

ECE C147/247 HW4 Q3: Dropout

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

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

```
In [3]: ## Imports and setups

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

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

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

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

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

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

X_train: (49000, 3, 32, 32)
Y_train: (49000,)
X_val: (1000, 3, 32, 32)
Y_val: (1000,)
X_test: (1000, 3, 32, 32)
Y_test: (1000,)

Dropout forward pass

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

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

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

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

Running tests with p = 0.3
Mean of input: 9.9985010214567
Mean of train-time output: 9.97354144793436
Mean of test-time output: 9.9985010214567
Fraction of train-time output set to zero: 0.301636
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.6
Mean of input: 9.9985010214567
Mean of train-time output: 9.995998684721517
Mean of test-time output: 9.9985010214567
Fraction of train-time output set to zero: 0.599936
Fraction of test-time output set to zero: 0.0
Running tests with p = 0.75
Mean of input: 9.9985010214567
Mean of train-time output: 10.00209523824298
Mean of test-time output: 9.99450214567
Fraction of train-time output set to zero: 0.749776
Fraction of test-time output set to zero: 0.0
```

Dropout backward pass

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

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

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

print('dx relative error: {}'.format(dx_num))

dx relative error: 1.892906988930252e-11
```

Implement a fully connected neural network with dropout layers

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

(1) In the forward pass, you will need to incorporate a dropout layer after every ReLU layer.

(2) In the backward pass, you will need to incorporate a dropout backward pass layer.

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

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

for dropout in [0, 0.25, 0.5]:
    print('Running check with dropout = {}'.format(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: {}'.format(loss))

    for name in sorted(grads):
        f = lambda _: model.loss(X, Y)
        grad_num = eval_numerical_gradient(f, model.params[name], verbose=False, h=1e-5)
        print('({}) relative error: {}'.format(name, rel_error(grad_num, grads[name])))
    print('\n')
```

Running check with dropout = 0
Initial loss: 2.3051948273987857
W1 relative error: 2.3272570382129597e-07
W2 relative error: 1.50348484922161239e-05
W3 relative error: 1.847669394519802e-07
b1 relative error: 2.3369579596391204e-06
b2 relative error: 5.0013939785120424e-08
b3 relative error: 1.1749467839205477e-08

Running check with dropout = 0.25
Initial loss: 2.3052077546540826
W1 relative error: 2.4038469533439284e-07
W2 relative error: 5.02205655497157e-07
W3 relative error: 4.4563160427353564e-08
b1 relative error: 1.30711745000869e-08
b2 relative error: 7.151679404650099e-10
b3 relative error: 1.003974732116764e-10

Running check with dropout = 0.5
Initial loss: 2.303567586595423
W1 relative error: 1.1401257489261862e-06
W2 relative error: 6.5966195253431734e-09
W3 relative error: 6.5966195253431734e-09
b1 relative error: 7.16395914021063e-08
b2 relative error: 1.175510493920166e-09
b3 relative error: 1.4558471033827801e-10

Dropout as a regularizer

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

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

num_train = 500
small_data = {
    'X_train': data['X_train'][:num_train],
    'Y_train': data['Y_train'][:num_train],
    'X_val': data['X_val'],
    'Y_val': data['Y_val'],
}

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

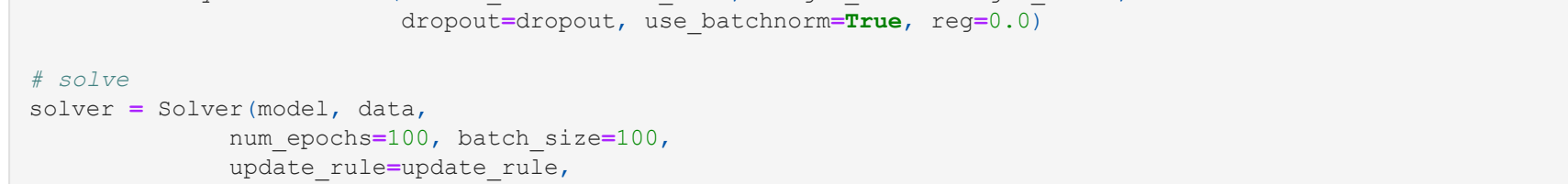
    solver = Solver(model, small_data,
                    num_epochs=25, batch_size=100,
                    update_rule='adam',
                    optim_config={
                        'learning_rate': 5e-4,
                    },
                    verbose=True, print_every=100)
    solver.train()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

