

Batch-Normalization

February 10, 2021

1 Batch Normalization

In this notebook, you will implement the batch normalization layers of a neural network to increase its performance. Please review the details of batch normalization from the lecture notes.

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

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

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

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

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

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

```
[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.

```
[12]: # 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: [ 8.59269312 -5.23135849 26.86727096]
```

```

stds: [32.94804157 30.64593955 30.65345025]
After batch normalization (gamma=1, beta=0)
mean: [-7.10542736e-17  4.30211422e-17  1.60982339e-16]
std: [1.          0.99999999  0.99999999]
After batch normalization (nontrivial gamma, beta)
means: [11. 12. 13.]
stds: [1.          1.99999999  2.99999998]

```

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

```

[13]: # 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: [-0.04389632  0.02964754  0.0984425 ]
stds: [0.99818918  0.94906423  1.1064687 ]

```

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.

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

N, D = 4, 5
x = 5 * np.random.randn(N, D) + 12
gamma = np.random.randn(D)
beta = np.random.randn(D)
dout = np.random.randn(N, D)

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

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

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

```
dx error:  1.1248049066356201e-09
dgamma error:  1.3401129392003245e-11
dbeta error:  3.275629256263503e-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 $1e-4$.

```
[30]: 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: 2.305933157828989
W1 relative error: 7.7477420586017e-05
W2 relative error: 4.196165108723515e-06
W3 relative error: 4.560476976762307e-10
b1 relative error: 0.002221689499037893
b2 relative error: 1.2434497875801753e-06
b3 relative error: 1.20353142917276e-10
beta1 relative error: 9.529613661272138e-09
beta2 relative error: 9.025460213519346e-09
gamma1 relative error: 7.124807552229981e-09
gamma2 relative error: 4.016190136037638e-09

```

```

Running check with reg = 3.14
Initial loss: 7.215262772233972
W1 relative error: 1.5879900801850465e-07
W2 relative error: 1.2270139496151539e-05
W3 relative error: 4.642431841796509e-09
b1 relative error: 4.163336342344337e-09
b2 relative error: 8.881784197001252e-08
b3 relative error: 3.8085525338151507e-10
beta1 relative error: 3.214152745609324e-08
beta2 relative error: 1.2542086077580079e-08
gamma1 relative error: 3.715365067200385e-08
gamma2 relative error: 2.2708724620180818e-08

```

1.4 Training a deep fully connected network with batch normalization.

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

```
[31]: # 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: 2.315647
(Epoch 0 / 10) train acc: 0.149000; val_acc: 0.149000
(Epoch 1 / 10) train acc: 0.348000; val_acc: 0.273000
(Epoch 2 / 10) train acc: 0.442000; val_acc: 0.336000
(Epoch 3 / 10) train acc: 0.501000; val_acc: 0.329000
(Epoch 4 / 10) train acc: 0.572000; val_acc: 0.344000
(Epoch 5 / 10) train acc: 0.616000; val_acc: 0.334000
(Epoch 6 / 10) train acc: 0.659000; val_acc: 0.324000
(Epoch 7 / 10) train acc: 0.741000; val_acc: 0.349000
(Epoch 8 / 10) train acc: 0.712000; val_acc: 0.317000
(Epoch 9 / 10) train acc: 0.787000; val_acc: 0.331000
(Epoch 10 / 10) train acc: 0.809000; val_acc: 0.333000
```

```
(Iteration 1 / 200) loss: 2.302030
(Epoch 0 / 10) train acc: 0.121000; val_acc: 0.124000
(Epoch 1 / 10) train acc: 0.240000; val_acc: 0.210000
(Epoch 2 / 10) train acc: 0.266000; val_acc: 0.243000
(Epoch 3 / 10) train acc: 0.325000; val_acc: 0.271000
(Epoch 4 / 10) train acc: 0.366000; val_acc: 0.283000
(Epoch 5 / 10) train acc: 0.415000; val_acc: 0.311000
(Epoch 6 / 10) train acc: 0.484000; val_acc: 0.327000
(Epoch 7 / 10) train acc: 0.522000; val_acc: 0.297000
(Epoch 8 / 10) train acc: 0.611000; val_acc: 0.310000
(Epoch 9 / 10) train acc: 0.583000; val_acc: 0.312000
(Epoch 10 / 10) train acc: 0.635000; val_acc: 0.325000
```

