

Optimization

February 16, 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`.

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

import time
import numpy as np
import matplotlib.pyplot as plt
from nndl.fc_net import *
from utils.data_utils import get_CIFAR10_data
from utils.gradient_check import eval_numerical_gradient, 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))))
```

```
[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

```
[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: $9.769848888397517e-10$

If `affine_backward` is working, error should be less than $1e-9$::
 dx error: $3.1350052789632175e-10$
 dw error: $3.304426457498361e-10$
 db error: $5.904796441023863e-12$

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: $3.2756200311074745e-12$

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

dx error: $1.311080410589162e-10$

dw error: $2.468488919627008e-10$

db error: $2.367137183126194e-11$

Running check with `reg = 0`

Initial loss: 2.304896869438183

W1 relative error: $4.003267597777456e-07$

W2 relative error: $3.4233399144768428e-06$

W3 relative error: $3.950178418014187e-07$

b1 relative error: $2.1096629320130773e-08$

b2 relative error: $3.785720248684948e-09$

b3 relative error: $1.7125759588652542e-10$

Running check with `reg = 3.14`

Initial loss: 6.919525031482678

W1 relative error: $1.7140271958755806e-08$

W2 relative error: $2.3993413030278107e-08$

W3 relative error: $5.434866705091047e-08$

b1 relative error: $4.234547395026865e-08$

b2 relative error: $9.652542799698091e-09$

b3 relative error: $3.0170199340259497e-10$

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.

```
[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`.

```

[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.

```
[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)

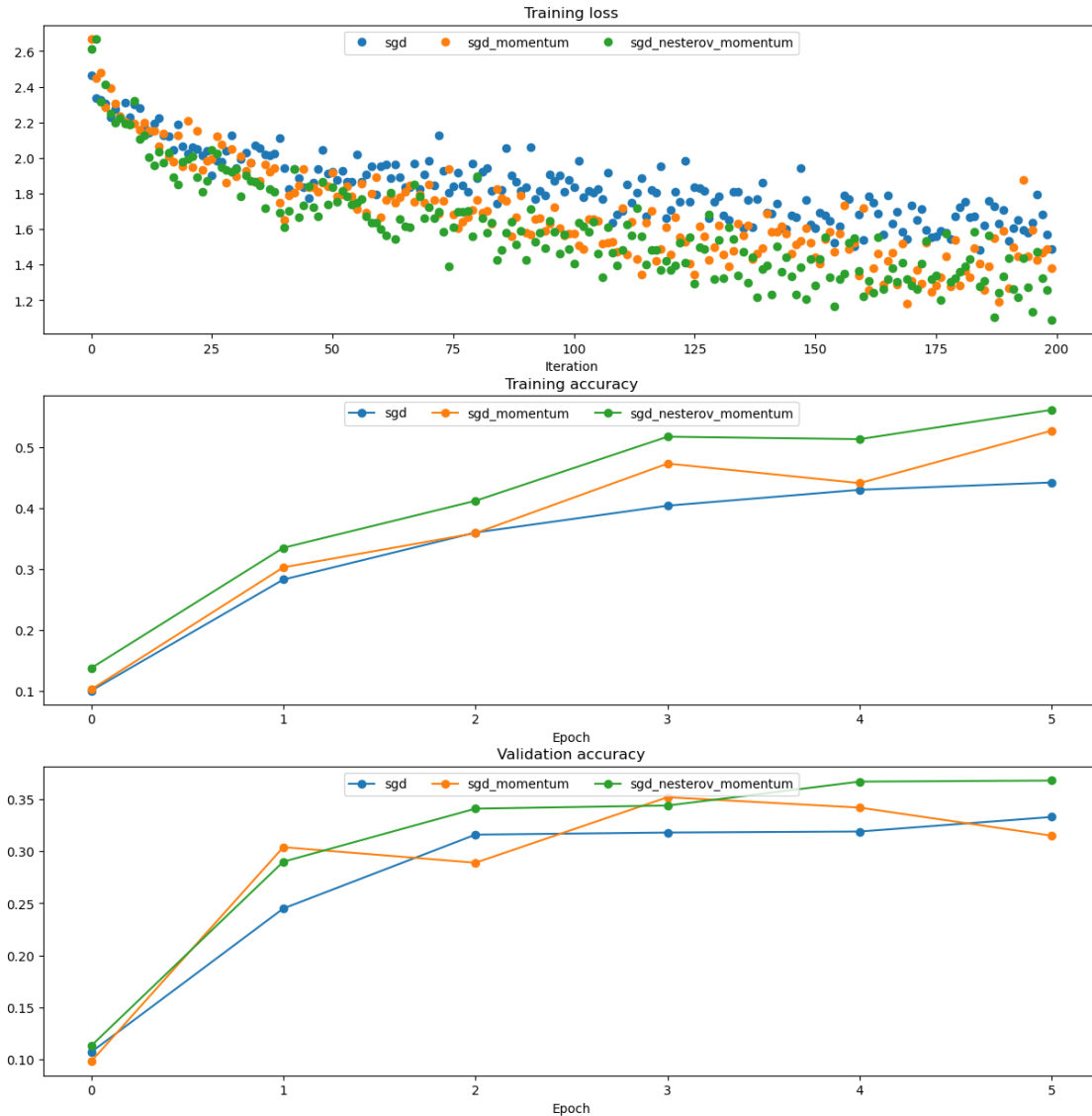
    ax = axes[1]
    ax.plot(solver.train_acc_history, '-o', label=update_rule)

    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.

```
[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.

```

[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.

```

[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)

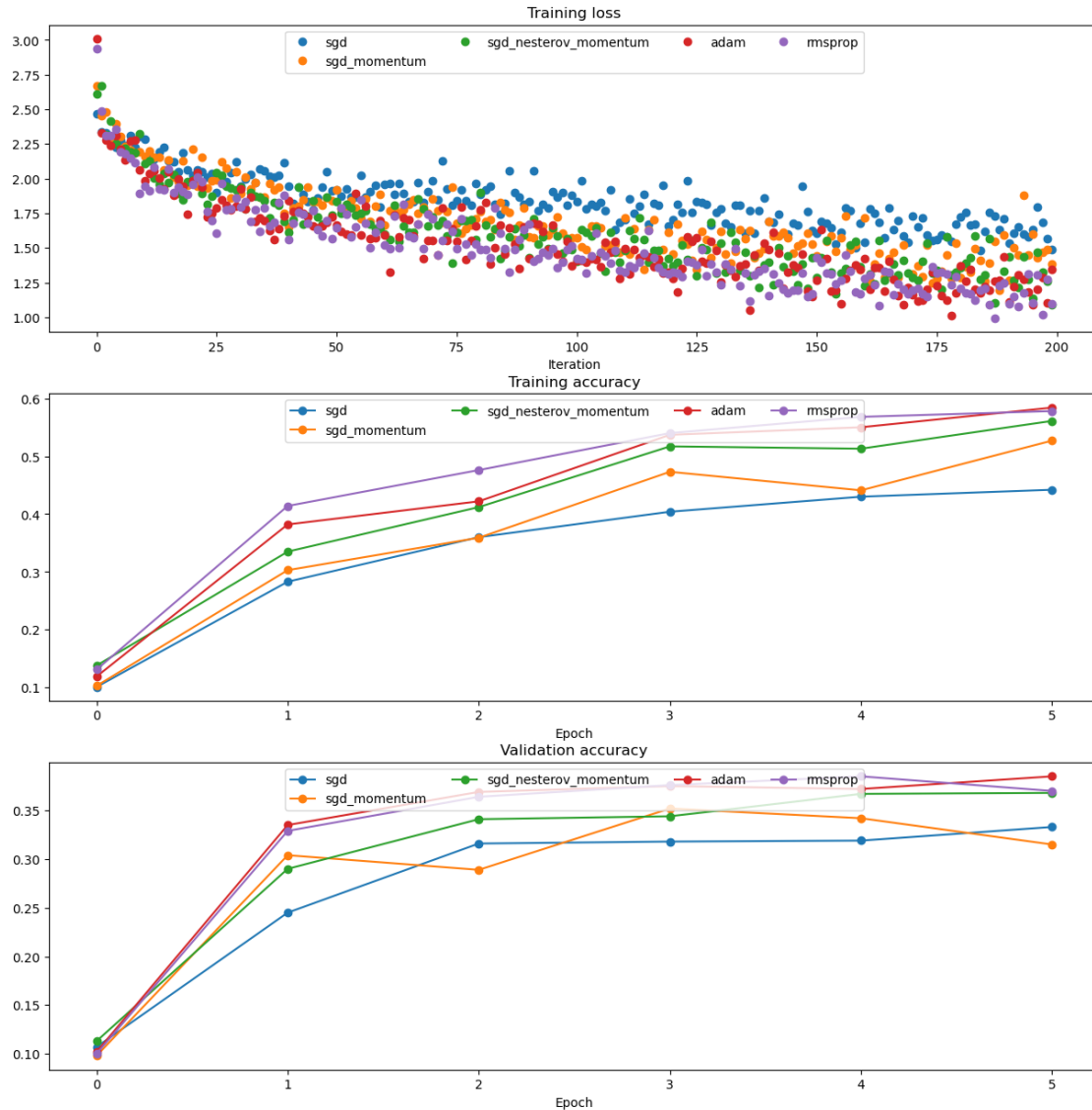
    ax = axes[2]
    ax.plot(solver.val_acc_history, '-o', label=update_rule)

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.

