

BatchNormalization

December 12, 2024

```
[39]: # This mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive')

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'cs6353/assignments/assignment3/'
FOLDERNAME = 'cs6353/assignments/assignment3/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

# Now that we've mounted your Drive, this ensures that
# the Python interpreter of the Colab VM can load
# python files from within it.
import sys
sys.path.append('/content/drive/My Drive/{}'.format(FOLDERNAME))

# This downloads the CIFAR-10 dataset to your Drive
# if it doesn't already exist.

%cd /content/drive/My\ Drive/$FOLDERNAME/cs6353/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME

# Install requirements from colab_requirements.txt
# TODO: Please change your path below to the colab_requirements.txt file
! python -m pip install -r /content/drive/My\ Drive/$FOLDERNAME/requirements.txt
```

```
Drive already mounted at /content/drive; to attempt to forcibly remount, call
drive.mount("/content/drive", force_remount=True).
/content/drive/My Drive/cs6353/assignments/assignment3/cs6353/datasets
--2024-12-11 21:39:17-- https://www.cs.toronto.edu/~kriz/cifar-10-python.tar.gz
Resolving www.cs.toronto.edu (www.cs.toronto.edu)... 128.100.3.30
Connecting to www.cs.toronto.edu (www.cs.toronto.edu)|128.100.3.30|:443...
connected.
HTTP request sent, awaiting response... 200 OK
Length: 170498071 (163M) [application/x-gzip]
Saving to: 'cifar-10-python.tar.gz'
```

```
cifar-10-python.tar 100%[=====>] 162.60M 34.9MB/s in 4.8s
```