(iteration 1 / 125) loss: 2.300804
Epoch 0 / 25) train acc: 0.220000; val_acc: 0.168000
Epoch 1 / 25) train acc: 0.188000; val_acc: 0.170000
Epoch 2 / 25) train acc: 0.266000; val_acc: 0.200000
Epoch 3 / 25) train acc: 0.398000; val_acc: 0.260000
Epoch 4 / 25) train acc: 0.780000; val_acc: 0.278000
Epoch 5 / 25) train acc: 0.428000; val_acc: 0.297000
Epoch 6 / 25) train acc: 0.468000; val_acc: 0.323000
Epoch 7 / 25) train acc: 0.494000; val_acc: 0.287000
Epoch 8 / 25) train acc: 0.566000; val_acc: 0.328000
Epoch 9 / 25) train acc: 0.570000; val_acc: 0.320000
Epoch 10 / 25) train acc: 0.622000; val_acc: 0.324000
Epoch 11 / 25) train acc: 0.670000; val_acc: 0.278000
Epoch 12 / 25) train acc: 0.710000; val_acc: 0.338000
Epoch 13 / 25) train acc: 0.746000; val_acc: 0.318000
Epoch 14 / 25) train acc: 0.792000; val_acc: 0.307000
Epoch 15 / 25) train acc: 0.834000; val_acc: 0.297000
Epoch 16 / 25) train acc: 0.876000; val_acc: 0.327000
Epoch 17 / 25) train acc: 0.886000; val_acc: 0.320000
Epoch 18 / 25) train acc: 0.918000; val_acc: 0.314000
Epoch 19 / 25) train acc: 0.922000; val_acc: 0.296000
Epoch 20 / 25) train acc: 0.948000; val_acc: 0.306000
Epoch 21 / 25) train acc: 0.948000; val_acc: 0.302000
Epoch 22 / 25) train acc: 0.978000; val_acc: 0.302000
Epoch 23 / 25) train acc: 0.976000; val_acc: 0.289000
Epoch 24 / 25) train acc: 0.986000; val_acc: 0.285000
Epoch 25 / 25) train acc: 0.978000; val_acc: 0.311000
(iteration 1 / 125) loss: 2.29816
Epoch 0 / 25) train acc: 0.320000; val_acc: 0.146000
Epoch 1 / 25) train acc: 0.118000; val_acc: 0.131000
Epoch 2 / 25) train acc: 0.220000; val_acc: 0.214000
Epoch 3 / 25) train acc: 0.206000; val_acc: 0.180000
Epoch 4 / 25) train acc: 0.220000; val_acc: 0.193000
Epoch 5 / 25) train acc: 0.264000; val_acc: 0.239000
Epoch 6 / 25) train acc: 0.268000; val_acc: 0.203000
Epoch 7 / 25) train acc: 0.266000; val_acc: 0.212000
Epoch 8 / 25) train acc: 0.282000; val_acc: 0.236000
Epoch 9 / 25) train acc: 0.310000; val_acc: 0.250000
Epoch 10 / 25) train acc: 0.320000; val_acc: 0.267000
Epoch 11 / 25) train acc: 0.338000; val_acc: 0.273000
Epoch 12 / 25) train acc: 0.346000; val_acc: 0.278000
Epoch 13 / 25) train acc: 0.332000; val_acc: 0.273000
Epoch 14 / 25) train acc: 0.328000; val_acc: 0.284000
Epoch 15 / 25) train acc: 0.354000; val_acc: 0.271000
Epoch 16 / 25) train acc: 0.386000; val_acc: 0.277000
Epoch 17 / 25) train acc: 0.388000; val_acc: 0.297000
Epoch 18 / 25) train acc: 0.402000; val_acc: 0.280000
Epoch 19 / 25) train acc: 0.388000; val_acc: 0.274000
Epoch 20 / 25) train acc: 0.386000; val_acc: 0.274000
(iteration 101 / 125) loss: 1.313649
Epoch 21 / 25) train acc: 0.402000; val_acc: 0.272000
Epoch 22 / 25) train acc: 0.440000; val_acc: 0.286000
Epoch 23 / 25) train acc: 0.458000; val_acc: 0.295000
Epoch 24 / 25) train acc: 0.462000; val_acc: 0.310000
Epoch 25 / 25) train acc: 0.446000; val_acc: 0.297000
```

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

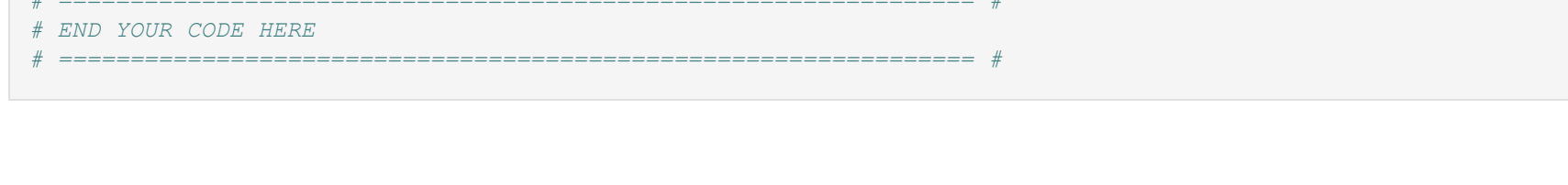
train_accs = []
val_accs = []
for dropout in dropout_choices:
    solver = solvers[dropout]
    train_accs.append(solver.train_acc_history[-1])
    val_accs.append(solver.val_acc_history[-1])

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



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

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



Question

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

Answer:

Yes, the dropout layer does perform regularization. Although the model with and without dropout have similar validation accuracies, the model without dropout (blue) has significantly higher training accuracy compared to the model with dropout (orange). This means that the additional training accuracy the model without dropout has is overfitting and the model with dropout, in fact, regularize it. Also, dropout can be thought of regularizing each hidden unit to work well in many different contexts.

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:

$$\min(\text{floor}((X - 32\%)) / 28\%, 1)$$

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

```
In [9]: # ===== #
# YOUR CODE HERE:
# Implement a FC-net that achieves at least 55% validation accuracy
# on CIFAR-10.
# ===== #

hidden_dims = [600, 600, 600, 600]
learning_rate = 2e-3
weight_scale = 0.01
lr_decay = 0.95
dropout = 0.55
update_rule = 'adam'

# create FullyConnectedNet
model = FullyConnectedNet(hidden_dims=hidden_dims, weight_scale=weight_scale,
                          dropout=dropout, use_batchnorm=True, reg=0.0)

# solve
solver = Solver(model, data,
                num_epochs=100, batch_size=100,
                update_rule=update_rule,
                optim_config={
                    'learning_rate': learning_rate,
                    'lr_decay': lr_decay,
                },
                verbose=True, print_every=100)
solver.train()

# print out the validation accuracy
y_test_pred = np.argmax(model.loss(data['X_test']), axis=1)
y_val_pred = np.argmax(model.loss(data['Y_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'])))