```
[11]: 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=False)

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: 2.283027
(Epoch 0 / 10) train acc: 0.109000; val_acc: 0.086000
(Iteration 51 / 4900) loss: 1.961103
(Iteration 101 / 4900) loss: 1.847143
(Iteration 151 / 4900) loss: 1.758200
(Iteration 201 / 4900) loss: 1.596741
(Iteration 251 / 4900) loss: 1.602392
(Iteration 301 / 4900) loss: 1.735744
(Iteration 351 / 4900) loss: 1.514897
(Iteration 401 / 4900) loss: 1.739453
(Iteration 451 / 4900) loss: 1.566036
(Epoch 1 / 10) train acc: 0.414000; val_acc: 0.450000
(Iteration 501 / 4900) loss: 1.431211
(Iteration 551 / 4900) loss: 1.274565
(Iteration 601 / 4900) loss: 1.249859
(Iteration 651 / 4900) loss: 1.278171
(Iteration 701 / 4900) loss: 1.496156
(Iteration 751 / 4900) loss: 1.510840
(Iteration 801 / 4900) loss: 1.639393
(Iteration 851 / 4900) loss: 1.440012
(Iteration 901 / 4900) loss: 1.682800
(Iteration 951 / 4900) loss: 1.367477
(Epoch 2 / 10) train acc: 0.454000; val_acc: 0.461000
(Iteration 1001 / 4900) loss: 1.453448
(Iteration 1051 / 4900) loss: 1.462146
(Iteration 1101 / 4900) loss: 1.309036
(Iteration 1151 / 4900) loss: 1.158621
(Iteration 1201 / 4900) loss: 1.402987
(Iteration 1251 / 4900) loss: 1.367415
(Iteration 1301 / 4900) loss: 1.385290
(Iteration 1351 / 4900) loss: 1.498346
(Iteration 1401 / 4900) loss: 1.429587
(Iteration 1451 / 4900) loss: 1.510713
(Epoch 3 / 10) train acc: 0.515000; val_acc: 0.492000

```

(Iteration 1501 / 4900) loss: 1.452781
(Iteration 1551 / 4900) loss: 1.339167
(Iteration 1601 / 4900) loss: 1.534219
(Iteration 1651 / 4900) loss: 1.183165
(Iteration 1701 / 4900) loss: 1.288240
(Iteration 1751 / 4900) loss: 1.165124
(Iteration 1801 / 4900) loss: 1.178671
(Iteration 1851 / 4900) loss: 1.404452
(Iteration 1901 / 4900) loss: 1.464458
(Iteration 1951 / 4900) loss: 1.283009
(Epoch 4 / 10) train acc: 0.558000; val_acc: 0.502000
(Iteration 2001 / 4900) loss: 1.388784
(Iteration 2051 / 4900) loss: 1.219097
(Iteration 2101 / 4900) loss: 1.136359
(Iteration 2151 / 4900) loss: 1.217704
(Iteration 2201 / 4900) loss: 1.206258
(Iteration 2251 / 4900) loss: 1.249017
(Iteration 2301 / 4900) loss: 1.107065
(Iteration 2351 / 4900) loss: 1.076624
(Iteration 2401 / 4900) loss: 1.339996
(Epoch 5 / 10) train acc: 0.566000; val_acc: 0.540000
(Iteration 2451 / 4900) loss: 1.159536
(Iteration 2501 / 4900) loss: 1.186217
(Iteration 2551 / 4900) loss: 1.343755
(Iteration 2601 / 4900) loss: 1.315703
(Iteration 2651 / 4900) loss: 1.093518
(Iteration 2701 / 4900) loss: 1.244090
(Iteration 2751 / 4900) loss: 1.094191
(Iteration 2801 / 4900) loss: 1.133828
(Iteration 2851 / 4900) loss: 1.169841
(Iteration 2901 / 4900) loss: 1.291944
(Epoch 6 / 10) train acc: 0.590000; val_acc: 0.523000
(Iteration 2951 / 4900) loss: 1.091541
(Iteration 3001 / 4900) loss: 1.083723
(Iteration 3051 / 4900) loss: 1.057854
(Iteration 3101 / 4900) loss: 1.108180
(Iteration 3151 / 4900) loss: 1.110797
(Iteration 3201 / 4900) loss: 1.129499
(Iteration 3251 / 4900) loss: 1.242609
(Iteration 3301 / 4900) loss: 0.835631
(Iteration 3351 / 4900) loss: 1.024364
(Iteration 3401 / 4900) loss: 1.113593
(Epoch 7 / 10) train acc: 0.632000; val_acc: 0.526000
(Iteration 3451 / 4900) loss: 0.984042
(Iteration 3501 / 4900) loss: 1.023786
(Iteration 3551 / 4900) loss: 1.044522
(Iteration 3601 / 4900) loss: 1.090103
(Iteration 3651 / 4900) loss: 1.169117

```

(Iteration 3701 / 4900) loss: 1.003608
(Iteration 3751 / 4900) loss: 1.203745
(Iteration 3801 / 4900) loss: 0.925290
(Iteration 3851 / 4900) loss: 0.990175
(Iteration 3901 / 4900) loss: 0.957027
(Epoch 8 / 10) train acc: 0.658000; val_acc: 0.519000
(Iteration 3951 / 4900) loss: 0.818857
(Iteration 4001 / 4900) loss: 1.048305
(Iteration 4051 / 4900) loss: 1.012329
(Iteration 4101 / 4900) loss: 1.020251
(Iteration 4151 / 4900) loss: 0.877174
(Iteration 4201 / 4900) loss: 0.962462
(Iteration 4251 / 4900) loss: 1.225517
(Iteration 4301 / 4900) loss: 0.957060
(Iteration 4351 / 4900) loss: 0.907629
(Iteration 4401 / 4900) loss: 0.958417
(Epoch 9 / 10) train acc: 0.660000; val_acc: 0.533000
(Iteration 4451 / 4900) loss: 0.759993
(Iteration 4501 / 4900) loss: 0.918283
(Iteration 4551 / 4900) loss: 0.868897
(Iteration 4601 / 4900) loss: 0.937191
(Iteration 4651 / 4900) loss: 1.010424
(Iteration 4701 / 4900) loss: 0.965550
(Iteration 4751 / 4900) loss: 0.992829
(Iteration 4801 / 4900) loss: 0.964204
(Iteration 4851 / 4900) loss: 0.860855
(Epoch 10 / 10) train acc: 0.693000; val_acc: 0.545000

```