2024-12-11 21:39:22 (34.2 MB/s) - 'cifar-10-python.tar.gz' saved
[170498071/170498071]

```
cifar-10-batches-py/  
cifar-10-batches-py/data_batch_4  
cifar-10-batches-py/readme.html  
cifar-10-batches-py/test_batch  
cifar-10-batches-py/data_batch_3  
cifar-10-batches-py/batches.meta  
cifar-10-batches-py/data_batch_2  
cifar-10-batches-py/data_batch_5  
cifar-10-batches-py/data_batch_1  
/content/drive/My Drive/cs6353/assignments/assignment3  
Collecting attrs==19.1.0 (from -r /content/drive/My  
Drive/cs6353/assignments/assignment3//requirements.txt (line 1))  
  Using cached attrs-19.1.0-py2.py3-none-any.whl.metadata (10 kB)  
Collecting backcall==0.1.0 (from -r /content/drive/My  
Drive/cs6353/assignments/assignment3//requirements.txt (line 2))  
  Using cached backcall-0.1.0.zip (11 kB)  
  Preparing metadata (setup.py) ... done  
Collecting bleach==3.1.0 (from -r /content/drive/My  
Drive/cs6353/assignments/assignment3//requirements.txt (line 3))  
  Using cached bleach-3.1.0-py2.py3-none-any.whl.metadata (19 kB)  
Collecting certifi==2019.6.16 (from -r /content/drive/My  
Drive/cs6353/assignments/assignment3//requirements.txt (line 4))  
  Using cached certifi-2019.6.16-py2.py3-none-any.whl.metadata (2.5 kB)  
Collecting cycycler==0.10.0 (from -r /content/drive/My  
Drive/cs6353/assignments/assignment3//requirements.txt (line 5))  
  Using cached cycycler-0.10.0-py2.py3-none-any.whl.metadata (722 bytes)  
Collecting decorator==4.4.0 (from -r /content/drive/My  
Drive/cs6353/assignments/assignment3//requirements.txt (line 6))  
  Using cached decorator-4.4.0-py2.py3-none-any.whl.metadata (3.7 kB)  
Collecting defusedxml==0.6.0 (from -r /content/drive/My  
Drive/cs6353/assignments/assignment3//requirements.txt (line 7))  
  Using cached defusedxml-0.6.0-py2.py3-none-any.whl.metadata (31 kB)  
Collecting entrypoints==0.3 (from -r /content/drive/My  
Drive/cs6353/assignments/assignment3//requirements.txt (line 8))  
  Using cached entrypoints-0.3-py2.py3-none-any.whl.metadata (1.4 kB)  
Collecting future==0.17.1 (from -r /content/drive/My  
Drive/cs6353/assignments/assignment3//requirements.txt (line 9))  
  Using cached future-0.17.1.tar.gz (829 kB)  
  Preparing metadata (setup.py) ... done  
Collecting imageio==2.5.0 (from -r /content/drive/My  
Drive/cs6353/assignments/assignment3//requirements.txt (line 10))  
  Using cached imageio-2.5.0-py3-none-any.whl.metadata (2.8 kB)  
Collecting ipykernel==5.1.2 (from -r /content/drive/My  
Drive/cs6353/assignments/assignment3//requirements.txt (line 11))
```

```

Using cached ipykernel-5.1.2-py3-none-any.whl.metadata (919 bytes)
Collecting ipython==7.8.0 (from -r /content/drive/My
Drive/cs6353/assignments/assignment3//requirements.txt (line 12))
Using cached ipython-7.8.0-py3-none-any.whl.metadata (4.3 kB)
Requirement already satisfied: ipython-genutils==0.2.0 in
/usr/local/lib/python3.10/dist-packages (from -r /content/drive/My
Drive/cs6353/assignments/assignment3//requirements.txt (line 13)) (0.2.0)
Collecting ipywidgets==7.5.1 (from -r /content/drive/My
Drive/cs6353/assignments/assignment3//requirements.txt (line 14))
Using cached ipywidgets-7.5.1-py2.py3-none-any.whl.metadata (1.8 kB)
Collecting jedi==0.15.1 (from -r /content/drive/My
Drive/cs6353/assignments/assignment3//requirements.txt (line 15))
Using cached jedi-0.15.1-py2.py3-none-any.whl.metadata (15 kB)
Collecting Jinja2==2.10.1 (from -r /content/drive/My
Drive/cs6353/assignments/assignment3//requirements.txt (line 16))
Using cached Jinja2-2.10.1-py2.py3-none-any.whl.metadata (2.2 kB)
Collecting jsonschema==3.0.2 (from -r /content/drive/My
Drive/cs6353/assignments/assignment3//requirements.txt (line 17))
Using cached jsonschema-3.0.2-py2.py3-none-any.whl.metadata (7.4 kB)
Collecting jupyter==1.0.0 (from -r /content/drive/My
Drive/cs6353/assignments/assignment3//requirements.txt (line 18))
Using cached jupyter-1.0.0-py2.py3-none-any.whl.metadata (995 bytes)
Collecting jupyter-client==5.3.1 (from -r /content/drive/My
Drive/cs6353/assignments/assignment3//requirements.txt (line 19))
Using cached jupyter_client-5.3.1-py2.py3-none-any.whl.metadata (3.6 kB)
Collecting jupyter-console==6.0.0 (from -r /content/drive/My
Drive/cs6353/assignments/assignment3//requirements.txt (line 20))
Using cached jupyter_console-6.0.0-py2.py3-none-any.whl.metadata (955 bytes)
Collecting jupyter-core==4.5.0 (from -r /content/drive/My
Drive/cs6353/assignments/assignment3//requirements.txt (line 21))
Using cached jupyter_core-4.5.0-py2.py3-none-any.whl.metadata (884 bytes)
Collecting kiwisolver==1.1.0 (from -r /content/drive/My
Drive/cs6353/assignments/assignment3//requirements.txt (line 22))
Using cached kiwisolver-1.1.0.tar.gz (30 kB)
Preparing metadata (setup.py) ... done
Collecting MarkupSafe==1.1.1 (from -r /content/drive/My
Drive/cs6353/assignments/assignment3//requirements.txt (line 23))
Using cached MarkupSafe-1.1.1.tar.gz (19 kB)
Preparing metadata (setup.py) ... done
Collecting matplotlib==3.1.1 (from -r /content/drive/My
Drive/cs6353/assignments/assignment3//requirements.txt (line 24))
Using cached matplotlib-3.1.1.tar.gz (37.8 MB)
Preparing metadata (setup.py) ... done
Collecting mistune==0.8.4 (from -r /content/drive/My
Drive/cs6353/assignments/assignment3//requirements.txt (line 25))
Using cached mistune-0.8.4-py2.py3-none-any.whl.metadata (8.5 kB)

```

```
ERROR: Could not find a version that satisfies the requirement mkl-  
fft==1.0.6 (from versions: 1.3.6, 1.3.8, 1.3.11)  
ERROR: No matching distribution found for mkl-fft==1.0.6
```

1 Batch Normalization

One way to make deep networks easier to train is to use more sophisticated optimization procedures such as SGD+momentum, and RMSProp. Another strategy is to change the architecture of the network to make it easier to train. One idea along these lines is batch normalization which was proposed by [3] in 2015.