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


Iteration 1 / 49000) loss: 2.368804
train acc: 0.610000; val_acc: 0.167000
(Iteration 101 / 49000) loss: 1.862784
(Iteration 201 / 49000) loss: 1.807409
(Iteration 301 / 49000) loss: 1.755246
(Iteration 401 / 49000) loss: 1.755246
(Epoch 2 / 100) train acc: 0.489000; val_acc: 0.436000
(Iteration 601 / 49000) loss: 1.741667
(Iteration 701 / 49000) loss: 1.707795
(Iteration 801 / 49000) loss: 1.446650
(Iteration 901 / 49000) loss: 1.553775
(Epoch 3 / 100) train acc: 0.490000; val_acc: 0.458000
(Iteration 1001 / 49000) loss: 1.891812
(Iteration 1101 / 49000) loss: 1.532933
(Iteration 1201 / 49000) loss: 1.562373
(Iteration 1301 / 49000) loss: 1.481514
(Iteration 1401 / 49000) loss: 1.563634
(Epoch 3 / 100) train acc: 0.507000; val_acc: 0.488000
(Iteration 1501 / 49000) loss: 1.463753
(Iteration 1601 / 49000) loss: 1.453057
(Iteration 1701 / 49000) loss: 1.414905
(Iteration 1801 / 49000) loss: 1.456929
(Iteration 1901 / 49000) loss: 1.476521
(Epoch 4 / 100) train acc: 0.511000; val_acc: 0.497000
(Iteration 2001 / 49000) loss: 1.435261
(Iteration 2101 / 49000) loss: 1.435045
(Iteration 2201 / 49000) loss: 1.594414
(Iteration 2301 / 49000) loss: 1.454643
(Iteration 2401 / 49000) loss: 1.424646
(Iteration 2501 / 49000) loss: 1.502963; val_acc: 0.514000
(Iteration 2601 / 49000) loss: 1.365315
(Iteration 2701 / 49000) loss: 1.543499
(Iteration 2801 / 49000) loss: 1.629832
(Iteration 2901 / 49000) loss: 1.533493
(Epoch 5 / 100) train acc: 0.523000; val_acc: 0.524000
(Iteration 3001 / 49000) loss: 1.651523
(Iteration 3101 / 49000) loss: 1.303552
(Iteration 3201 / 49000) loss: 1.321157
(Iteration 3301 / 49000) loss: 1.326001
(Iteration 3401 / 49000) loss: 1.405216
(Epoch 5 / 100) train acc: 0.566000; val_acc: 0.524000
(Iteration 3501 / 49000) loss: 1.351739
(Iteration 3601 / 49000) loss: 1.660494
(Iteration 3701 / 49000) loss: 1.394698
(Iteration 3801 / 49000) loss: 1.402928
(Iteration 3901 / 49000) loss: 1.292914
(Epoch 6 / 100) train acc: 0.553000; val_acc: 0.539000
(Iteration 4001 / 49000) loss: 1.430378
(Iteration 4101 / 49000) loss: 1.417105
(Iteration 4201 / 49000) loss: 1.382712
(Iteration 4301 / 49000) loss: 1.368239
(Iteration 4401 / 49000) loss: 1.436760
(Epoch 6 / 100) train acc: 0.627000; val_acc: 0.537000
(Iteration 4501 / 49000) loss: 1.305237
(Iteration 4601 / 49000) loss: 1.323998
(Iteration 4701 / 49000) loss: 1.114894
(Iteration 4801 / 49000) loss: 1.219810
(Epoch 7 / 100) train acc: 0.575000; val_acc: 0.543000
(Iteration 4901 / 49000) loss: 1.393963
(Iteration 5001 / 49000) loss: 1.321516
(Iteration 5101 / 49000) loss: 1.302963
(Iteration 5201 / 49000) loss: 1.327811
(Iteration 5301 / 49000) loss: 1.231963
(Epoch 7 / 100) train acc: 0.620000; val_acc: 0.549000
(Iteration 5401 / 49000) loss: 1.380101
(Iteration 5501 / 49000) loss: 1.201157
(Iteration 5601 / 49000) loss: 1.035552
(Iteration 5701 / 49000) loss: 1.353704
(Epoch 8 / 100) train acc: 0.593000; val_acc: 0.539000
(Iteration 5801 / 49000) loss: 1.211575
(Iteration 5901 / 49000) loss: 1.817611
(Iteration 6001 / 49000) loss: 1.416619
(Iteration 6101 / 49000) loss: 1.416619
(Iteration 6201 / 49000) loss: 1.247955
(Iteration 6301 / 49000) loss: 1.516268
(Epoch 8 / 100) train acc: 0.607000; val_acc: 0.542000
(Iteration 6401 / 49000) loss: 1.058277
(Iteration 6501 / 49000) loss: 1.306083
(Iteration 6601 / 49000) loss: 1.263972
(Iteration 6701 / 49000) loss: 1.784866
(Iteration 6801 / 49000) loss: 1.301901
(Epoch 9 / 100) train acc: 0.598000; val_acc: 0.560000
(Iteration 6901 / 49000) loss: 1.190909
(Iteration 7001 / 49000) loss: 1.060947
(Iteration 7101 / 49000) loss: 1.060947
(Iteration 7201 / 49000) loss: 1.264882
(Iteration 7301 / 49000) loss: 1.417238
(Epoch 9 / 100) train acc: 0.627000; val_acc: 0.557000
(Iteration 7401 / 49000) loss: 1.127628
(Iteration 7501 / 49000) loss: 1.223679
(Iteration 7601 / 49000) loss: 1.302084
(Iteration 7701 / 49000) loss: 1.223759
(Iteration 7801 / 49000) loss: 1.063215
(Epoch 10 / 100) train acc: 0.626000; val_acc: 0.552000
(Iteration 7901 / 49000) loss: 1.247635
(Iteration 8001 / 49000) loss: 1.050236
(Iteration 8101 / 49000) loss: 0.995238
(Iteration 8201 / 49000) loss: 0.995238
(Iteration 8301 / 49000) loss: 1.235742
(Epoch 10 / 100) train acc: 0.641000; val_acc: 0.571000
(Iteration 8401 / 49000) loss: 1.058277
(Iteration 8501 / 49000) loss: 1.204427
(Iteration 8601 / 49000) loss: 1.075615
(Iteration 8701 / 49000) loss: 1.516268
(Iteration 8801 / 49000) loss: 1.217487
(Epoch 11 / 100) train acc: 0.674000; val_acc: 0.564000
(Iteration 8901 / 49000) loss: 1.307140
(Iteration 9001 / 49000) loss: 1.055895
(Iteration 9101 / 49000) loss: 1.055895
(Iteration 9201 / 49000) loss: 1.181401
(Iteration 9301 / 49000) loss: 1.010250
(Epoch 11 / 100) train acc: 0.650000; val_acc: 0.570000
(Iteration 9401 / 49000) loss: 1.164261
(Iteration 9501 / 49000) loss: 1.243444
(Iteration 9601 / 49000) loss: 1.068270
(Iteration 9701 / 49000) loss: 0.994167; val_acc: 0.571000
(Epoch 12 / 100) train acc: 0.665000; val_acc: 0.571000
(Iteration 9801 / 49000) loss: 1.331903
(Iteration 9901 / 49000) loss: 1.118267
(Iteration 10001 / 49000) loss: 1.172186
(Epoch 12 / 100) train acc: 0.646000; val_acc: 0.563000
(Iteration 10101 / 49000) loss: 1.291211
(Iteration 10201 / 49000) loss: 1.069844
(Iteration 10301 / 49000) loss: 1.030236
(Iteration 10401 / 49000) loss: 0.939506
(Iteration 10501 / 49000) loss: 1.282977
(Epoch 12 / 100) train acc: 0.673000; val_acc: 0.575000
(Iteration 10601 / 49000) loss: 1.071739
(Iteration 10701 / 49000) loss: 1.163702
(Iteration 10801 / 49000) loss: 1.026890
(Iteration 10901 / 49000) loss: 1.056204
(Iteration 11001 / 49000) loss: 1.041666
(Epoch 13 / 100) train acc: 0.670000; val_acc: 0.576000
(Iteration 11101 / 49000) loss: 0.980277
(Iteration 11201 / 49000) loss: 1.024924
(Iteration 11301 / 49000) loss: 0.970555
(Iteration 11401 / 49000) loss: 1.035091