```

[12]: 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'])))

```

```

(3072, 500)
(500, 500)
(500, 500)
(3072, 500)
(500, 500)
(500, 500)
Validation set accuracy: 0.545
Test set accuracy: 0.531

```

```

[ ]:

```

Batch-Normalization

February 16, 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`.

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

import time
import numpy as np
import matplotlib.pyplot as plt
from nndl.fc_net import *
from nndl.layers import *
from utils.data_utils import get_CIFAR10_data
from utils.gradient_check import eval_numerical_gradient, \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))))
```

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

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

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

1.1 Batchnorm forward pass

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

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

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

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

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

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

```
Before batch normalization:
  means: [-10.00457198  16.04505515  19.23080314]
  stds:  [37.41826382 31.96563484 36.48977956]
```


After batch normalization (gamma=1, beta=0)
mean: [-1.60982339e-17 3.61100039e-16 4.03010958e-16]
std: [1. 1. 1.]

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

```
[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: [-0.07583125 -0.08963107 -0.09446156]
stds: [0.94402115 0.96717243 1.05437211]

1.2 Batchnorm backward pass

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

```
[5]: # Gradient check batchnorm backward pass

N, D = 4, 5
```

```

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

bn_param = {'mode': 'train'}
fx = lambda x: batchnorm_forward(x, gamma, beta, bn_param)[0]
fg = lambda 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:  1.0286924753847968e-08
dgamma error:  1.051318155351602e-11
dbeta error:  3.282273415383073e-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 $1e-4$.

```

[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: 2.557708804606235
W1 relative error: 2.439478963329463e-05
W2 relative error: 3.36117016126612e-06
W3 relative error: 1.5816378924460456e-08
b1 relative error: 4.440892098500626e-08
b2 relative error: 4.440892098500626e-08
b3 relative error: 1.2203109853432995e-10
beta1 relative error: 5.576873117898441e-08
beta2 relative error: 3.6237143231206144e-09
gamma1 relative error: 5.3111037572313395e-08
gamma2 relative error: 1.0401735996362805e-08

```

```

Running check with reg = 3.14
Initial loss: 6.93961610307317
W1 relative error: 5.949456134759869e-06
W2 relative error: 0.00030555164665159657
W3 relative error: 8.816515307898475e-08
b1 relative error: 2.220446049250313e-08
b2 relative error: 2.886579864025407e-07
b3 relative error: 3.3990316943449e-10
beta1 relative error: 6.0586315013477285e-09
beta2 relative error: 5.185544391298469e-09
gamma1 relative error: 6.007944811542892e-09
gamma2 relative error: 7.777982374209474e-09

```

1.4 Training a deep fully connected network with batch normalization.

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

```

[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: 2.300539
(Epoch 0 / 10) train acc: 0.132000; val_acc: 0.119000
(Epoch 1 / 10) train acc: 0.346000; val_acc: 0.266000
(Epoch 2 / 10) train acc: 0.434000; val_acc: 0.294000
(Epoch 3 / 10) train acc: 0.517000; val_acc: 0.319000
(Epoch 4 / 10) train acc: 0.557000; val_acc: 0.305000
(Epoch 5 / 10) train acc: 0.614000; val_acc: 0.330000
(Epoch 6 / 10) train acc: 0.661000; val_acc: 0.336000
(Epoch 7 / 10) train acc: 0.723000; val_acc: 0.304000
(Epoch 8 / 10) train acc: 0.753000; val_acc: 0.327000
(Epoch 9 / 10) train acc: 0.784000; val_acc: 0.350000
(Epoch 10 / 10) train acc: 0.787000; val_acc: 0.333000
(Iteration 1 / 200) loss: 2.301798
(Epoch 0 / 10) train acc: 0.170000; val_acc: 0.182000
(Epoch 1 / 10) train acc: 0.256000; val_acc: 0.243000

```

```
(Epoch 2 / 10) train acc: 0.303000; val_acc: 0.251000
(Epoch 3 / 10) train acc: 0.325000; val_acc: 0.264000
(Epoch 4 / 10) train acc: 0.395000; val_acc: 0.285000
(Epoch 5 / 10) train acc: 0.414000; val_acc: 0.299000
(Epoch 6 / 10) train acc: 0.432000; val_acc: 0.295000
(Epoch 7 / 10) train acc: 0.523000; val_acc: 0.293000
(Epoch 8 / 10) train acc: 0.583000; val_acc: 0.314000
(Epoch 9 / 10) train acc: 0.644000; val_acc: 0.326000
(Epoch 10 / 10) train acc: 0.616000; val_acc: 0.321000
```

```
[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')

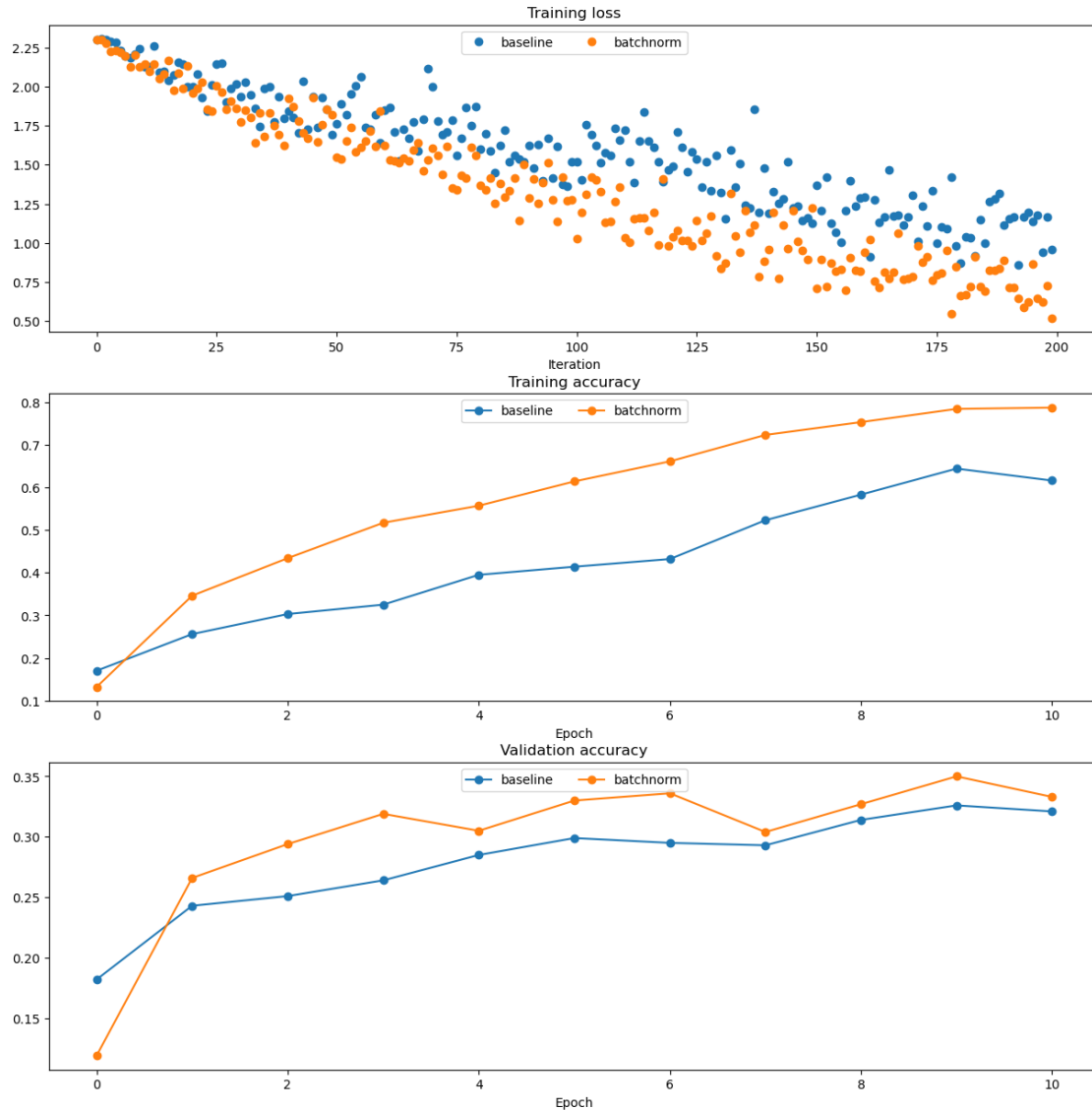
ax = axes[2]
ax.set_title('Validation accuracy')
ax.set_xlabel('Epoch')

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.

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

