T, dout)

```

```

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

return dx, dw, db

def relu_forward(x):
    """
    Computes the forward pass for a layer of rectified linear units (ReLU).

    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(x, 0)
    # ===== #
    # 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 (ReLU).

    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
    # ===== #

```

```

x[x<0] = 0
x[x>0] = 1
dx = np.multiply(x,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.

    Input:
    - x: Data of shape (N, D)
    - gamma: Scale parameter of shape (D,)
    - beta: Shift paremeter 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.
    # ===== #
    #(1)
    sample_mean = np.sum(x, axis = 0) / N
    sample_var = np.var(x, axis = 0)
    #(2)
    x_mean = x - sample_mean
    x_normalized = x_mean / np.sqrt(sample_var + eps)
    #(3)
    out = gamma * x_normalized + beta
    #(4)
    cache = (x, x_normalized, gamma, beta, eps, sample_var, sample_mean)
    running_mean = momentum * running_mean + (1 - momentum) * sample_mean
    running_var = momentum * running_var + (1 - momentum) * sample_var

    pass

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

    # ===== #
    # 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.
    # ===== #
    N, D = dout.shape
    x, x_normalized, gamma, beta, eps, sample_var, sample_mean = cache
    dbeta = np.sum(dout, axis = 0)
    dgamma = np.sum(dout * x_normalized, axis = 0)
    dxnorm = dout * gamma

```



```

#implement db and da written on lecture slides
b = 1 / np.sqrt(eps + sample_var)
a = x - sample_mean
db = a * dxnorm
da = dxnorm * b

#dx has 3 sources totally
d_sample_var = np.sum(dxnorm * (-0.5 * (x - sample_mean) / ((sample_var +
↳eps)**1.5)), axis = 0)
d_sample_mean_1 = d_sample_var * np.sum(-2*(x - sample_mean), axis = 0) / N
d_sample_mean_2 = np.sum(-1 * da, axis = 0)
d_sample_mean = d_sample_mean_1 + d_sample_mean_2
dx1 = da
dx2 = d_sample_mean / N
dx3 = 2 * (x - sample_mean) * d_sample_var / N
dx = dx1 + dx2 + dx3

# ===== #
# 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 keep 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']
    assert (0<p<=1), "Dropout probability is not in (0,1]"
    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.rand(*x.shape) < p ) / p
    out = x * mask
    cache = (dropout_param, mask)

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

elif mode == 'test':

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

    pass
    # ===== #
    # 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

    pass
    # ===== #
    # END YOUR CODE HERE
    # ===== #
elif mode == 'test':
    # ===== #
    # YOUR CODE HERE:
    # Implement the inverted dropout backward pass during test time.
    # ===== #
    dx = dout

    pass
    # ===== #
    # 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(np.maximum(probs[np.arange(N), y], 1e-8))) / N
    dx = probs.copy()
    dx[np.arange(N), y] -= 1
    dx /= N
    return loss, dx

```

Optimization

February 11, 2023

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.

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

```
[ ]: ## 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, u
    ↪ 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))))
```

```
[ ]: # 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

```
[ ]: 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: 9.7698500479884e-10

If affine_backward is working, error should be less than 1e-9::
dx error: 9.921339690291491e-11
dw error: 5.118769370670294e-10
db error: 1.8832862889766754e-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: 3.275615306486218e-12

If `affine_relu_forward` and `affine_relu_backward` are working, error should be less than $1e-9$:
dx error: 1.9920118497528304e-09
dw error: 5.427422203418818e-10
db error: 7.826637978569618e-12

Running check with `reg = 0`
Initial loss: 2.3001838187602317
W1 relative error: 9.736428459439352e-07
W2 relative error: 1.9764259447160983e-05
W3 relative error: 3.7829964729163316e-07
b1 relative error: 1.4214626486789504e-08
b2 relative error: 2.101726270429346e-09
b3 relative error: 1.3010201989833611e-10
Running check with `reg = 3.14`
Initial loss: 6.833209685694318
W1 relative error: 1.6706600879849896e-08
W2 relative error: 4.024192483235043e-08
W3 relative error: 1.767462785238029e-08
b1 relative error: 5.647835445221942e-08
b2 relative error: 7.919412260379965e-09
b3 relative error: 1.7379827755901627e-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.