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

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

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

plt.subplot(3, 1, 1)
plt.plot(solver.loss_history, 'o', label='baseline')
plt.plot(bn_solver.loss_history, 'o', label='batchnorm')

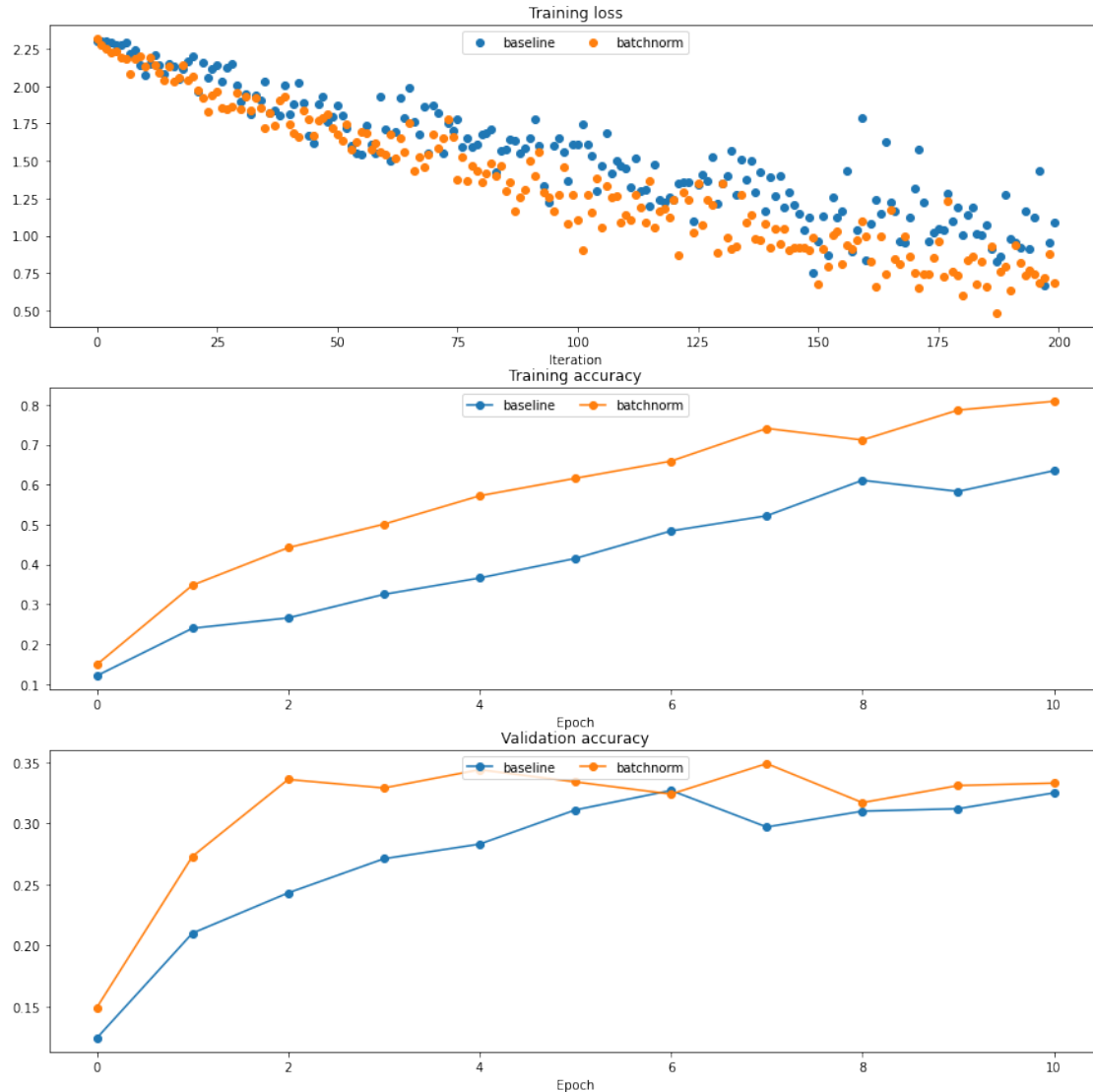
plt.subplot(3, 1, 2)
plt.plot(solver.train_acc_history, '-o', label='baseline')
plt.plot(bn_solver.train_acc_history, '-o', label='batchnorm')