```

```

[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()

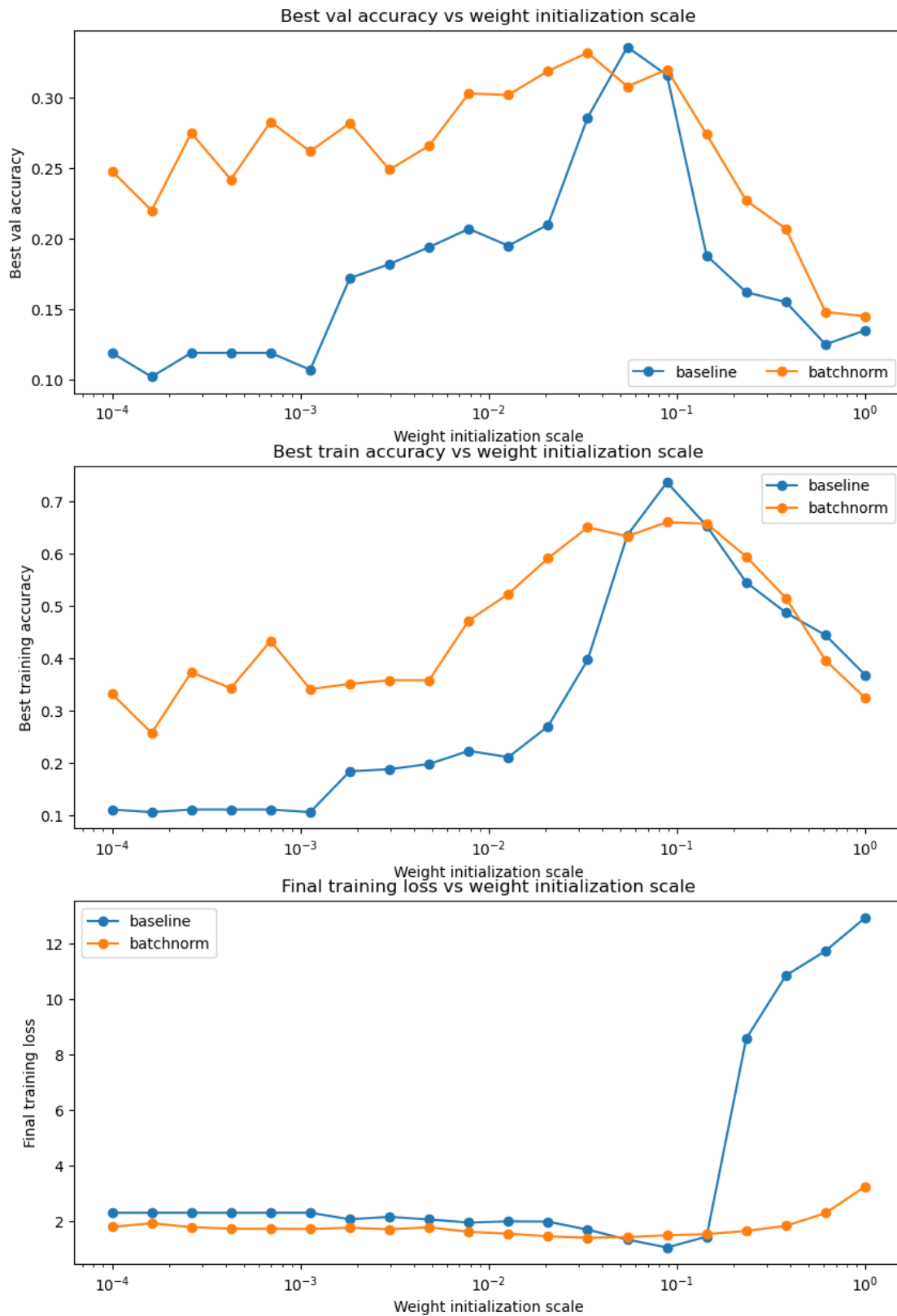
plt.subplot(3, 1, 3)
plt.title('Final training loss vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Final training loss')
plt.semilogx(weight_scales, final_train_loss, '-o', label='baseline')
plt.semilogx(weight_scales, bn_final_train_loss, '-o', label='batchnorm')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

Using a batchnorm in a neural network makes our network more robust to changes in weight initialization. If we do not use batchnorm, layer's activation output can become zero as we proceed to deeper layers(last layer is the deepest layer) depending on the initialization. Therefore, when we do not use batchnorm, neural networks are sensitive to weight initialization.

[]:

Dropout

February 16, 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`.

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

import time
import numpy as np
import matplotlib.pyplot as plt
from nndl.fc_net import *
from nndl.layers import *
from utils.data_utils import get_CIFAR10_data
from utils.gradient_check import eval_numerical_gradient, \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))))
```

```
[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.

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

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

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

```
Running tests with p = 0.3
Mean of input: 10.002160636198084
Mean of train-time output: 10.034498191614126
Mean of test-time output: 10.002160636198084
Fraction of train-time output set to zero: 0.698856
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.6
Mean of input: 10.002160636198084
Mean of train-time output: 10.014191558610936
Mean of test-time output: 10.002160636198084
Fraction of train-time output set to zero: 0.398984
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.75
Mean of input: 10.002160636198084
Mean of train-time output: 9.991626951850632
Mean of test-time output: 10.002160636198084
Fraction of train-time output set to zero: 0.250844
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:

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

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

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

dx relative error: 5.4456129818149784e-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).

```
[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: 7.439246277174223e-08
W3 relative error: 1.910371079177692e-08
b1 relative error: 4.112891126518e-09
b2 relative error: 5.756219669730665e-10
b3 relative error: 1.3204470857080166e-10

```

```

Running check with dropout = 0.75
Initial loss: 2.306133548427975
W1 relative error: 8.729860959912199e-08
W2 relative error: 2.977730805018163e-07
W3 relative error: 1.8832781050769266e-08
b1 relative error: 5.379484812384277e-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: 1.2744094629312605e-06
W2 relative error: 4.678743287158363e-07
W3 relative error: 6.915242131708905e-08
b1 relative error: 4.085353839090859e-08
b2 relative error: 1.9513419724381706e-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.

[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

```

```

[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='%0.2f 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:

```

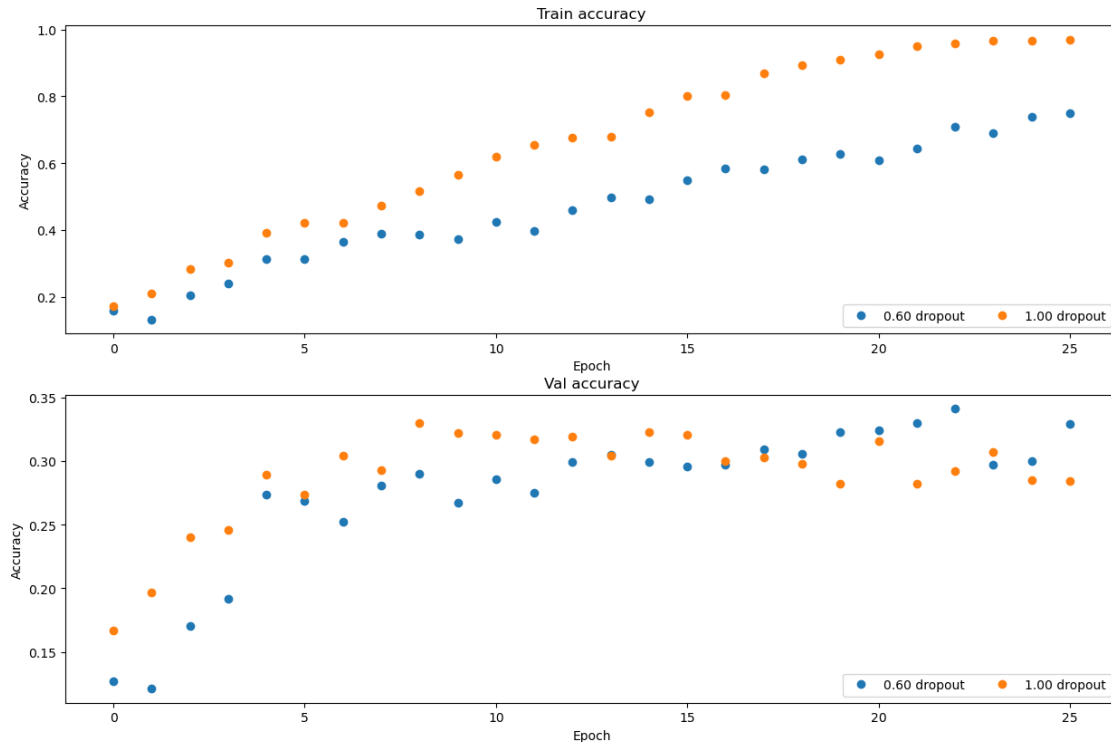


```

plt.plot(solvers[dropout].val_acc_history, 'o', label='%.2f dropout' %L
↪ 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:

Based off the results, it is clear to see that using dropout performs regularization for our neural network. It can be seen from the plot 1 that training accuracy reaches to 1 quicker than plot 2. Both validation accuracies are nearly end up to 0.33. Therefore, we get training accuracy close to 1, whereas we get validation accuracy close 0.33. This indicates that when we do not use dropout model overfits to the data. When we look at the results of dropout used neural network, we get training accuracy close to 0.75, whereas we get validation accuracy close to 0.33. Therefore, we do not overfit like we did in the first case. In result, we can conclude that dropout method can be used for regularization purposes.

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.

```
[17]: # ===== #
# YOUR CODE HERE:
#   Implement a FC-net that achieves at least 55% validation accuracy
#   on CIFAR-10.
# ===== #
hidden_dims = [300,300,300]
weight_scale = 5e-2
dropout = 0.8
epochs = 20
batch_size = 200
optimizer = "adam"
lr = 5e-3
lr_dec = 0.925
reg = 0
best_net = FullyConnectedNet(hidden_dims, weight_scale = weight_scale,
                             dropout=dropout,use_batchnorm = True, reg = reg)

solver = Solver(best_net, data,
                num_epochs=epochs, batch_size=batch_size,
                update_rule='adam',
                optim_config={
                    'learning_rate': lr,
                },
                lr_decay = lr_dec,
                verbose=True, print_every= 100)

solver.train()

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

```
(Iteration 1 / 4900) loss: 2.530187
(Epoch 0 / 20) train acc: 0.164000; val_acc: 0.179000
(Iteration 101 / 4900) loss: 1.665730
(Iteration 201 / 4900) loss: 1.452605
(Epoch 1 / 20) train acc: 0.451000; val_acc: 0.470000
(Iteration 301 / 4900) loss: 1.446262
(Iteration 401 / 4900) loss: 1.347595
(Epoch 2 / 20) train acc: 0.544000; val_acc: 0.492000
(Iteration 501 / 4900) loss: 1.465576
(Iteration 601 / 4900) loss: 1.523093
(Iteration 701 / 4900) loss: 1.306868
(Epoch 3 / 20) train acc: 0.594000; val_acc: 0.531000
```

(Iteration 801 / 4900) loss: 1.418239
(Iteration 901 / 4900) loss: 1.311072
(Epoch 4 / 20) train acc: 0.561000; val_acc: 0.534000
(Iteration 1001 / 4900) loss: 1.282176
(Iteration 1101 / 4900) loss: 1.291150
(Iteration 1201 / 4900) loss: 1.152161
(Epoch 5 / 20) train acc: 0.600000; val_acc: 0.536000
(Iteration 1301 / 4900) loss: 1.213567
(Iteration 1401 / 4900) loss: 1.125476
(Epoch 6 / 20) train acc: 0.647000; val_acc: 0.569000
(Iteration 1501 / 4900) loss: 1.169649
(Iteration 1601 / 4900) loss: 1.096392
(Iteration 1701 / 4900) loss: 1.087360
(Epoch 7 / 20) train acc: 0.628000; val_acc: 0.559000
(Iteration 1801 / 4900) loss: 1.128816
(Iteration 1901 / 4900) loss: 1.165210
(Epoch 8 / 20) train acc: 0.669000; val_acc: 0.558000
(Iteration 2001 / 4900) loss: 1.057295
(Iteration 2101 / 4900) loss: 1.169359
(Iteration 2201 / 4900) loss: 1.048152
(Epoch 9 / 20) train acc: 0.675000; val_acc: 0.569000
(Iteration 2301 / 4900) loss: 1.144551
(Iteration 2401 / 4900) loss: 1.033312
(Epoch 10 / 20) train acc: 0.690000; val_acc: 0.585000
(Iteration 2501 / 4900) loss: 1.019493
(Iteration 2601 / 4900) loss: 0.967580
(Epoch 11 / 20) train acc: 0.724000; val_acc: 0.573000
(Iteration 2701 / 4900) loss: 0.911838
(Iteration 2801 / 4900) loss: 0.983276
(Iteration 2901 / 4900) loss: 0.887045
(Epoch 12 / 20) train acc: 0.738000; val_acc: 0.578000
(Iteration 3001 / 4900) loss: 1.001920
(Iteration 3101 / 4900) loss: 1.081322
(Epoch 13 / 20) train acc: 0.740000; val_acc: 0.586000
(Iteration 3201 / 4900) loss: 0.802459
(Iteration 3301 / 4900) loss: 0.896141
(Iteration 3401 / 4900) loss: 0.959170
(Epoch 14 / 20) train acc: 0.738000; val_acc: 0.588000
(Iteration 3501 / 4900) loss: 0.884596
(Iteration 3601 / 4900) loss: 0.873856
(Epoch 15 / 20) train acc: 0.794000; val_acc: 0.581000
(Iteration 3701 / 4900) loss: 0.815538
(Iteration 3801 / 4900) loss: 0.905308
(Iteration 3901 / 4900) loss: 0.771382
(Epoch 16 / 20) train acc: 0.767000; val_acc: 0.600000
(Iteration 4001 / 4900) loss: 0.793065
(Iteration 4101 / 4900) loss: 0.771485
(Epoch 17 / 20) train acc: 0.796000; val_acc: 0.588000

```
(Iteration 4201 / 4900) loss: 0.771555
(Iteration 4301 / 4900) loss: 0.923525
(Iteration 4401 / 4900) loss: 0.776838
(Epoch 18 / 20) train acc: 0.810000; val_acc: 0.591000
(Iteration 4501 / 4900) loss: 0.740192
(Iteration 4601 / 4900) loss: 0.829750
(Epoch 19 / 20) train acc: 0.833000; val_acc: 0.593000
(Iteration 4701 / 4900) loss: 0.868081
(Iteration 4801 / 4900) loss: 0.775580
(Epoch 20 / 20) train acc: 0.815000; val_acc: 0.593000
```