```
[ ]: 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`.

```

[ ]: 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
```

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.

```
[ ]: 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)

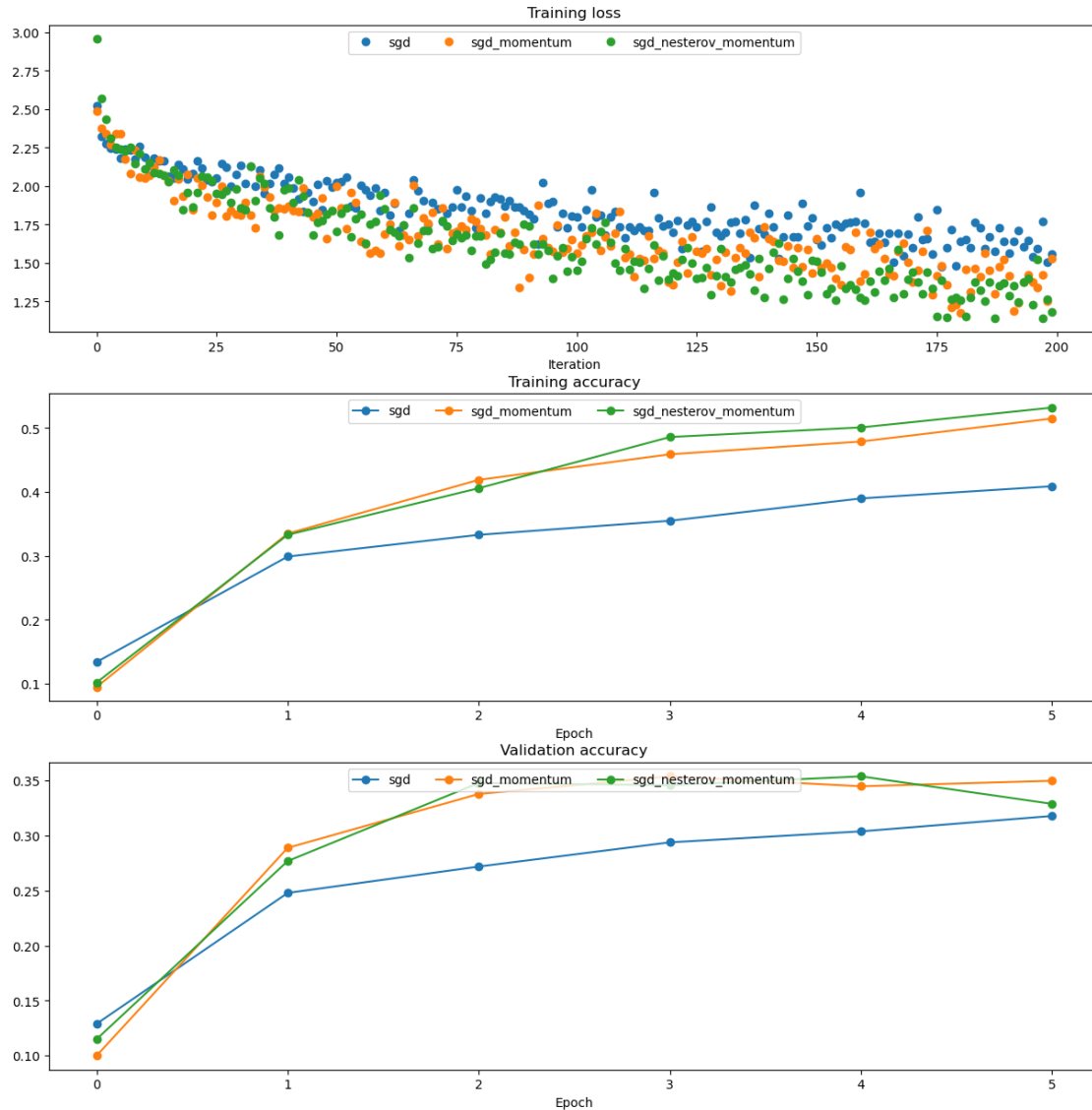
    ax = axes[1]
    ax.plot(solver.train_acc_history, '-o', label=update_rule)

    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.

```
[ ]: 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.

```

[ ]: # 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.

```

[ ]: 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)