(Iteration 11501 / 49000) loss: 1.153221
(Epoch 13 / 100) train acc: 0.665000; val_acc: 0.572000
(Iteration 11601 / 49000) loss: 1.030236
(Iteration 11701 / 49000) loss: 1.166782
(Iteration 11801 / 49000) loss: 1.157622
(Iteration 11901 / 49000) loss: 0.965574
(Epoch 13 / 100) train acc: 0.697000; val_acc: 0.561000
(Iteration 12001 / 49000) loss: 1.010246
(Iteration 12101 / 49000) loss: 1.002679
(Iteration 12201 / 49000) loss: 1.069844
(Iteration 12301 / 49000) loss: 1.169410
(Iteration 12401 / 49000) loss: 1.136060
(Iteration 12501 / 49000) loss: 0.954982
(Epoch 14 / 100) train acc: 0.697000; val_acc: 0.561000
(Iteration 12601 / 49000) loss: 0.989076
(Iteration 12701 / 49000) loss: 1.062014
(Iteration 12801 / 49000) loss: 0.971786
(Iteration 12901 / 49000) loss: 1.056204
(Iteration 13001 / 49000) loss: 1.105684
(Epoch 14 / 100) train acc: 0.690000; val_acc: 0.569000
(Iteration 13101 / 49000) loss: 1.058277
(Iteration 13201 / 49000) loss: 1.059711
(Iteration 13301 / 49000) loss: 1.174263
(Iteration 13401 / 49000) loss: 1.094023
(Iteration 13501 / 49000) loss: 1.054292
(Epoch 14 / 100) train acc: 0.710000; val_acc: 0.572000
(Iteration 13601 / 49000) loss: 1.114179
(Iteration 13701 / 49000) loss: 1.053979
(Iteration 13801 / 49000) loss: 1.257944
(Iteration 13901 / 49000) loss: 1.086937
(Iteration 14001 / 49000) loss: 1.113557
(Epoch 15 / 100) train acc: 0.740000; val_acc: 0.561000
(Iteration 14101 / 49000) loss: 0.946266
(Iteration 14201 / 49000) loss: 0.964021
(Iteration 14301 / 49000) loss: 1.047627
(Iteration 14401 / 49000) loss: 1.037640
(Iteration 14501 / 49000) loss: 0.938981
(Iteration 14601 / 49000) loss: 1.200740
(Epoch 15 / 100) train acc: 0.730000; val_acc: 0.573000
(Iteration 14701 / 49000) loss: 1.037305
(Iteration 14801 / 49000) loss: 1.272926
(Iteration 14901 / 49000) loss: 0.902014
(Iteration 15001 / 49000) loss: 1.175724
(Iteration 15101 / 49000) loss: 0.953093
(Epoch 15 / 100) train acc: 0.690000; val_acc: 0.570000
(Iteration 15201 / 49000) loss: 0.997450
(Iteration 15301 / 49000) loss: 0.990114
(Iteration 15401 / 49000) loss: 1.151598
(Iteration 15501 / 49000) loss: 1.018678
(Iteration 15601 / 49000) loss: 0.954982
(Epoch 16 / 100) train acc: 0.734000; val_acc: 0.569000
(Iteration 15701 / 49000) loss: 0.964021
(Iteration 15801 / 49000) loss: 1.041394
(Iteration 15901 / 49000) loss: 1.104018
(Iteration 16001 / 49000) loss: 0.973629
(Iteration 16101 / 49000) loss: 1.168481
(Epoch 16 / 100) train acc: 0.712000; val_acc: 0.575000
(Iteration 16201 / 49000) loss: 0.956470
(Iteration 16301 / 49000) loss: 1.182668
(Iteration 16401 / 49000) loss: 0.964021
(Iteration 16501 / 49000) loss: 1.102167
(Iteration 16601 / 49000) loss: 1.113402
(Epoch 16 / 100) train acc: 0.740000; val_acc: 0.574000
(Iteration 16701 / 49000) loss: 0.987856
(Iteration 16801 / 49000) loss: 1.012500
(Iteration 16901 / 49000) loss: 1.237228
(Iteration 17001 / 49000) loss: 1.072272
(Iteration 17101 / 49000) loss: 0.995394
(Epoch 16 / 100) train acc: 0.713000; val_acc: 0.578000
(Iteration 17201 / 49000) loss: 0.975236
(Iteration 17301 / 49000) loss: 0.996633
(Iteration 17401 / 49000) loss: 1.102627
(Iteration 17501 / 49000) loss: 1.117850
(Iteration 17601 / 49000) loss: 1.139715
(Epoch 17 / 100) train acc: 0.702000; val_acc: 0.569000
(Iteration 17701 / 49000) loss: 1.058277
(Iteration 17801 / 49000) loss: 1.030979
(Iteration 17901 / 49000) loss: 1.035874
(Iteration 18001 / 49000) loss: 0.953979
(Iteration 18101 / 49000) loss: 0.959946
(Epoch 17 / 100) train acc: 0.719000; val_acc: 0.568000
(Iteration 18201 / 49000) loss: 1.050895
(Iteration 18301 / 49000) loss: 0.887651
(Iteration 18401 / 49000) loss: 1.120229
(Iteration 18501 / 49000) loss: 1.087794
(Iteration 18601 / 49000) loss: 0.935568
(Epoch 17 / 100) train acc: 0.710000; val_acc: 0.564000
(Iteration 18701 / 49000) loss: 1.130590
(Iteration 18801 / 49000) loss: 1.035197
(Iteration 18901 / 49000) loss: 0.916587
(Iteration 19001 / 49000) loss: 1.075615
(Epoch 18 / 100) train acc: 0.759000; val_acc: 0.577000
(Iteration 19101 / 49000) loss: 0.930315
(Iteration 19201 / 49000) loss: 1.287657
(Iteration 19301 / 49000) loss: 0.977760
(Epoch 18 / 100) train acc: 0.740000; val_acc: 0.570000
(Iteration 19401 / 49000) loss: 0.934075
(Iteration 19501 / 49000) loss: 0.940475
(Iteration 19601 / 49000) loss: 0.864108
(Iteration 19701 / 49000) loss: 0.993289
(Iteration 19801 / 49000) loss: 0.988361
(Iteration 19901 / 49000) loss: 0.873222
(Iteration 20001 / 49000) loss: 0.950216
(Epoch 19 / 100) train acc: 0.757000; val_acc: 0.575000
(Iteration 20101 / 49000) loss: 0.940822
(Iteration 20201 / 49000) loss: 0.963686
(Iteration 20301 / 49000) loss: 1.298711
(Iteration 20401 / 49000) loss: 0.953220
(Iteration 20501 / 49000) loss: 0.973895
(Epoch 19 / 100) train acc: 0.723000; val_acc: 0.568000
(Iteration 20601 / 49000) loss: 1.042928
(Iteration 20701 / 49000) loss: 1.039326
(Iteration 20801 / 49000) loss: 1.031932
(Iteration 20901 / 49000) loss: 1.050236
(Epoch 19 / 100) train acc: 0.741000; val_acc: 0.591000
(Iteration 21001 / 49000) loss: 1.091723
(Iteration 21101 / 49000) loss: 0.972820
(Iteration 21201 / 49000) loss: 0.996474
(Iteration 21301 / 49000) loss: 0.909890
(Epoch 19 / 100) train acc: 0.750000; val_acc: 0.588000
(Iteration 21401 / 49000) loss: 1.004290
(Iteration 21501 / 49000) loss: 1.069518
(Iteration 21601 / 49000) loss: 0.977466
(Iteration 21701 / 49000) loss: 1.054688
(Iteration 21801 / 49000) loss: 0.923615
(Epoch 20 / 100) train acc: 0.772000; val_acc: 0.588000
(Iteration 21901 / 49000) loss: 0.955362
(Iteration 22001 / 49000) loss: 1.040334
(Iteration 22101 / 49000) loss: 0.963634
(Iteration 22201 / 49000) loss: 0.955362
(Iteration 22301 / 49000) loss: 0.955362
(Iteration 22401 / 49000) loss: 0.962926
(Epoch 20 / 100) train acc: 0.740000; val_acc: 0.592000
(Iteration 22501 / 49000) loss: 0.950236
(Iteration 22601 / 49000) loss: 1.030281
(Iteration 22701 / 49000) loss: 1.030281