[]:

```

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

```

```

68     v = config['momentum'] * v - config['learning_rate'] * dw
69     next_w = v + w
70     # ===== #
71     # END YOUR CODE HERE
72     # ===== #
73
74     config['velocity'] = v
75
76     return next_w, config
77
78 def sgd_nesterov_momentum(w, dw, config=None):
79     """
80     Performs stochastic gradient descent with Nesterov momentum.
81
82     config format:
83     - learning_rate: Scalar learning rate.
84     - momentum: Scalar between 0 and 1 giving the momentum value.
85       Setting momentum = 0 reduces to sgd.
86     - velocity: A numpy array of the same shape as w and dw used to store a moving
87       average of the gradients.
88     """
89     if config is None: config = {}
90     config.setdefault('learning_rate', 1e-2)
91     config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
92     v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.
93
94     # ===== #
95     # YOUR CODE HERE:
96     #   Implement the momentum update formula. Return the updated weights
97     #   as next_w, and the updated velocity as v.
98     # ===== #
99     v_prev = v
100    v = config['momentum'] * v_prev - config['learning_rate'] * dw
101    next_w = v + w + config['momentum'] * (v - v_prev)
102    # ===== #
103    # END YOUR CODE HERE
104    # ===== #
105
106    config['velocity'] = v
107
108    return next_w, config
109
110 def rmsprop(w, dw, config=None):
111     """
112     Uses the RMSProp update rule, which uses a moving average of squared gradient
113     values to set adaptive per-parameter learning rates.
114
115     config format:
116     - learning_rate: Scalar learning rate.
117     - decay_rate: Scalar between 0 and 1 giving the decay rate for the squared
118       gradient cache.
119     - epsilon: Small scalar used for smoothing to avoid dividing by zero.
120     - beta: Moving average of second moments of gradients.
121     """
122     if config is None: config = {}
123     config.setdefault('learning_rate', 1e-2)
124     config.setdefault('decay_rate', 0.99)
125     config.setdefault('epsilon', 1e-8)
126     config.setdefault('a', np.zeros_like(w))
127
128     next_w = None
129
130     # ===== #
131     # YOUR CODE HERE:
132     #   Implement RMSProp. Store the next value of w as next_w. You need
133     #   to also store in config['a'] the moving average of the second
134     #   moment gradients, so they can be used for future gradients. Concretely,

```

```

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

```

```

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

```



```

67     x D after reshaping
68     # db should be M; it is just the sum over dout examples
69     # ===== #
70     flattened_x = x.reshape(x.shape[0],-1)
71     dx = np.dot(dout,w.T).reshape(x.shape)
72     dw = np.dot(flattened_x.T,dout)
73     db = np.sum(dout, axis = 0)
74     # ===== #
75     # END YOUR CODE HERE
76     # ===== #
77
78     return dx, dw, db
79
80 def relu_forward(x):
81     """
82     Computes the forward pass for a layer of rectified linear units (ReLU).
83
84     Input:
85     - x: Inputs, of any shape
86
87     Returns a tuple of:
88     - out: Output, of the same shape as x
89     - cache: x
90     """
91     # ===== #
92     # YOUR CODE HERE:
93     # Implement the ReLU forward pass.
94     # ===== #
95     relu = lambda x: x * (x > 0)
96     out = relu(x)
97     # ===== #
98     # END YOUR CODE HERE
99     # ===== #
100
101     cache = x
102     return out, cache
103
104
105 def relu_backward(dout, cache):
106     """
107     Computes the backward pass for a layer of rectified linear units (ReLU).
108
109     Input:
110     - dout: Upstream derivatives, of any shape
111     - cache: Input x, of same shape as dout
112
113     Returns:
114     - dx: Gradient with respect to x
115     """
116     x = cache
117
118     # ===== #
119     # YOUR CODE HERE:
120     # Implement the ReLU backward pass
121     # ===== #
122     dx = dout * (x.reshape(x.shape[0],-1) > 0)
123
124     # ===== #
125     # END YOUR CODE HERE
126     # ===== #
127
128     return dx
129
130 def batchnorm_forward(x, gamma, beta, bn_param):
131     """
132     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 parameter of shape (D,)
- bn_param: Dictionary with the following keys:
 - mode: 'train' or 'test'; required
 - eps: Constant for numeric stability
 - momentum: Constant for running mean / variance.
 - running_mean: Array of shape (D,) giving running mean of features
 - running_var: Array of shape (D,) giving running variance of features

Returns a tuple of:

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

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

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

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

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

```
    # ===== #
    # YOUR CODE HERE:
    #   A few steps here:
    #   (1) Calculate the running mean and variance of the minibatch.
    #   (2) Normalize the activations with the running mean and variance.
    #   (3) Scale and shift the normalized activations. Store this
    #       as the variable 'out'
    #   (4) Store any variables you may need for the backward pass in
    #       the 'cache' variable.
    # ===== #
    mean_x = np.mean(x, axis = 0)
    var_x = np.var(x, axis = 0)

    running_mean = momentum * running_mean + ( 1 - momentum ) * mean_x
    running_var = momentum * running_var + ( 1 - momentum ) * var_x

    standard_x = ( x - mean_x ) / ( np.sqrt(var_x + eps))

    out = gamma * standard_x + beta
    cache = (mean_x, var_x, standard_x, gamma, x, eps)
```

```

200         # ===== #
201         # END YOUR CODE HERE
202         # ===== #
203     elif mode == 'test':
204         # ===== #
205         # YOUR CODE HERE:
206         #     Calculate the testing time normalized activation. Normalize using
207         #     the running mean and variance, and then scale and shift appropriately.
208         #     Store the output as 'out'.
209         # ===== #
210
211         standard_x = ( x - running_mean) / (np.sqrt(running_var))
212         out = gamma * standard_x + beta
213         #cache = []
214         # ===== #
215         # END YOUR CODE HERE
216         # ===== #
217     else:
218         raise ValueError('Invalid forward batchnorm mode "%s"' % mode)
219
220     # Store the updated running means back into bn_param
221     bn_param['running_mean'] = running_mean
222     bn_param['running_var'] = running_var
223
224     return out, cache
225
226 def batchnorm_backward(dout, cache):
227     """
228     Backward pass for batch normalization.
229
230     For this implementation, you should write out a computation graph for
231     batch normalization on paper and propagate gradients backward through
232     intermediate nodes.
233
234     Inputs:
235     - dout: Upstream derivatives, of shape (N, D)
236     - cache: Variable of intermediates from batchnorm_forward.
237
238     Returns a tuple of:
239     - dx: Gradient with respect to inputs x, of shape (N, D)
240     - dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
241     - dbeta: Gradient with respect to shift parameter beta, of shape (D,)
242     """
243     dx, dgamma, dbeta = None, None, None
244
245     # ===== #
246     # YOUR CODE HERE:
247     #     Implement the batchnorm backward pass, calculating dx, dgamma, and dbeta.
248     # ===== #
249     (mean_x, var_x, standard_x, gamma, x, eps) = cache
250     sample_size = x.shape[0]
251     sigma_x = np.sqrt(var_x + eps)
252
253     dgamma = np.sum(standard_x * dout,axis = 0)
254     dbeta = np.sum(dout, axis = 0)
255
256     dL_dx_st = dout * gamma
257
258     dx_st_da = 1 / sigma_x
259     dL_da = dx_st_da * dL_dx_st
260     da_dx = 1
261
262     dx_st_de = -0.5 * ( dx_st_da ** 3) * (x - mean_x)
263     dL_de = dx_st_de * dL_dx_st
264
265     dL_dvar = np.sum(dL_de, axis = 0)
266     dvar_dx = (2 * ( x - mean_x)) / sample_size

```