The idea is relatively straightforward. Machine learning methods tend to work better when their input data consists of uncorrelated features with zero mean and unit variance. When training a neural network, we can preprocess the data before feeding it to the network to explicitly decorrelate its features; this will ensure that the first layer of the network sees data that follows a nice distribution. However, even if we preprocess the input data, the activations at deeper layers of the network will likely no longer be decorrelated and will no longer have zero mean or unit variance since they are output from earlier layers in the network. Even worse, during the training process the distribution of features at each layer of the network will shift as the weights of each layer are updated.

The authors of [1] hypothesize that the shifting distribution of features inside deep neural networks may make training deep networks more difficult. To overcome this problem, [3] proposes to insert batch normalization layers into the network. At training time, a batch normalization layer uses a minibatch of data to estimate the mean and standard deviation of each feature. These estimated means and standard deviations are then used to center and normalize the features of the minibatch. A running average of these means and standard deviations is kept during training, and at test time these running averages are used to center and normalize features.

It is possible that this normalization strategy could reduce the representational power of the network, since it may sometimes be optimal for certain layers to have features that are not zero-mean or unit variance. To this end, the batch normalization layer includes learnable shift and scale parameters for each feature dimension.

[1] Sergey Ioffe and Christian Szegedy, “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift”, ICML 2015.

```
[40]: # As usual, a bit of setup  
import time  
import numpy as np  
import matplotlib.pyplot as plt  
from cs6353.classifiers.fc_net import *  
from cs6353.data_utils import get_CIFAR10_data  
from cs6353.gradient_check import eval_numerical_gradient,   
    ↪ eval_numerical_gradient_array  
from cs6353.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'

%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))))

def print_mean_std(x,axis=0):
    print(' means: ', x.mean(axis=axis))
    print(' stds: ', x.std(axis=axis))
    print()

```

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

```
%reload_ext autoreload
```

```

[41]: # Load the (preprocessed) CIFAR10 data.
data = get_CIFAR10_data()
for k, v in data.items():
    print('%s: ' % k, v.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 Batch Normalization: Forward

In the file `cs6353/layers.py`, implement the batch normalization forward pass in the function `batchnorm_forward`. Once you have done so, run the following to test your implementation.

Referencing the paper linked to above would be helpful!

```

[42]: # 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
np.random.seed(231)
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_mean_std(a,axis=0)

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

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

```

Before batch normalization:

```

means:  [ -2.3814598 -13.18038246  1.91780462]
stds:   [27.18502186 34.21455511 37.68611762]

```

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

```

means:  [5.32907052e-17 7.04991621e-17 1.85962357e-17]
stds:   [0.99999999 1.          1.          ]

```

After batch normalization (gamma= [1. 2. 3.] , beta= [11. 12. 13.])

```

means:  [11. 12. 13.]
stds:   [0.99999999 1.99999999 2.99999999]

```

[43]: *# 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.*

```

np.random.seed(231)
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 range(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_mean_std(a_norm,axis=0)

```

```

After batch normalization (test-time):
means:  [-0.03927354 -0.04349152 -0.10452688]
stds:   [1.01531428 1.01238373 0.97819988]

```

2 Batch normalization: Backward Pass

Now implement the backward pass for batch normalization in the function `batchnorm_backward`.

To derive the backward pass you should write out the computation graph for batch normalization and backprop through each of the intermediate nodes. Some intermediates may have multiple outgoing branches; make sure to sum gradients across these branches in the backward pass.

Once you have finished, run the following to numerically check your backward pass.