(Iteration 22801 / 49000) loss: 0.933981
(Iteration 22901 / 49000) loss: 0.933981
(Epoch 20 / 100) train acc: 0.752000; val_acc: 0.561000
(Iteration 23001 / 49000) loss: 0.997412
(Iteration 23101 / 49000) loss: 1.040023
(Iteration 23201 / 49000) loss: 0.904662
(Iteration 23301 / 49000) loss: 0.938277
(Iteration 23401 / 49000) loss: 0.943500
(Epoch 20 / 100) train acc: 0.763000; val_acc: 0.581000
(Iteration 23501 / 49000) loss: 1.061302
(Iteration 23601 / 49000) loss: 0.980078
(Iteration 23701 / 49000) loss: 0.956021
(Iteration 23801 / 49000) loss: 1.124112
(Iteration 23901 / 49000) loss: 0.984222
(Iteration 24001 / 49000) loss: 0.939633
(Epoch 21 / 100) train acc: 0.773000; val_acc: 0.585000
(Iteration 24101 / 49000) loss: 0.999506
(Iteration 24201 / 49000) loss: 0.955090
(Iteration 24301 / 49000) loss: 0.961211
(Epoch 21 / 100) train acc: 0.743000; val_acc: 0.578000
(Iteration 24401 / 49000) loss: 1.068633
(Iteration 24501 / 49000) loss: 1.003164
(Iteration 24601 / 49000) loss: 0.937819
(Iteration 24701 / 49000) loss: 1.090954
(Iteration 24801 / 49000) loss: 0.918728
(Iteration 24901 / 49000) loss: 0.930378
(Epoch 21 / 100) train acc: 0.758000; val_acc: 0.568000
(Iteration 25001 / 49000) loss: 0.918385
(Iteration 25101 / 49000) loss: 0.912790
(Iteration 25201 / 49000) loss: 1.038813
(Iteration 25301 / 49000) loss: 0.983297
(Iteration 25401 / 49000) loss: 0.796176
(Epoch 21 / 100) train acc: 0.752000; val_acc: 0.585000
(Iteration 25501 / 49000) loss: 1.144894
(Iteration 25601 / 49000) loss: 0.944377
(Iteration 25701 / 49000) loss: 1.114894
(Iteration 25801 / 49000) loss: 0.990313
(Iteration 25901 / 49000) loss: 0.915563
(Epoch 21 / 100) train acc: 0.770000; val_acc: 0.577000
(Iteration 26001 / 49000) loss: 1.039719
(Iteration 26101 / 49000) loss: 1.076385
(Iteration 26201 / 49000) loss: 0.966788
(Iteration 26301 / 49000) loss: 0.930168
(Epoch 22 / 100) train acc: 0.772000; val_acc: 0.581000
(Iteration 26401 / 49000) loss: 1.050236
(Iteration 26501 / 49000) loss: 1.025582
(Iteration 26601 / 49000) loss: 1.120333
(Iteration 26701 / 49000) loss: 0.921766
(Iteration 26801 / 49000) loss: 0.977084
(Iteration 26901 / 49000) loss: 0.977084
(Epoch 22 / 100) train acc: 0.746000; val_acc: 0.583000
(Iteration 27001 / 49000) loss: 0.926161
(Iteration 27101 / 49000) loss: 0.910321
(Iteration 27201 / 49000) loss: 0.772592
(Iteration 27301 / 49000) loss: 0.919287
(Iteration 27401 / 49000) loss: 0.923547
(Epoch 22 / 100) train acc: 0.746000; val_acc: 0.576000
(Iteration 27501 / 49000) loss: 0.984666
(Iteration 27601 / 49000) loss: 0.927114
(Iteration 27701 / 49000) loss: 0.960834
(Iteration 27801 / 49000) loss: 1.1156081
(Iteration 27901 / 49000) loss: 0.960951
(Epoch 22 / 100) train acc: 0.750000; val_acc: 0.581000
(Iteration 28001 / 49000) loss: 0.937897
(Iteration 28101 / 49000) loss: 0.909295
(Iteration 28201 / 49000) loss: 0.960213
(Iteration 28301 / 49000) loss: 0.853147
(Epoch 22 / 100) train acc: 0.745000; val_acc: 0.579000
(Iteration 28401 / 49000) loss: 0.928216
(Iteration 28501 / 49000) loss: 0.966800
(Iteration 28601 / 49000) loss: 0.973815
(Iteration 28701 / 49000) loss: 1.003696
(Iteration 28801 / 49000) loss: 0.988931
(Iteration 28901 / 49000) loss: 0.946939
(Epoch 23 / 100) train acc: 0.777000; val_acc: 0.584000
(Iteration 29001 / 49000) loss: 0.965242
(Iteration 29101 / 49000) loss: 1.001254
(Iteration 29201 / 49000) loss: 0.932011
(Epoch 23 / 100) train acc: 0.773000; val_acc: 0.586000
(Iteration 29301 / 49000) loss: 0.943500
(Iteration 29401 / 49000) loss: 0.877197
(Iteration 29501 / 49000) loss: 0.950236
(Iteration 29601 / 49000) loss: 0.950236
(Iteration 29701 / 49000) loss: 0.950236
(Iteration 29801 / 49000) loss: 0.898424
(Epoch 23 / 100) train acc: 0.769000; val_acc: 0.586000
(Iteration 29901 / 49000) loss: 0.917788
(Iteration 30001 / 49000) loss: 1.029092
(Iteration 30101 / 49000) loss: 0.826729
(Iteration 30201 / 49000) loss: 0.877080
(Iteration 30301 / 49000) loss: 0.888737
(Epoch 23 / 100) train acc: 0.741000; val_acc: 0.569000
(Iteration 30401 / 49000) loss: 1.044946
(Iteration 30501 / 49000) loss: 1.008125
(Iteration 30601 / 49000) loss: 1.079879
(Iteration 30701 / 49000) loss: 0.938385
(Iteration 30801 / 49000) loss: 1.110585
(Epoch 23 / 100) train acc: 0.760000; val_acc: 0.583000
(Iteration 30901 / 49000) loss: 1.069518
(Iteration 31001 / 49000) loss: 0.907207
(Iteration 31101 / 49000) loss: 0.972556
(Iteration 31201 / 49000) loss: 1.061106
(Iteration 31301 / 49000) loss: 0.932104
(Epoch 24 / 100) train acc: 0.759000; val_acc: 0.587000
(Iteration 31401 / 49000) loss: 0.984279
(Iteration 31501 / 49000) loss: 1.069518
(Iteration 31601 / 49000) loss: 0.886769
(Iteration 31701 / 49000) loss: 0.702446
(Iteration 31801 / 49000) loss: 0.938222
(Epoch 24 / 100) train acc: 0.777000; val_acc: 0.591000
(Iteration 31901 / 49000) loss: 0.966800
(Iteration 32001 / 49000) loss: 0.698999
(Iteration 32101 / 49000) loss: 1.005880
(Iteration 32201 / 49000) loss: 0.976777
(Iteration 32301 / 49000) loss: 1.008020
(Epoch 24 / 100) train acc: 0.753000; val_acc: 0.567000
(Iteration 32401 / 49000) loss: 0.936593
(Iteration 32501 / 49000) loss: 1.028871
(Iteration 32601 / 49000) loss: 0.953634
(Iteration 32701 / 49000) loss: 0.924877
(Iteration 32801 / 49000) loss: 0.997265
(Epoch 24 / 100) train acc: 0.771000; val_acc: 0.589000
(Iteration 32901 / 49000) loss: 0.916587
(Iteration 33001 / 49000) loss: 0.926986
(Iteration 33101 / 49000) loss: 0.916362
(Iteration 33201 / 49000) loss: 1.008122
(Iteration 33301 / 49000) loss: 0.916862
(Epoch 24 / 100) train acc: 0.771000; val_acc: 0.585000
(Iteration 33401 / 49000) loss: 1.410132
(Iteration 33501 / 49000) loss: 0.849892
(Iteration 33601 / 49000) loss: 0.777145
(Iteration 33701 / 49000) loss: 0.973194
(Iteration 33801 / 49000) loss: 0.689751
(Epoch 24 / 100) train acc: 0.765000; val_acc: 0.585000
(Iteration 33901 / 49000) loss: 0.962608
(Iteration 34001 / 49000


```