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

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

<ipython-input-32-8e49aa315b6d>:13: MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

```
plt.subplot(3, 1, 1)
<ipython-input-32-8e49aa315b6d>:17: MatplotlibDeprecationWarning: Adding an axes
using the same arguments as a previous axes currently reuses the earlier
instance. In a future version, a new instance will always be created and
returned. Meanwhile, this warning can be suppressed, and the future behavior
ensured, by passing a unique label to each axes instance.
plt.subplot(3, 1, 2)
<ipython-input-32-8e49aa315b6d>:21: MatplotlibDeprecationWarning: Adding an axes
using the same arguments as a previous axes currently reuses the earlier
instance. In a future version, a new instance will always be created and
returned. Meanwhile, this warning can be suppressed, and the future behavior
ensured, by passing a unique label to each axes instance.
plt.subplot(3, 1, 3)
<ipython-input-32-8e49aa315b6d>:26: MatplotlibDeprecationWarning: Adding an axes
using the same arguments as a previous axes currently reuses the earlier
instance. In a future version, a new instance will always be created and
returned. Meanwhile, this warning can be suppressed, and the future behavior
ensured, by passing a unique label to each axes instance.
plt.subplot(3, 1, i)
```

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.

```
[33]: # 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

C:\Users\Ashwin\Desktop\UCLA\current classes\ece
247\hw4\HW4-code\nndl\layers.py:444: RuntimeWarning: divide by zero encountered
in log
    loss = -np.sum(np.log(probs[np.arange(N), y])) / N

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