```

[44]: # Gradient check batchnorm backward pass
np.random.seed(231)
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, a, beta, bn_param)[0]
fb = lambda b: batchnorm_forward(x, gamma, b, bn_param)[0]

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

_, cache = batchnorm_forward(x, gamma, beta, bn_param)
dx, dgamma, dbeta = batchnorm_backward(dout, cache)

```

```
#You should expect to see relative errors between 1e-13 and 1e-8
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:  1.7029261167605239e-09
dgamma error:  7.420414216247087e-13
dbeta error:  2.8795057655839487e-12
```

3 Batch Normalization: Alternative Backward

In class we talked about two different implementations for the sigmoid backward pass. One strategy is to write out a computation graph composed of simple operations and backprop through all intermediate values. Another strategy is to work out the derivatives on paper. For example, you can derive a very simple formula for the sigmoid function's backward pass by simplifying gradients on paper.

Surprisingly, it turns out that you can do a similar simplification for the batch normalization backward pass too.

Given a set of inputs $X = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_N \end{bmatrix}$, we first calculate the mean $\mu = \frac{1}{N} \sum_{k=1}^N x_k$ and variance

$$v = \frac{1}{N} \sum_{k=1}^N (x_k - \mu)^2.$$

With μ and v calculated, we can calculate the standard deviation $\sigma = \sqrt{v + \epsilon}$ and normalized data Y with $y_i = \frac{x_i - \mu}{\sigma}$.

The meat of our problem is to get $\frac{\partial L}{\partial X}$ from the upstream gradient $\frac{\partial L}{\partial Y}$. It might be challenging to directly reason about the gradients over X and Y - try reasoning about it in terms of x_i and y_i first.

You will need to come up with the derivations for $\frac{\partial L}{\partial x_i}$, by relying on the Chain Rule to first calculate the intermediate $\frac{\partial \mu}{\partial x_i}, \frac{\partial v}{\partial x_i}, \frac{\partial \sigma}{\partial x_i}$, then assemble these pieces to calculate $\frac{\partial y_i}{\partial x_i}$. You should make sure each of the intermediary steps are all as simple as possible.

After doing so, implement the simplified batch normalization backward pass in the function `batchnorm_backward_alt` and compare the two implementations by running the following. Your two implementations should compute nearly identical results, but the alternative implementation should be a bit faster.

```
[45]: np.random.seed(231)
N, D = 100, 500
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'}
out, cache = batchnorm_forward(x, gamma, beta, bn_param)
```



```

t1 = time.time()
dx1, dgamma1, dbeta1 = batchnorm_backward(dout, cache)
t2 = time.time()
dx2, dgamma2, dbeta2 = batchnorm_backward_alt(dout, cache)
t3 = time.time()

print('dx difference: ', rel_error(dx1, dx2))
print('dgamma difference: ', rel_error(dgamma1, dgamma2))
print('dbeta difference: ', rel_error(dbeta1, dbeta2))
print('speedup: %.2fx' % ((t2 - t1) / (t3 - t2)))

```

```

dx difference:  6.117609279344672e-13
dgamma difference:  0.0
dbeta difference:  0.0
speedup: 2.30x

```

4 Fully Connected Networks with Batch Normalization

Now that you have a working implementation for batch normalization, go back to your `FullyConnectedNet` in the file `cs6353/classifiers/fc_net.py`. Modify your implementation to add batch normalization.

Concretely, when the `normalization` flag is set to `"batchnorm"` in the constructor, you should insert a batch normalization layer before each ReLU nonlinearity. The outputs from the last layer of the network should not be normalized. Once you are done, run the following to gradient-check your implementation.

HINT: You might find it useful to define an additional helper layer similar to those in the file `cs6353/layer_utils.py`.

```

[46]: np.random.seed(231)
N, D, H1, H2, C = 2, 15, 20, 30, 10
X = np.random.randn(N, D)
y = np.random.randint(C, size=(N,))

# You should expect losses between 1e-4~1e-10 for W,
# losses between 1e-08~1e-10 for b,
# and losses between 1e-08~1e-09 for beta and gammas.
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,
                              normalization='batchnorm')

    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('%s relative error: %.2e' % (name, rel_error(grad_num, grads[name])))
if reg == 0: print()

```

```

Running check with reg = 0
Initial loss: 2.2611955101340957
W1 relative error: 1.10e-04
W2 relative error: 2.85e-06
W3 relative error: 4.05e-10
b1 relative error: 4.44e-08
b2 relative error: 2.22e-08
b3 relative error: 1.01e-10
beta1 relative error: 7.33e-09
beta2 relative error: 1.89e-09
gamma1 relative error: 6.96e-09
gamma2 relative error: 1.96e-09

```