```

267
268 dL_dmean = np.sum(-dL_da, axis = 0)
269 dmean_dx = 1 / sample_size
270
271 dx = da_dx * dL_da + dvar_dx * dL_dvar + dmean_dx * dL_dmean
272 # ===== #
273 # END YOUR CODE HERE
274 # ===== #
275
276 return dx, dgamma, dbeta
277
278 def dropout_forward(x, dropout_param):
279     """
280     Performs the forward pass for (inverted) dropout.
281
282     Inputs:
283     - x: Input data, of any shape
284     - dropout_param: A dictionary with the following keys:
285       - p: Dropout parameter. We keep each neuron output with probability p.
286       - mode: 'test' or 'train'. If the mode is train, then perform dropout;
287         if the mode is test, then just return the input.
288       - seed: Seed for the random number generator. Passing seed makes this
289         function deterministic, which is needed for gradient checking but not in
290         real networks.
291
292     Outputs:
293     - out: Array of the same shape as x.
294     - cache: A tuple (dropout_param, mask). In training mode, mask is the dropout
295       mask that was used to multiply the input; in test mode, mask is None.
296     """
297     p, mode = dropout_param['p'], dropout_param['mode']
298     assert (0 < p <= 1), "Dropout probability is not in (0,1]"
299     if 'seed' in dropout_param:
300         np.random.seed(dropout_param['seed'])
301
302     mask = None
303     out = None
304
305     if mode == 'train':
306         # ===== #
307         # YOUR CODE HERE:
308         # Implement the inverted dropout forward pass during training time.
309         # Store the masked and scaled activations in out, and store the
310         # dropout mask as the variable mask.
311         # ===== #
312
313         mask = (np.random.rand(x.shape[0], x.shape[1]) < p) / p
314         out = x * mask
315         # ===== #
316         # END YOUR CODE HERE
317         # ===== #
318
319     elif mode == 'test':
320
321         # ===== #
322         # YOUR CODE HERE:
323         # Implement the inverted dropout forward pass during test time.
324         # ===== #
325
326         out = x
327
328         # ===== #
329         # END YOUR CODE HERE
330         # ===== #
331
332     cache = (dropout_param, mask)
333     out = out.astype(x.dtype, copy=False)

```

```

334
335     return out, cache
336
337 def dropout_backward(dout, cache):
338     """
339     Perform the backward pass for (inverted) dropout.
340
341     Inputs:
342     - dout: Upstream derivatives, of any shape
343     - cache: (dropout_param, mask) from dropout_forward.
344     """
345     dropout_param, mask = cache
346     mode = dropout_param['mode']
347
348     dx = None
349     if mode == 'train':
350         # ===== #
351         # YOUR CODE HERE:
352         # Implement the inverted dropout backward pass during training time.
353         # ===== #
354
355         dx = dout * mask
356
357         # ===== #
358         # END YOUR CODE HERE
359         # ===== #
360     elif mode == 'test':
361         # ===== #
362         # YOUR CODE HERE:
363         # Implement the inverted dropout backward pass during test time.
364         # ===== #
365
366         dx = dout
367
368         # ===== #
369         # END YOUR CODE HERE
370         # ===== #
371     return dx
372
373 def svm_loss(x, y):
374     """
375     Computes the loss and gradient using for multiclass SVM classification.
376
377     Inputs:
378     - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
379         for the ith input.
380     - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
381         0 <= y[i] < C
382
383     Returns a tuple of:
384     - loss: Scalar giving the loss
385     - dx: Gradient of the loss with respect to x
386     """
387     N = x.shape[0]
388     correct_class_scores = x[np.arange(N), y]
389     margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
390     margins[np.arange(N), y] = 0
391     loss = np.sum(margins) / N
392     num_pos = np.sum(margins > 0, axis=1)
393     dx = np.zeros_like(x)
394     dx[margins > 0] = 1
395     dx[np.arange(N), y] -= num_pos
396     dx /= N
397     return loss, dx
398
399
400 def softmax_loss(x, y):

```

```

401 """
402 Computes the loss and gradient for softmax classification.
403
404 Inputs:
405 - x: Input data, of shape (N, C) where x[i, j] is the score for the jth class
406     for the ith input.
407 - y: Vector of labels, of shape (N,) where y[i] is the label for x[i] and
408     0 <= y[i] < C
409
410 Returns a tuple of:
411 - loss: Scalar giving the loss
412 - dx: Gradient of the loss with respect to x
413 """
414
415 probs = np.exp(x - np.max(x, axis=1, keepdims=True))
416 probs /= np.sum(probs, axis=1, keepdims=True)
417 N = x.shape[0]
418 loss = -np.sum(np.log(np.maximum(probs[np.arange(N), y], 1e-8))) / N
419 dx = probs.copy()
420 dx[np.arange(N), y] -= 1
421 dx /= N
422 return loss, dx
423

```

```

1  from .layers import *
2
3  def affine_relu_forward(x, w, b):
4      """
5          Convenience layer that performs an affine transform followed by a ReLU
6
7          Inputs:
8          - x: Input to the affine layer
9          - w, b: Weights for the affine layer
10
11         Returns a tuple of:
12         - out: Output from the ReLU
13         - cache: Object to give to the backward pass
14         """
15         a, fc_cache = affine_forward(x, w, b)
16         out, relu_cache = relu_forward(a)
17         cache = (fc_cache, relu_cache)
18         return out, cache
19
20
21 def affine_relu_backward(dout, cache):
22     """
23     Backward pass for the affine-relu convenience layer
24     """
25     fc_cache, relu_cache = cache
26     da = relu_backward(dout, relu_cache)
27     dx, dw, db = affine_backward(da, fc_cache)
28     return dx, dw, db
29
30 def affine_batchnorm_relu_forward(x, w, b, gamma, beta, bn_param):
31     a_out, a_cache = affine_forward(x, w, b)
32     batch_out, batch_cache = batchnorm_forward(a_out, gamma, beta, bn_param)
33     out, relu_cache = relu_forward(batch_out)
34     cache = (a_cache, relu_cache, batch_cache)
35     return out, cache
36
37 def affine_batchnorm_relu_backward(dout, cache):
38     a_cache, relu_cache, batch_cache = cache
39     dbatch = relu_backward(dout, relu_cache)
40     da, dgamma, dbeta = batchnorm_backward(dbatch, batch_cache)
41     dx, dw, db = affine_backward(da, a_cache)
42     return dx, dw, db, dgamma, dbeta

```