```

```

[34]: # 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')

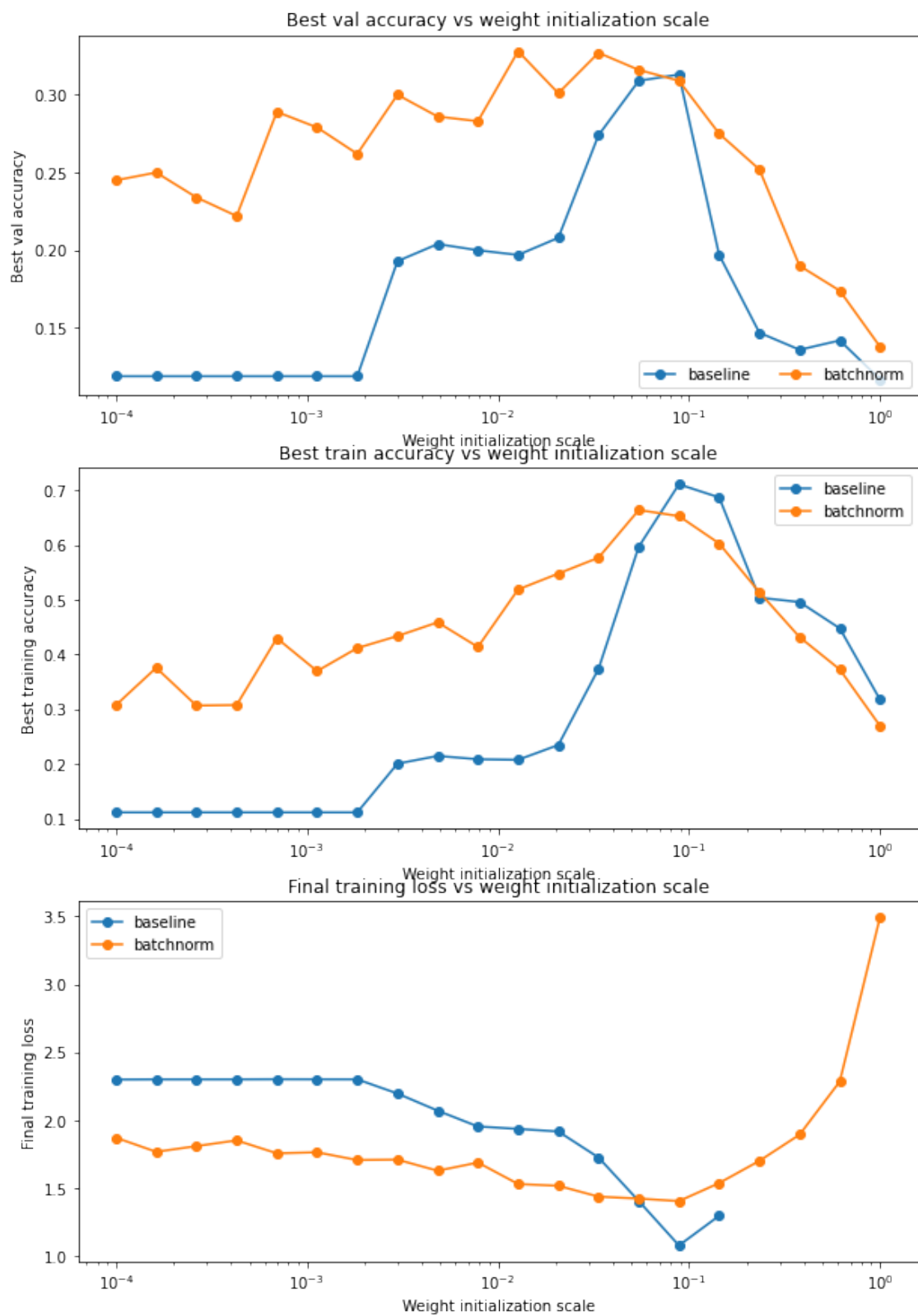
plt.subplot(3, 1, 2)
plt.title('Best train accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Best training accuracy')
plt.semilogx(weight_scales, best_train_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_train_accs, '-o', label='batchnorm')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

(note there's a divide by 0 error, and adding a small epsilon value to the softmax function should solve the problem. However my graphs looked correct regardless, so I figured this wasn't necessary in my case)

The whole point of batchnorm is to reduce extreme changes based on initial weights propagating through the hidden layers. Hence, we expect the batchnorm neural network to be more resistant to change than the baseline neural network due to different initial weights.

Indeed, this is what we see. In particular, batchnorm prevents explosion (weights going to infinity), and vanishing (weights going to 0); however, it looks like $\sim 10^{-3}$, the baseline model's weights started to vanish, leading to low training accuracy, and similarly around $10^0 = 1$ for exploding weights. Note the batchnorm neural network is resistant to these low/high weights due to its enforcement of unit variance + zero mean.

Additionally, gradient descent converges faster (less zigzagging) when batchnorm is employed, hence leading to higher performance; we can see that batchnorm mostly outperforms the baseline model over all weights.

1.7.1 fc_net.py:

```
[ ]: import numpy as np
import pdb

from .layers import *
from .layer_utils import *

"""
This code was originally written for CS 231n at Stanford University
(cs231n.stanford.edu). It has been modified in various areas for use in the
ECE 239AS class at UCLA. This includes the descriptions of what code to
implement as well as some slight potential changes in variable names to be
consistent with class nomenclature. We thank Justin Johnson & Serena Yeung for
permission to use this code. To see the original version, please visit
cs231n.stanford.edu.
"""