```

Running check with reg = 3.14
Initial loss: 6.996533220108303
W1 relative error: 1.98e-06
W2 relative error: 2.28e-06
W3 relative error: 1.11e-08
b1 relative error: 2.78e-09
b2 relative error: 2.22e-08
b3 relative error: 1.73e-10
beta1 relative error: 6.65e-09
beta2 relative error: 3.48e-09
gamma1 relative error: 8.80e-09
gamma2 relative error: 5.28e-09

```

5 Batch Normalization for Deep Networks

Run the following to train a six-layer network on a subset of 1000 training examples both with and without batch normalization.

```

[47]: np.random.seed(231)
      # 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,
    ↪normalization='batchnorm')
model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale,
    ↪normalization=None)

bn_solver = Solver(bn_model, small_data,
    num_epochs=10, batch_size=50,
    update_rule='sgd_momentum',
    optim_config={
        'learning_rate': 1e-3,
    },
    verbose=True, print_every=20)
bn_solver.train()

solver = Solver(model, small_data,
    num_epochs=10, batch_size=50,
    update_rule='sgd_momentum',
    optim_config={
        'learning_rate': 1e-3,
    },
    verbose=True, print_every=20)
solver.train()

```

```

(Iteration 1 / 200) loss: 2.340974
(Epoch 0 / 10) train acc: 0.104000; val_acc: 0.103000
(Epoch 1 / 10) train acc: 0.216000; val_acc: 0.160000
(Iteration 21 / 200) loss: 2.210432
(Epoch 2 / 10) train acc: 0.339000; val_acc: 0.232000
(Iteration 41 / 200) loss: 2.200850
(Epoch 3 / 10) train acc: 0.385000; val_acc: 0.236000
(Iteration 61 / 200) loss: 2.080494
(Epoch 4 / 10) train acc: 0.430000; val_acc: 0.261000
(Iteration 81 / 200) loss: 1.860840
(Epoch 5 / 10) train acc: 0.469000; val_acc: 0.270000
(Iteration 101 / 200) loss: 1.844514
(Epoch 6 / 10) train acc: 0.543000; val_acc: 0.277000
(Iteration 121 / 200) loss: 1.728642
(Epoch 7 / 10) train acc: 0.557000; val_acc: 0.295000
(Iteration 141 / 200) loss: 1.678779
(Epoch 8 / 10) train acc: 0.644000; val_acc: 0.299000
(Iteration 161 / 200) loss: 1.539485
(Epoch 9 / 10) train acc: 0.699000; val_acc: 0.307000
(Iteration 181 / 200) loss: 1.416017

```

```

(Epoch 10 / 10) train acc: 0.746000; val_acc: 0.309000
(Iteration 1 / 200) loss: 2.302332
(Epoch 0 / 10) train acc: 0.082000; val_acc: 0.093000
(Epoch 1 / 10) train acc: 0.107000; val_acc: 0.111000
(Iteration 21 / 200) loss: 2.302091
(Epoch 2 / 10) train acc: 0.107000; val_acc: 0.129000
(Iteration 41 / 200) loss: 2.303817
(Epoch 3 / 10) train acc: 0.118000; val_acc: 0.126000
(Iteration 61 / 200) loss: 2.300032
(Epoch 4 / 10) train acc: 0.128000; val_acc: 0.132000
(Iteration 81 / 200) loss: 2.302165
(Epoch 5 / 10) train acc: 0.112000; val_acc: 0.123000
(Iteration 101 / 200) loss: 2.302931
(Epoch 6 / 10) train acc: 0.112000; val_acc: 0.119000
(Iteration 121 / 200) loss: 2.300168
(Epoch 7 / 10) train acc: 0.113000; val_acc: 0.119000
(Iteration 141 / 200) loss: 2.302988
(Epoch 8 / 10) train acc: 0.113000; val_acc: 0.119000
(Iteration 161 / 200) loss: 2.299930
(Epoch 9 / 10) train acc: 0.112000; val_acc: 0.119000
(Iteration 181 / 200) loss: 2.299801
(Epoch 10 / 10) train acc: 0.112000; val_acc: 0.119000

```

Run the following to visualize the results from two networks trained above. You should find that using batch normalization helps the network to converge much faster.