```

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

```



```

68     If y is None, then run a test-time forward pass of the model and return:
69     - scores: Array of shape (N, C) giving classification scores, where
70       scores[i, c] is the classification score for X[i] and class c.
71
72     If y is not None, then run a training-time forward and backward pass and
73     return a tuple of:
74     - loss: Scalar value giving the loss
75     - grads: Dictionary with the same keys as self.params, mapping parameter
76       names to gradients of the loss with respect to those parameters.
77     """
78     scores = None
79
80     # ===== #
81     # YOUR CODE HERE:
82     # Implement the forward pass of the two-layer neural network. Store
83     # the class scores as the variable 'scores'. Be sure to use the layers
84     # you prior implemented.
85     # ===== #
86     h, cache_h = affine_relu_forward(X, self.params["W1"], self.params["b1"])
87     scores, cache_scores = affine_forward(h, self.params["W2"], self.params["b2"])
88     # ===== #
89     # END YOUR CODE HERE
90     # ===== #
91
92     # If y is None then we are in test mode so just return scores
93     if y is None:
94         return scores
95
96     loss, grads = 0, {}
97     # ===== #
98     # YOUR CODE HERE:
99     # Implement the backward pass of the two-layer neural net. Store
100    # the loss as the variable 'loss' and store the gradients in the
101    # 'grads' dictionary. For the grads dictionary, grads['W1'] holds
102    # the gradient for W1, grads['b1'] holds the gradient for b1, etc.
103    # i.e., grads[k] holds the gradient for self.params[k].
104    #
105    # Add L2 regularization, where there is an added cost 0.5*self.reg*W^2
106    # for each W. Be sure to include the 0.5 multiplying factor to
107    # match our implementation.
108    #
109    # And be sure to use the layers you prior implemented.
110    # ===== #
111    loss, d_softmax = softmax_loss(scores, y)
112    loss = loss + 0.5 * self.reg * (np.sum(self.params["W1"]**2) + np.sum(self.params
113    ["W2"]**2))
114
115    d_h, d_w2, d_b2 = affine_backward(d_softmax, cache_scores)
116    _, d_w1, d_b1 = affine_relu_backward(d_h, cache_h)
117
118    grads["W1"] = (self.reg * self.params["W1"]) + d_w1
119    grads["b1"] = d_b1
120
121    grads["W2"] = (self.reg * self.params["W2"]) + d_w2
122    grads["b2"] = d_b2
123    # ===== #
124    # END YOUR CODE HERE
125    # ===== #
126
127    return loss, grads
128
129 class FullyConnectedNet(object):
130     """
131     A fully-connected neural network with an arbitrary number of hidden layers,
132     ReLU nonlinearities, and a softmax loss function. This will also implement
133     dropout and batch normalization as options. For a network with L layers,

```

```

134 the architecture will be
135
136 {affine - [batch norm] - relu - [dropout]} x (L - 1) - affine - softmax
137
138 where batch normalization and dropout are optional, and the {...} block is
139 repeated L - 1 times.
140
141 Similar to the TwoLayerNet above, learnable parameters are stored in the
142 self.params dictionary and will be learned using the Solver class.
143 """
144
145 def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,
146             dropout=1, use_batchnorm=False, reg=0.0,
147             weight_scale=1e-2, dtype=np.float32, seed=None):
148     """
149     Initialize a new FullyConnectedNet.
150
151     Inputs:
152     - hidden_dims: A list of integers giving the size of each hidden layer.
153     - input_dim: An integer giving the size of the input.
154     - num_classes: An integer giving the number of classes to classify.
155     - dropout: Scalar between 0 and 1 giving dropout strength. If dropout=1 then
156       the network should not use dropout at all.
157     - use_batchnorm: Whether or not the network should use batch normalization.
158     - reg: Scalar giving L2 regularization strength.
159     - weight_scale: Scalar giving the standard deviation for random
160       initialization of the weights.
161     - dtype: A numpy datatype object; all computations will be performed using
162       this datatype. float32 is faster but less accurate, so you should use
163       float64 for numeric gradient checking.
164     - seed: If not None, then pass this random seed to the dropout layers. This
165       will make the dropout layers deterministic so we can gradient check the
166       model.
167     """
168     self.use_batchnorm = use_batchnorm
169     self.use_dropout = dropout < 1
170     self.reg = reg
171     self.num_layers = 1 + len(hidden_dims)
172     self.dtype = dtype
173     self.params = {}
174
175     # ===== #
176     # YOUR CODE HERE:
177     # Initialize all parameters of the network in the self.params dictionary.
178     # The weights and biases of layer 1 are W1 and b1; and in general the
179     # weights and biases of layer i are Wi and bi. The
180     # biases are initialized to zero and the weights are initialized
181     # so that each parameter has mean 0 and standard deviation weight_scale.
182     #
183     # BATCHNORM: Initialize the gammas of each layer to 1 and the beta
184     # parameters to zero. The gamma and beta parameters for layer 1 should
185     # be self.params['gamma1'] and self.params['beta1']. For layer 2, they
186     # should be gamma2 and beta2, etc. Only use batchnorm if self.use_batchnorm
187     # is true and DO NOT do batch normalize the output scores.
188     # ===== #
189     for i in range(1, self.num_layers + 1):
190         if i == 1:
191             self.params["W" + str(i)] = weight_scale * np.random.randn(input_dim,
192                                 hidden_dims[i - 1])
193             self.params["b" + str(i)] = np.zeros(hidden_dims[i - 1])
194             if self.use_batchnorm:
195                 self.params['gamma' + str(i)] = np.ones(hidden_dims[i-1])
196                 self.params['beta' + str(i)] = np.zeros(hidden_dims[i-1])
197             elif i == self.num_layers:
198                 self.params["W" + str(i)] = weight_scale * np.random.randn(hidden_dims[i
199                                     - 2], num_classes)
200                 self.params["b" + str(i)] = np.zeros(num_classes)

```

```

199         else:
200             self.params["W" + str(i)] = weight_scale * np.random.randn(hidden_dims[i
- 2], hidden_dims[i - 1])
201             self.params["b" + str(i)] = np.zeros(hidden_dims[i - 1])
202             if self.use_batchnorm:
203                 self.params['gamma' + str(i)] = np.ones(hidden_dims[i-1])
204                 self.params['beta' + str(i)] = np.zeros(hidden_dims[i-1])
205             # ===== #
206             # END YOUR CODE HERE
207             # ===== #
208
209             # When using dropout we need to pass a dropout_param dictionary to each
210             # dropout layer so that the layer knows the dropout probability and the mode
211             # (train / test). You can pass the same dropout_param to each dropout layer.
212             self.dropout_param = {}
213             if self.use_dropout:
214                 self.dropout_param = {'mode': 'train', 'p': dropout}
215             if seed is not None:
216                 self.dropout_param['seed'] = seed
217
218             # With batch normalization we need to keep track of running means and
219             # variances, so we need to pass a special bn_param object to each batch
220             # normalization layer. You should pass self.bn_params[0] to the forward pass
221             # of the first batch normalization layer, self.bn_params[1] to the forward
222             # pass of the second batch normalization layer, etc.
223             self.bn_params = []
224             if self.use_batchnorm:
225                 self.bn_params = [{'mode': 'train'} for i in np.arange(self.num_layers - 1)]
226
227             # Cast all parameters to the correct datatype
228             for k, v in self.params.items():
229                 self.params[k] = v.astype(dtype)
230
231
232     def loss(self, X, y=None):
233         """
234         Compute loss and gradient for the fully-connected net.
235
236         Input / output: Same as TwoLayerNet above.
237         """
238         X = X.astype(self.dtype)
239         mode = 'test' if y is None else 'train'
240
241         # Set train/test mode for batchnorm params and dropout param since they
242         # behave differently during training and testing.
243         if self.dropout_param is not None:
244             self.dropout_param['mode'] = mode
245         if self.use_batchnorm:
246             for bn_param in self.bn_params:
247                 bn_param['mode'] = mode
248
249         scores = None
250
251         # ===== #
252         # YOUR CODE HERE:
253         # Implement the forward pass of the FC net and store the output
254         # scores as the variable "scores".
255         #
256         # BATCHNORM: If self.use_batchnorm is true, insert a bathnorm layer
257         # between the affine_forward and relu_forward layers. You may
258         # also write an affine_batchnorm_relu() function in layer_utils.py.
259         #
260         # DROPOUT: If dropout is non-zero, insert a dropout layer after
261         # every ReLU layer.
262         # ===== #
263         cache_h = []
264         cache_dropout = []

```

```

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

```

```
322         else:
323             d_h_tmp, grads["W" + str(i)], grads["b" + str(i)] =
324                 affine_relu_backward(d_h_tmp, cache_h[i - 1])
325
326             grads["W" + str(i)] += self.reg * self.params["W" + str(i)]
327
328
329             # ===== #
330             # END YOUR CODE HERE
331             # ===== #
332
333     return loss, grads
334
```