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

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

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

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

```
"""
def __init__(self, input_dim=3*32*32, hidden_dims=100, num_classes=10,
              dropout=0, weight_scale=1e-3, reg=0.0):
    """
    Initialize a new network.

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

    # ===== #
    # YOUR CODE HERE:
    # Initialize W1, W2, b1, and b2. Store these as self.params['W1'],
    # self.params['W2'], self.params['b1'] and self.params['b2']. The
    # biases are initialized to zero and the weights are initialized
    # so that each parameter has mean 0 and standard deviation weight_scale.
    # The dimensions of W1 should be (input_dim, hidden_dim) and the
    # dimensions of W2 should be (hidden_dims, num_classes)
    # ===== #
    self.params['W1'] = np.random.normal(scale=weight_scale, size = (input_dim,
→hidden_dims), loc=0.0)
    self.params['W2'] = np.random.normal(scale=weight_scale, size =
→(hidden_dims, num_classes), loc=0.0)
    self.params['b1'] = np.zeros(hidden_dims)
    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, \dots, 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, W2, b1, b2 = self.params['W1'], self.params['W2'], self.params['b1'],
↪ self.params['b2']

hidden, hidden_cache = affine_relu_forward(X, W1, b1)

scores, scores_cache = affine_forward(hidden, W2, b2) #forward prop

=====

END YOUR CODE HERE

=====

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

if y is None:

 return scores

loss, grads = 0, {}

=====

YOUR CODE HERE:

Implement the backward pass of the two-layer neural net. Store
the loss as the variable 'loss' and store the gradients in the
'grads' dictionary. For the grads dictionary, grads['W1'] holds
the gradient for W1, grads['b1'] holds the gradient for b1, etc.
i.e., grads[k] holds the gradient for self.params[k].

#


```

# Add L2 regularization, where there is an added cost 0.5*self.reg*W^2
# for each W. Be sure to include the 0.5 multiplying factor to
# match our implementation.
#
# And be sure to use the layers you prior implemented.
# ===== #
loss, dscores = softmax_loss(scores, y)

sum_term = np.sum(W1*W1) + np.sum(W2*W2)
reg_term = 0.5 * self.reg * sum_term
loss += reg_term

dhidden, dW2, db2 = affine_backward(dscores, scores_cache)
dX, dW1, db1 = affine_relu_backward(dhidden, hidden_cache)

#update gradient

grads['W1'] = self.reg * W1 + dW1
grads['W2'] = self.reg * W2 + dW2
grads['b1'] = db1
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.
- input_dim: An integer giving the size of the input.
- num_classes: An integer giving the number of classes to classify.
- dropout: Scalar between 0 and 1 giving dropout strength. If dropout=0 then
  the network should not use dropout at all.
- use_batchnorm: Whether or not the network should use batch normalization.
- reg: Scalar giving L2 regularization strength.
- weight_scale: Scalar giving the standard deviation for random
  initialization of the weights.
- dtype: A numpy datatype object; all computations will be performed using
  this datatype. float32 is faster but less accurate, so you should use
  float64 for numeric gradient checking.
- seed: If not None, then pass this random seed to the dropout layers. This
  will make the dropout layers deterministic so we can gradient check the
  model.
"""
self.use_batchnorm = use_batchnorm
self.use_dropout = dropout > 0
self.reg = reg
self.num_layers = 1 + len(hidden_dims)
self.dtype = dtype
self.params = {}

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

    #add batchnorm

```

```

        if self.use_batchnorm and i != self.num_layers: #if i = self.
→ num_layers, index is out of range?
            gam = "gamma" + str(i)
            beta = "beta" + str(i)
            self.params[gam] = np.ones(hidden_dims[i-1])
            self.params[beta] = np.zeros(hidden_dims[i-1]) #0, 1 wouldn't work
→ lol