import numpy as np

def affine_forward(x, w, b):
    """
    Computes the forward pass for an affine (fully-connected) layer.

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

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

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

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

    x_res = x.reshape((x.shape[0], w.shape[0])) # Shape of N * D
    out = x_res.dot(w) + b.reshape((1, b.shape[0])) # Shape of N * M

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

    cache = (x, w, b, out)
    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:
    #   - 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 is just the sum over dout examples
    #   - db should be M; it is just the sum over dout examples
    # ===== #

    x_res = np.reshape(x, (x.shape[0], w.shape[0]))
    dx_res = np.dot(dout, w.T)
    dx = np.reshape(dx_res, x.shape) # Shape of N * D
    dw = x_res.T.dot(dout) # Shape of D * M
    db = dout.T.dot(np.ones(x.shape[0])) # Shape of M * 1

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

    return dx, dw, db

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

    Input:
    - x: Inputs, of any shape

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

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

    out = np.maximum(0, x)

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

    cache = (x, out)
    return out, cache

def relu_backward(dout, cache):
    """
    Computes the backward pass for a layer of rectified linear units (ReLUs).

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

    Returns:
    - dx: Gradient with respect to x
    - cache: None

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

    dx = (x > 0) * (dout)

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

    return dx

def batchnorm_forward(x, gamma, beta, bn_param):
    """
    Forward pass for batch normalization.

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

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

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

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

    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_minibatch = np.mean(x, axis=0)
        var_minibatch = np.var(x, axis=0)
        x_normalize = (x - mean_minibatch) / np.sqrt(var_minibatch + eps)
        out = gamma * x_normalize + beta

        mean_running = running_mean
        var_running = running_var

        mean_running = momentum * mean_running + (1 - momentum) * mean_minibatch
        var_running = momentum * var_running + (1 - momentum) * var_minibatch
        bn_param['running_mean'] = mean_running
        bn_param['running_var'] = var_running

        cache = (
            'minibatch_mean': var_minibatch,
            'x_centralize': (x - mean_minibatch),
            'x_normalize': x_normalize,
            'gamma': gamma,
            'eps': eps
        )

    # ===== #
    # 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

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

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

# Store the updated running mean 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 = dout.shape[0]
    minibatch_var = cache.get('minibatch_var')
    x_centralize = cache.get('x_centralize')
    x_normalize = cache.get('x_normalize')
    gamma = cache.get('gamma')
    eps = cache.get('eps')

    # calculate dx
    ddxat = dout * np.sqrt(minibatch_var + eps)
    sqrt_var = np.sqrt(minibatch_var + eps)
    dsqrt_var = np.sum(ddxat * x_centralize, axis=0) / (sqrt_var**2)
    dvar = dsqrt_var * 0.5 / sqrt_var
    dx1 = ddxat + dsqrt_var * dvar * np.ones_like(dout) / N
    dx2 = np.sum(dx1, axis=0) * np.ones_like(dout) / N
    dx = dx1 + dx2

    # calculate dgamma and dgamma
    dbeta = np.sum(dout, axis=0)
    dgamma = np.sum(dout * x_normalize, axis=0)

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

    return dx, dgamma, dbeta

def dropout_forward(x, dropout_param):
    """
    Performs the forward pass for (inverted) dropout.

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

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

    """
    p, mode = dropout_param['p'], dropout_param['mode']
    if 'seed' in dropout_param:
        np.random.seed(dropout_param['seed'])

    mask = None
    out = None

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

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

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

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

        out = x

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

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

    return out, cache

def dropout_backward(dout, cache):
    """
    Perform the backward pass for (inverted) dropout.

    Inputs:
    - dout: Upstream derivatives, of any shape
    - cache: (dropout_param, mask) from dropout_forward.
    """
    dropout_param, mask = cache
    mode = dropout_param['mode']

    dx = None
    if mode == 'train':
        # ===== #
        # YOUR CODE HERE:
        # Implement the inverted dropout backward pass during training time.
        # ===== #

        dx = dout * mask

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

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

        dx = dout

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

    return dx

def svm_loss(x, y):
    """
    Computes the loss and gradient using multiclass SVM classification.

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

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

    """
    probs = np.exp(x - np.max(x, axis=1, keepdims=True))
    probs /= np.sum(probs, axis=1, keepdims=True)
    N = x.shape[0]
    loss = -np.sum(np.log(probs[np.arange(N), y])) / N
    dx = probs.copy()
    dx[np.arange(N), y] -= 1
    dx /= N
    return loss, dx