```

[48]: def plot_training_history(title, label, baseline, bn_solvers, plot_fn,
    ↪bl_marker='.', bn_marker='.', labels=None):
    """utility function for plotting training history"""
    plt.title(title)
    plt.xlabel(label)
    bn_plots = [plot_fn(bn_solver) for bn_solver in bn_solvers]
    bl_plot = plot_fn(baseline)
    num_bn = len(bn_plots)
    for i in range(num_bn):
        label='with_norm'
        if labels is not None:
            label += str(labels[i])
        plt.plot(bn_plots[i], bn_marker, label=label)
    label='baseline'
    if labels is not None:
        label += str(labels[0])
    plt.plot(bl_plot, bl_marker, label=label)
    plt.legend(loc='lower center', ncol=num_bn+1)

plt.subplot(3, 1, 1)
plot_training_history('Training loss', 'Iteration', solver, [bn_solver], \

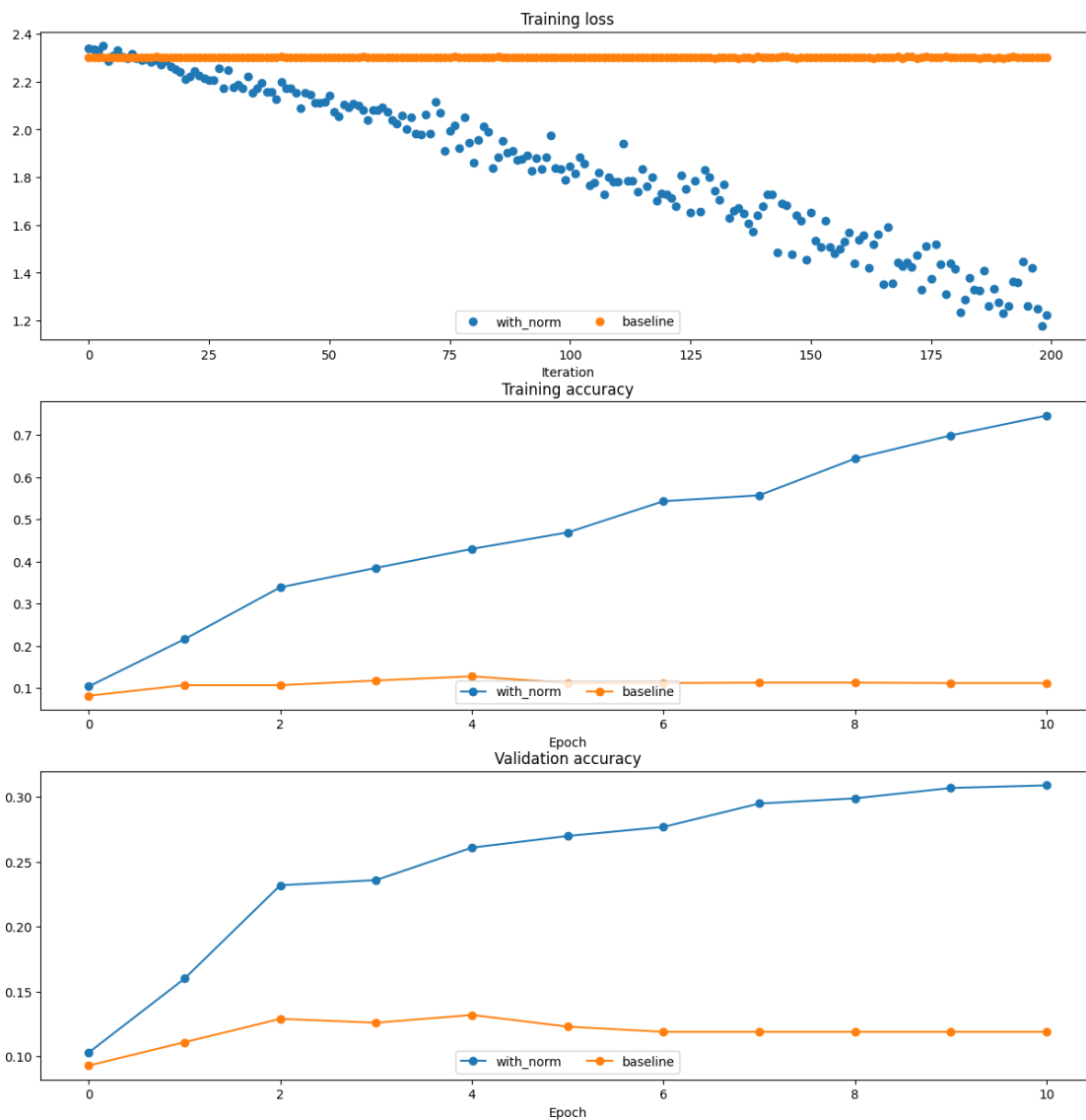
```

```

        lambda x: x.loss_history, bl_marker='o', bn_marker='o')
plt.subplot(3, 1, 2)
plot_training_history('Training accuracy', 'Epoch', solver, [bn_solver], \
        lambda x: x.train_acc_history, bl_marker='-o', \
        bn_marker='-o')
plt.subplot(3, 1, 3)
plot_training_history('Validation accuracy', 'Epoch', solver, [bn_solver], \
        lambda x: x.val_acc_history, bl_marker='-o', \
        bn_marker='-o')

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

```



6 Batch Normalization and Initialization

We will now run a small experiment to study the interaction of batch normalization and weight initialization.

The first cell will train 8-layer networks both with and without batch normalization using different scales for weight initialization. The second layer will plot training accuracy, validation set accuracy, and training loss as a function of the weight initialization scale.

```
[49]: np.random.seed(231)
      # 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_ws = {}
      solvers_ws = {}
      weight_scales = np.logspace(-4, 0, num=20)
      for i, weight_scale in enumerate(weight_scales):
          print('Running weight scale %d / %d' % (i + 1, len(weight_scales)))
          bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale,
          ↪normalization='batchnorm')
          model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale,
          ↪normalization=None)

          bn_solver = Solver(bn_model, small_data,