        #make our layers' names
        W = 'W' + str(i) # 'W' + i wouldn't work
        b = 'b' + str(i)

        size, zeros_size = 0, 0

        #last layer
        if i == self.num_layers:
            size = (hidden_dims[i-2], num_classes)
            zeros_size = num_classes

        #first layer
        elif i == 1:
            size = (input_dim, hidden_dims[i-1])
            zeros_size = hidden_dims[i-1]

        #hidden layers
        else:
            size = (hidden_dims[i-2], hidden_dims[i-1])
            zeros_size = hidden_dims[i-1]

        #update params
        self.params[W] = np.random.normal(loc=0.0, scale=weight_scale,
→ size=size)
        self.params[b] = np.zeros(zeros_size)

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

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

```

```

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

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

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

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

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

    scores = None

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

```

```

#same as in the init function
for i in range(1, self.num_layers + 1):
    W = 'W' + str(i)
    b = 'b' + str(i)
    gam = 'gamma' + str(i)
    beta = "beta" + str(i)

    if i == self.num_layers:
        dup = affine_forward(hidden[i-2], self.params[W], self.params[b])
        scores, obj = dup

    else:

        #do hidden first
        if i == 1:
            hidden_param = X
        else:
            hidden_param = hidden[i-2]

        #try to prevent duplicate code
        if self.use_batchnorm :
            appendVar = affine_batchnorm_forward(hidden_param, self.
→params[W], self.params[b], self.params[gam], self.params[beta], self.
→bn_params[i-1])
        else:
            appendVar = affine_relu_forward(hidden_param, self.params[W],
→self.params[b])

        hidden.append(appendVar[0])

        #then do hidden_cache
        obj = appendVar[1]

        hidden_cache.append(obj)

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

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

loss, grads = 0.0, {}
# ===== #
# YOUR CODE HERE:

```

```

# Implement the backwards pass of the FC net and store the gradients
# in the grads dict, so that grads[k] is the gradient of self.params[k]
# Be sure your L2 regularization includes a 0.5 factor.
#
# BATCHNORM: Incorporate the backward pass of the batchnorm.
#
# DROPOUT: Incorporate the backward pass of dropout.
# ===== #
loss, dscores = softmax_loss(scores, y)
dhidden = [] #update this each iteration

#go backwards
for i in range(self.num_layers, 0, -1):
    W = 'W' + str(i)
    b = 'b' + str(i)
    gam = 'gamma' + str(i)
    beta = 'beta' + str(i)

    sum_term = self.params[W] * self.params[W]
    tot_term = 0.5 * self.reg * np.sum(sum_term)

    loss = loss + tot_term

    if i == self.num_layers:
        dh, grads[W], grads[b] = affine_backward(dscores, hidden_cache[self.
↪num_layers-1])
    else:
        if self.use_batchnorm:
            dh, grads[W], grads[b], grads[gam], grads[beta] =
↪affine_batchnorm_backward(dh, hidden_cache[i-1])

        else:
            dh, grads[W], grads[b] = affine_relu_backward(dh,
↪hidden_cache[i-1])

    dhidden.append(dh)

    grads[W] += self.reg * self.params[W] #need to have tuple unpacking
↪above for this to work

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

return loss, grads

```

1.7.2 layers_utils.py

```
[ ]: from .layers import *

"""
This code was originally written for CS 231n at Stanford University
(cs231n.stanford.edu). It has been modified in various areas for use in the
ECE 239AS class at UCLA. This includes the descriptions of what code to
implement as well as some slight potential changes in variable names to be
consistent with class nomenclature. We thank Justin Johnson & Serena Yeung for
permission to use this code. To see the original version, please visit
cs231n.stanford.edu.
"""

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

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

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

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

def affine_batchnorm_forward(x, w, b, gamma, beta, alpha): #batchnorm parameter
    ↪ -> alpha

    affine_out, affine_cache = affine_forward(x, w, b)
```

```

    batchnorm_out, batchnorm_cache = batchnorm_forward(affine_out, gamma, beta,
↪alpha)

    y, rcache = relu_forward(batchnorm_out) #y is what we return, also y =
↪gamma*x + beta