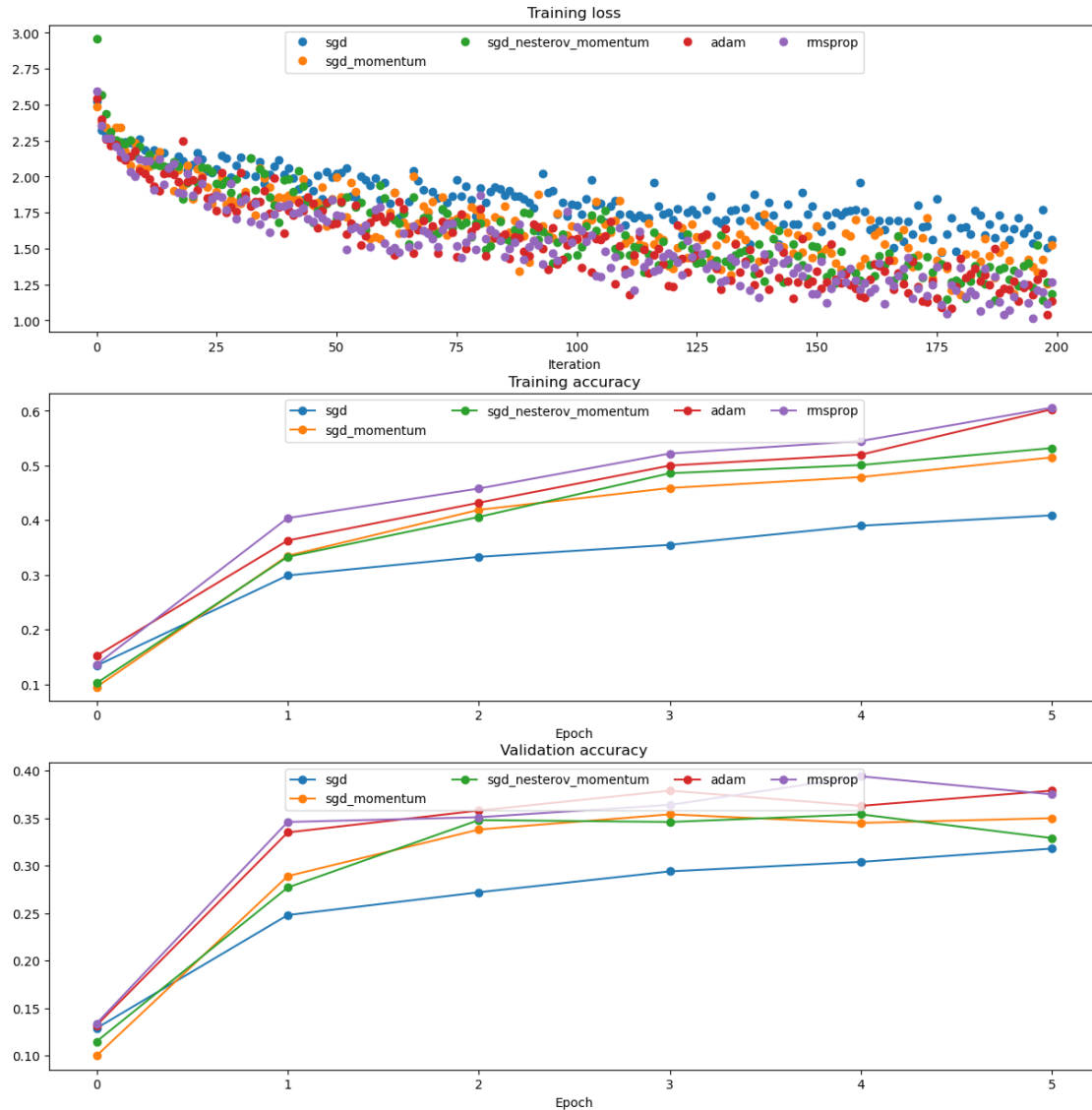
    ax = axes[1]
    ax.plot(solver.train_acc_history, '-o', label=update_rule)

    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.

```
[ ]: 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: 2.325989
(Epoch 0 / 10) train acc: 0.115000; val_acc: 0.127000
(Iteration 51 / 4900) loss: 1.870683
(Iteration 101 / 4900) loss: 1.749163
(Iteration 151 / 4900) loss: 1.645403
(Iteration 201 / 4900) loss: 1.786959
(Iteration 251 / 4900) loss: 1.493956
(Iteration 301 / 4900) loss: 1.709110
(Iteration 351 / 4900) loss: 1.603957
(Iteration 401 / 4900) loss: 1.669217
(Iteration 451 / 4900) loss: 1.806439
(Epoch 1 / 10) train acc: 0.402000; val_acc: 0.433000
(Iteration 501 / 4900) loss: 1.759044
(Iteration 551 / 4900) loss: 1.447526
(Iteration 601 / 4900) loss: 1.471226
(Iteration 651 / 4900) loss: 1.740063
(Iteration 701 / 4900) loss: 1.582539
(Iteration 751 / 4900) loss: 1.391246
(Iteration 801 / 4900) loss: 1.347183
(Iteration 851 / 4900) loss: 1.542118
(Iteration 901 / 4900) loss: 1.550798
(Iteration 951 / 4900) loss: 1.357640
(Epoch 2 / 10) train acc: 0.473000; val_acc: 0.449000
(Iteration 1001 / 4900) loss: 1.282699
(Iteration 1051 / 4900) loss: 1.382429
(Iteration 1101 / 4900) loss: 1.482079
(Iteration 1151 / 4900) loss: 1.375703
(Iteration 1201 / 4900) loss: 1.341312
(Iteration 1251 / 4900) loss: 1.434927
(Iteration 1301 / 4900) loss: 1.428094
(Iteration 1351 / 4900) loss: 1.514103
(Iteration 1401 / 4900) loss: 1.358413
(Iteration 1451 / 4900) loss: 1.282261