```

```

In [ ]:

import numpy as np

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 data augmentation or softmax; 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_dim=100, num_classes=10,
                  dropout=0, weight_scale=1e-3, reg=0.0):
        """
        Initialize a new network.

        Inputs:
        - input_dim: An integer giving the size of the input
        - hidden_dim: 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 weight_scale.
        # The dimensions of W1 should be (input_dim, hidden_dim) and the
        # dimensions of W2 should be (hidden_dim, num_classes)
        # ===== #

        W1_size = (input_dim, hidden_dim)
        W2_size = (hidden_dim, num_classes)

        self.params['W1'] = np.random.normal(loc=0.0, scale=weight_scale, size=W1_size)
        self.params['W2'] = np.random.normal(loc=0.0, scale=weight_scale, size=W2_size)
        self.params['b1'] = np.zeros(hidden_dim)
        self.params['b2'] = np.zeros(num_classes)

    # ===== #
    # 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.
        # ===== #

        W1 = self.params['W1']
        b1 = self.params['b1']
        W2 = self.params['W2']
        b2 = self.params['b2']

        H, cache_h = affine_relu_forward(X, W1, b1)
        D, cache_d = affine_relu_forward(H, W2, b2)

        scores = D

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

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

    # ===== #
    # 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*self.reg*W^2
    # for each W. Be sure to include the 0.5 multiplying factor.
    # And be sure to use the layers you prior implemented.
    # ===== #

    loss, dz = softmax_loss(scores, y)
    loss += 0.5*self.reg*(np.sum(W1**2) + np.sum(W2**2))

    dh, dw2, db2 = affine_backward(dz, cache_h)
    dx, dw1, db1 = affine_backward(dx, cache_d)

    grads['W1'] = dw1 + self.reg * W1
    grads['b1'] = db1
    grads['W2'] = dw2 + self.reg * W2
    grads['b2'] = db2

    # ===== #
    # 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=0, use_batchnorm=False, reg=0.0,
                  weight_scale=1e-2, dtype=np.float32, seed=None):
        """
        Initialize a new FullyConnectedNet.

        Inputs:
        - hidden_dims: A list of integers giving the size of each hidden layer.
        - num_classes: An integer giving the number of classes to classify.
        - dropout: Scalar between 0 and 1 giving dropout strength. If dropout=0 then
          the network should not use dropout at all.
        - use_batchnorm: Whether or not the network should use batch normalization.
        - reg: Scalar giving L2 regularization strength.
        - weight_scale: Scalar giving the standard deviation for random
          initialization of the weights.
        - dtype: A numpy datatype object, all computations will be performed using
          float64 for numeric gradient checking.
        - seed: If not None, then this value should be used to seed the random
          state. This will make the dropout layers deterministic so we can gradient check
          the model.

        """
        self.use_batchnorm = use_batchnorm
        self.use_dropout = dropout > 0
        self.reg = reg
        self.num_layers = 1 + len(hidden_dims)
        self.dtype = dtype
        self.params = {}

        # ===== #
        # YOUR CODE HERE:
        # Initialize all parameters of the network in the self.params dictionary.
        # The weights and biases of layer l are Wl and bl, and in general the
        # weights should be gamma[l+1] and self.params['beta'][l]. For layer L, 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 gamma and beta parameters for l and the beta
        # parameters to zero. The gamma and beta parameters for layer l should
        # be self.params['gamma'+str(l)] and self.params['beta'+str(l)]. For layer L,
        # should be gamma[L] and beta[L], etc. Only use batchnorm if self.use_batchnorm
        # is true and DO NOT do batch normalization the output scores.
        # ===== #

        cur_dim = input_dim
        for idx, hidden_dim in enumerate(hidden_dims):
            self.params['W'+str(idx+1)] = np.random.randn(cur_dim, hidden_dim) * weight_scale
            self.params['b'+str(idx+1)] = np.zeros(hidden_dim)

            if self.use_batchnorm:
                self.params['gamma'+str(idx+1)] = np.ones(hidden_dim)
                self.params['beta'+str(idx+1)] = np.zeros(hidden_dim)

            cur_dim = hidden_dim

        self.params['W'+str(self.num_layers)] = np.random.randn(cur_dim, num_classes) * weight_scale
        self.params['b'+str(self.num_layers)] = np.zeros(num_classes)

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

    # When using dropout we need to pass a dropout_param dictionary to each
    # (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:
            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 batchnorm 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.
        # ===== #

        fc_cache = {}
        relu_cache = {}
        dropout_cache = {}

        # flatten image
        X = np.reshape(X, (X.shape[0], -1))

        # go through all layers
        for i in range(self.num_layers - 1):
            # fc layer
            z = self.params['W'+str(i+1)] * X + self.params['b'+str(i+1)]

            # batchnorm layer
            relu_input = z
            if self.use_batchnorm:
                relu_cache[i] = batchnorm_forward(z, self.params['gamma'+str(i+1)], self.params['beta'+str(i+1)], self.bn_params[i])
            else:
                relu_input = z
            relu_output = relu_relu(relu_input)

            # dropout layer
            if self.use_dropout:
                relu_out, dropout_cache[i] = dropout_forward(relu_out, self.dropout_param)
            else:
                relu_out = relu_output

            # update X
            X = relu_out.copy()

        scores, final_cache = affine_forward(X, self.params['W'+str(self.num_layers)], self.params['b'+str(self.num_layers)])

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

    # If test mode return early
    if y is None:
        return scores

    # ===== #
    # 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, dz = softmax_loss(scores, y)
    loss += 0.5 * self.reg * (np.sum(np.square(self.params['W'+str(self.num_layers)])))
    dx, dx_back, dx_dropout, db, db_back, db_dropout = affine_backward(dz, final_cache)
    grads['W'+str(self.num_layers)] = dx + self.reg * self.params['W'+str(self.num_layers)]
    grads['b'+str(self.num_layers)] = db_back
    grads['b'+str(self.num_layers)] = db_dropout

    # go backward all layers and update weights, bias, gammas and betas
    for i in range(self.num_layers - 1, 0, -1):
        # dropout layer
        dx_dropout, dx_dropout_cache[i] = dropout_backward(dx_dropout, dropout_cache[i])
        dx_dropout = relu_backward(dx_dropout, relu_cache[i])

        # batchnorm layer
        affine_backward_input = dx_dropout
        dx_dropout, db_dropout, db_dropout_cache[i] = batchnorm_backward(dx_dropout, batchnorm_cache[i])
        dx_dropout = dx_dropout + db_dropout
        dx_dropout = dx_dropout + db_dropout_cache[i]

        grads['W'+str(i+1)] = dx_dropout + self.reg * self.params['W'+str(i+1)]
        grads['b'+str(i+1)] = db_dropout
        loss += 0.5 * self.reg * (np.sum(np.square(self.params['W'+str(i+1)])))

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

    return loss, grads