    #rewrite cache
    cache = affine_cache, batchnorm_cache, rcache

    return y, cache

def affine_batchnorm_backward(dout, cache):

    #following the affine_relu_forward and affine_relu_backward examples
    affine_cache, batchnorm_cache, rcache = cache
    da = relu_backward(dout, rcache)

    daff, dgam, dbet = batchnorm_backward(da, batchnorm_cache)
    dx, dw, db = affine_backward(daff, affine_cache)

    return dx, dw, db, dgam, dbet

```

1.7.3 layers.py

```

[ ]: import numpy as np
import pdb

"""
This code was originally written for CS 231n at Stanford University
(cs231n.stanford.edu). It has been modified in various areas for use in the
ECE 239AS class at UCLA. This includes the descriptions of what code to
implement as well as some slight potential changes in variable names to be
consistent with class nomenclature. We thank Justin Johnson & Serena Yeung for
permission to use this code. To see the original version, please visit
cs231n.stanford.edu.
"""

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

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

```



```

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

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

# ===== #
# YOUR CODE HERE:
# Calculate the output of the forward pass. Notice the dimensions
# of w are D x M, which is the transpose of what we did in earlier
# assignments.
# ===== #
out = np.dot(x.reshape(x.shape[0], -1), w) + b
# ===== #
# END YOUR CODE HERE
# ===== #

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

def affine_backward(dout, cache):
    """
    Computes the backward pass for an affine layer.

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

    Returns a tuple of:
    - dx: Gradient with respect to x, of shape (N, d_1, ..., d_k)
    - dw: Gradient with respect to w, of shape (D, M)
    - db: Gradient with respect to b, of shape (M,)
    """
    x, w, b = cache
    dx, dw, db = None, None, None

    # ===== #
    # YOUR CODE HERE:
    # Calculate the gradients for the backward pass.

```

```

# Notice:
#   dout is N x M
#   dx should be N x d1 x ... x dk; it relates to dout through multiplication
↳with w, which is D x M
#   dw should be D x M; it relates to dout through multiplication with x,
↳which is N x D after reshaping
#   db should be M; it is just the sum over dout examples
# ===== #
dx = np.dot(dout, w.T).reshape(x.shape)

db = np.sum(dout, axis=0)

dw = np.dot(x.reshape(x.shape[0], -1).T, dout)

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

return dx, dw, db

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

    Input:
    - x: Inputs, of any shape

    Returns a tuple of:
    - out: Output, of the same shape as x
    - cache: x
    """
    # ===== #
    # YOUR CODE HERE:
    #   Implement the ReLU forward pass.
    # ===== #
    relu = lambda x : x * (x > 0)
    out = relu(x)
    # ===== #
    # END YOUR CODE HERE
    # ===== #

    cache = x
    return out, cache

def relu_backward(dout, cache):
    """

```

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

Input:

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

Returns:

- *dx: Gradient with respect to x*
- """

x = cache

```
# ===== #
# YOUR CODE HERE:
#   Implement the ReLU backward pass
# ===== #
# ReLU directs linearly to those > 0
dx = dout * (x.reshape(x.shape[0], -1) >= 0)

# ===== #
# 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)) #key,
→value to be returned if key is not found
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 sample mean and variance.
    →(changed running -> sample)
    #   (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.
    # ===== #

    #what is the point of running_mean and running_var?? idk

    #from wk5_coding_dl video

    #step 1 -- calculate sample mean + variance
    mu = np.mean(x, axis=0)
    sigma_squared = np.var(x, axis=0)
```

```

#calculate running mean + variance
running_mean = momentum * running_mean + (1.0 - momentum) * mu
running_var = momentum * running_var + (1.0 - momentum) * sigma_squared

#step 2 -- normalize with the *sample* mean and variance (Piazza @368)
→ [what does this mean?]
xhat = (x - mu) / np.sqrt(sigma_squared + eps)