```


(Epoch 3 / 10) train acc: 0.532000; val_acc: 0.467000
(Iteration 1501 / 4900) loss: 1.287993
(Iteration 1551 / 4900) loss: 1.330630
(Iteration 1601 / 4900) loss: 1.464470
(Iteration 1651 / 4900) loss: 1.402027
(Iteration 1701 / 4900) loss: 1.356814
(Iteration 1751 / 4900) loss: 1.322569
(Iteration 1801 / 4900) loss: 1.296675
(Iteration 1851 / 4900) loss: 1.352026
(Iteration 1901 / 4900) loss: 1.162464
(Iteration 1951 / 4900) loss: 1.464830
(Epoch 4 / 10) train acc: 0.548000; val_acc: 0.489000
(Iteration 2001 / 4900) loss: 1.255898
(Iteration 2051 / 4900) loss: 1.250183
(Iteration 2101 / 4900) loss: 1.277716
(Iteration 2151 / 4900) loss: 1.298472
(Iteration 2201 / 4900) loss: 1.116835
(Iteration 2251 / 4900) loss: 1.354155
(Iteration 2301 / 4900) loss: 1.150200
(Iteration 2351 / 4900) loss: 1.208711
(Iteration 2401 / 4900) loss: 0.970558
(Epoch 5 / 10) train acc: 0.574000; val_acc: 0.496000
(Iteration 2451 / 4900) loss: 1.179215
(Iteration 2501 / 4900) loss: 1.390239
(Iteration 2551 / 4900) loss: 1.181364
(Iteration 2601 / 4900) loss: 1.210565
(Iteration 2651 / 4900) loss: 1.070150
(Iteration 2701 / 4900) loss: 1.052570
(Iteration 2751 / 4900) loss: 1.299136
(Iteration 2801 / 4900) loss: 1.183212
(Iteration 2851 / 4900) loss: 1.030881
(Iteration 2901 / 4900) loss: 1.125599
(Epoch 6 / 10) train acc: 0.600000; val_acc: 0.496000
(Iteration 2951 / 4900) loss: 1.108003
(Iteration 3001 / 4900) loss: 1.191459
(Iteration 3051 / 4900) loss: 1.125336
(Iteration 3101 / 4900) loss: 1.183858
(Iteration 3151 / 4900) loss: 1.193573
(Iteration 3201 / 4900) loss: 1.063729
(Iteration 3251 / 4900) loss: 1.044464
(Iteration 3301 / 4900) loss: 1.164088
(Iteration 3351 / 4900) loss: 0.986111
(Iteration 3401 / 4900) loss: 1.014930
(Epoch 7 / 10) train acc: 0.601000; val_acc: 0.523000
(Iteration 3451 / 4900) loss: 1.144177
(Iteration 3501 / 4900) loss: 1.106877
(Iteration 3551 / 4900) loss: 0.903499
(Iteration 3601 / 4900) loss: 1.282903

```

(Iteration 3651 / 4900) loss: 1.002898
(Iteration 3701 / 4900) loss: 0.842896
(Iteration 3751 / 4900) loss: 1.065477
(Iteration 3801 / 4900) loss: 1.031447
(Iteration 3851 / 4900) loss: 1.342536
(Iteration 3901 / 4900) loss: 1.151595
(Epoch 8 / 10) train acc: 0.656000; val_acc: 0.537000
(Iteration 3951 / 4900) loss: 0.921140
(Iteration 4001 / 4900) loss: 1.075565
(Iteration 4051 / 4900) loss: 1.020048
(Iteration 4101 / 4900) loss: 0.975876
(Iteration 4151 / 4900) loss: 1.109071
(Iteration 4201 / 4900) loss: 0.818720
(Iteration 4251 / 4900) loss: 1.149080
(Iteration 4301 / 4900) loss: 0.961885
(Iteration 4351 / 4900) loss: 1.029657
(Iteration 4401 / 4900) loss: 0.992939
(Epoch 9 / 10) train acc: 0.663000; val_acc: 0.538000
(Iteration 4451 / 4900) loss: 0.900818
(Iteration 4501 / 4900) loss: 0.878764
(Iteration 4551 / 4900) loss: 0.805838
(Iteration 4601 / 4900) loss: 0.921345
(Iteration 4651 / 4900) loss: 0.920808
(Iteration 4701 / 4900) loss: 0.990059
(Iteration 4751 / 4900) loss: 0.826135
(Iteration 4801 / 4900) loss: 1.000700
(Iteration 4851 / 4900) loss: 0.904633
(Epoch 10 / 10) train acc: 0.701000; val_acc: 0.544000

```

```

[ ]: 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: 0.544
Test set accuracy: 0.539

```