```



```
import numpy as np
```

```
"""
```

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

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

```
Inputs:
```

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

```
Returns:
```

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

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

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

```
"""
```

```
def sgd(w, dw, config=None):  
    """
```

Performs vanilla stochastic gradient descent.

```
    config format:
```

- learning_rate: Scalar learning rate.
- if config is None: config = {}
config.setdefault('learning_rate', 1e-2)

```
    w -= config['learning_rate'] * dw  
    return w, config
```

```
def sgd_momentum(w, dw, config=None):
```

Performs stochastic gradient descent with momentum.

```
    config format:
```

- learning_rate: Scalar learning rate.
 - momentum: Scalar between 0 and 1 giving the momentum value. Setting momentum = 0 reduces to sgd.
 - velocity: A numpy array of the same shape as w and dw used to store a moving average of the gradients.
- ```
 """
 if config is None: config = {}
 config.setdefault('learning_rate', 1e-2)
 config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there
 v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.

 # =====
 # YOUR CODE HERE:
 # Implement the momentum update formula. Return the updated weights
 # as next_w, and the updated velocity as v.
 # =====
 v = config['momentum'] * v + config['learning_rate'] * dw
 next_w = v + v
```

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

```
 config['velocity'] = v
```

```
 return next_w, config
```

```
def sgd_nesterov_momentum(w, dw, config=None):
```

Performs stochastic gradient descent with Nesterov momentum.

```
 config format:
```

- learning\_rate: Scalar learning rate.
  - momentum: Scalar between 0 and 1 giving the momentum value.
  - decay\_rate: Scalar between 0 and 1 reduces to sgd.
  - velocity: A numpy array of the same shape as w and dw used to store a moving average of the gradients.
- ```
    """  
    if config is None: config = {}  
    config.setdefault('learning_rate', 1e-2)  
    config.setdefault('momentum', 0.9) # set momentum to 0.9 if it wasn't there  
    v = config.get('velocity', np.zeros_like(w)) # gets velocity, else sets it to zero.  
  
    # =====  
    # YOUR CODE HERE:  
    # Implement the momentum update formula. Return the updated weights  
    # as next_w, and the updated velocity as v.  
    # =====  
    v_0 = v  
    v = config['momentum'] * v + config['learning_rate'] * dw  
    w += v + config['momentum'] * (v - v_0)  
    next_w = w
```

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

```
    config['velocity'] = v
```

```
    return next_w, config
```

```
def rmsprop(w, dw, config=None):
```

Uses the RMSProp update rule, which uses a moving average of squared gradient values to set adaptive per-parameter learning rates.

```
    config format:
```

- learning_rate: Scalar learning rate.
 - decay_rate: Decay rate for moving average of first moment of gradient.
 - epsilon: Small scalar used for smoothing to avoid dividing by zero.
 - m: Moving average of gradient.
 - v: Moving average of squared gradient.
 - t: Iteration number.
- ```
 """
 if config is None: config = {}
 config.setdefault('learning_rate', 1e-3)
 config.setdefault('beta1', 0.9)
 config.setdefault('beta2', 0.999)
 config.setdefault('epsilon', 1e-8)
 config.setdefault('v', np.zeros_like(w))
 config.setdefault('m', np.zeros_like(w))

 next_w = None
```

```
 # =====
 # YOUR CODE HERE:
 # Implement RMSProp. Store the next value of w as next_w. You need
 # to also store in config['a'] the moving average of the second
 # moment gradients, so they can be used for future gradients. Concretely,
 # config['a'] corresponds to "a" in the lecture notes.
 # =====
```

```
 config['a'] = config['decay_rate'] * config['a'] + (1 - config['decay_rate']) * (dw**2)
 next_w = w - config['learning_rate'] * dw / (np.sqrt(config['a']) + config['epsilon'])
```

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

```
 return next_w, config
```

```
def adam(w, dw, config=None):
```

Uses the Adam update rule, which incorporates moving averages of both the gradient and its square and a bias correction term.

```
 config format:
```

- learning\_rate: Scalar learning rate.
  - beta1: Decay rate for moving average of first moment of gradient.
  - beta2: Decay rate for moving average of second moment of gradient.
  - epsilon: Small scalar used for smoothing to avoid dividing by zero.
  - m: Moving average of gradient.
  - v: Moving average of squared gradient.
  - t: Iteration number.
- ```
    """  
    if config is None: config = {}  
    config.setdefault('learning_rate', 1e-3)  
    config.setdefault('beta1', 0.9)  
    config.setdefault('beta2', 0.999)  
    config.setdefault('epsilon', 1e-8)  
    config.setdefault('v', np.zeros_like(w))  
    config.setdefault('m', np.zeros_like(w))  
    config.setdefault('t', 0)  
  
    next_w = None
```

```
    # =====  
    # YOUR CODE HERE:  
    # Implement Adam. Store the next value of w as next_w. You need  
    # to also store in config['a'] the moving average of the second  
    # moment gradients, and in config['v'] the moving average of the  
    # first moments. Finally, store in config['t'] the increasing time.  
    # =====
```

```
    beta1 = config['beta1']  
    beta2 = config['beta2']  
    t = config['t'] + 1  
  
    v = beta1 * config['v'] + (1 - beta1) * dw  
    a = beta2 * config['a'] + (1 - beta2) * (dw**2)  
    v_corrected = v / (1 - beta1**t)  
    a_corrected = a / (1 - beta2**t)  
    next_w = w - config['learning_rate'] * v_corrected / (np.sqrt(a_corrected) + config['epsilon'])
```

```
    config['v'] = v  
    config['a'] = a  
    config['t'] = t
```

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

```
    return next_w, config
```