#step 3 -- scale and shift

#y -> out
out = gamma*xhat + beta

cache = (gamma, beta, running_var, running_mean, mu, sigma_squared, x, eps,
→ xhat)
# ===== #
# 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'.
# ===== #

#Piazza #343 -> only 1 vector, so sample mean -> running mean
xhat = (x - running_mean) / np.sqrt(running_var + eps)
out = gamma*xhat + beta

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

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

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

return out, cache

def batchnorm_backward(dout, cache):

```

```

"""
Backward pass for batch normalization.

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

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

Returns a tuple of:
- dx: Gradient with respect to inputs x, of shape (N, D)
- dgamma: Gradient with respect to scale parameter gamma, of shape (D,)
- dbeta: Gradient with respect to shift parameter beta, of shape (D,)
"""
dx, dgamma, dbeta = None, None, None

# ===== #
# YOUR CODE HERE:
# Implement the batchnorm backward pass, calculating dx, dgamma, and dbeta.
# ===== #
gamma, beta, running_var, running_mean, mu, sigma_squared, x, eps, xhat = cache
→cache #tuple unpacking

N, D = dout.shape #rows, cols

#from lecture slides: page 47 of lecture9_training, annotated

dL_dxhat = dout * gamma

temp1 = sigma_squared + eps

#calculating dL_dsigma_squared
first_part = - (x - mu) / (2 * temp1 ** (3/2))
inside_sum = first_part * dL_dxhat

#axis=0 means along the column, axis=1 means along the row; returns an array

#when I got rid of axis=0, my error went from 10^-12 -> 1 lol
dL_dsigma_squared = np.sum(inside_sum, axis=0)

#calculated dL_dmu
dL_dmu = - np.sum(dL_dxhat, axis=0) / np.sqrt(temp1)

#calculate dL_dx = dx
#in this case, m = N

```

```

term1 = dL_dxhat / np.sqrt(temp1)
term2 = 2 * (x - mu) * dL_dsigma_squared / N
term3 = dL_dmu / N

dx = term1 + term2 + term3

#dbeta and dgamma can be computed from page 44 of lecture9_training slides
dbeta = np.sum(dout, axis=0)
dgamma = np.sum(dout*xhat, axis=0) #is axis=0 even needed?

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

return dx, dgamma, dbeta

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

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

    Outputs:
    - out: Array of the same shape as x.
    - cache: A tuple (dropout_param, mask). In training mode, mask is the dropout
      mask that was used to multiply the input; in test mode, mask is None.
    """
    p, mode = dropout_param['p'], dropout_param['mode']
    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.
# ===== #
pass
# ===== #
# END YOUR CODE HERE
# ===== #

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

    # ===== #
    # 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.
        # ===== #
        pass
        # ===== #
        # END YOUR CODE HERE
        # ===== #
    elif mode == 'test':
        # ===== #
        # YOUR CODE HERE:

```



```

# Implement the inverted dropout backward pass during test time.
# ===== #
pass
# ===== #
# END YOUR CODE HERE
# ===== #
return dx

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

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

    Returns a tuple of:
    - loss: Scalar giving the loss
    - dx: Gradient of the loss with respect to x
    """
    N = x.shape[0]
    correct_class_scores = x[np.arange(N), y]
    margins = np.maximum(0, x - correct_class_scores[:, np.newaxis] + 1.0)
    margins[np.arange(N), y] = 0
    loss = np.sum(margins) / N
    num_pos = np.sum(margins > 0, axis=1)
    dx = np.zeros_like(x)
    dx[margins > 0] = 1
    dx[np.arange(N), y] -= num_pos
    dx /= N
    return loss, dx

def softmax_loss(x, y):
    """
    Computes the loss and gradient for softmax classification.

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

    Returns a tuple of:
    - loss: Scalar giving the loss

```

```

- dx: Gradient of the loss with respect to x
"""

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

```