                          num_epochs=10, batch_size=50,
                          update_rule='sgd_momentum',
                          optim_config={
                              'learning_rate': 1e-3,
                          },
                          verbose=False, print_every=200)
          bn_solver.train()
          bn_solvers_ws[weight_scale] = bn_solver

          solver = Solver(model, small_data,
                          num_epochs=10, batch_size=50,
                          update_rule='sgd_momentum',
                          optim_config={
                              'learning_rate': 1e-3,
                          },
                          verbose=False, print_every=200)
          solver.train()
```

```
solvers_ws[weight_scale] = solver
```

```
Running weight scale 1 / 20
Running weight scale 2 / 20
Running weight scale 3 / 20
Running weight scale 4 / 20
Running weight scale 5 / 20
Running weight scale 6 / 20
Running weight scale 7 / 20
Running weight scale 8 / 20
Running weight scale 9 / 20
Running weight scale 10 / 20
Running weight scale 11 / 20
Running weight scale 12 / 20
Running weight scale 13 / 20
Running weight scale 14 / 20
Running weight scale 15 / 20
Running weight scale 16 / 20
Running weight scale 17 / 20
```

```
/content/drive/MyDrive/cs6353/assignments/assignment3/cs6353/layers.py:491:
```

```
RuntimeWarning: invalid value encountered in subtract
```

```
    shifted_logits = x - np.max(x, axis=1, keepdims=True)
```

```
Running weight scale 18 / 20
```

```
Running weight scale 19 / 20
```

```
/content/drive/MyDrive/cs6353/assignments/assignment3/cs6353/classifiers/fc_net.
```

```
py:306: RuntimeWarning: invalid value encountered in multiply
```

```
    loss += 0.5 * self.reg * np.sum(W * W)
```

```
Running weight scale 20 / 20
```

```
[50]: # 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[ws].train_acc_history))
    bn_best_train_accs.append(max(bn_solvers_ws[ws].train_acc_history))

    best_val_accs.append(max(solvers_ws[ws].val_acc_history))
    bn_best_val_accs.append(max(bn_solvers_ws[ws].val_acc_history))

    final_train_loss.append(np.mean(solvers_ws[ws].loss_history[-100:]))
    bn_final_train_loss.append(np.mean(bn_solvers_ws[ws].loss_history[-100:]))

plt.subplot(3, 1, 1)
```

```

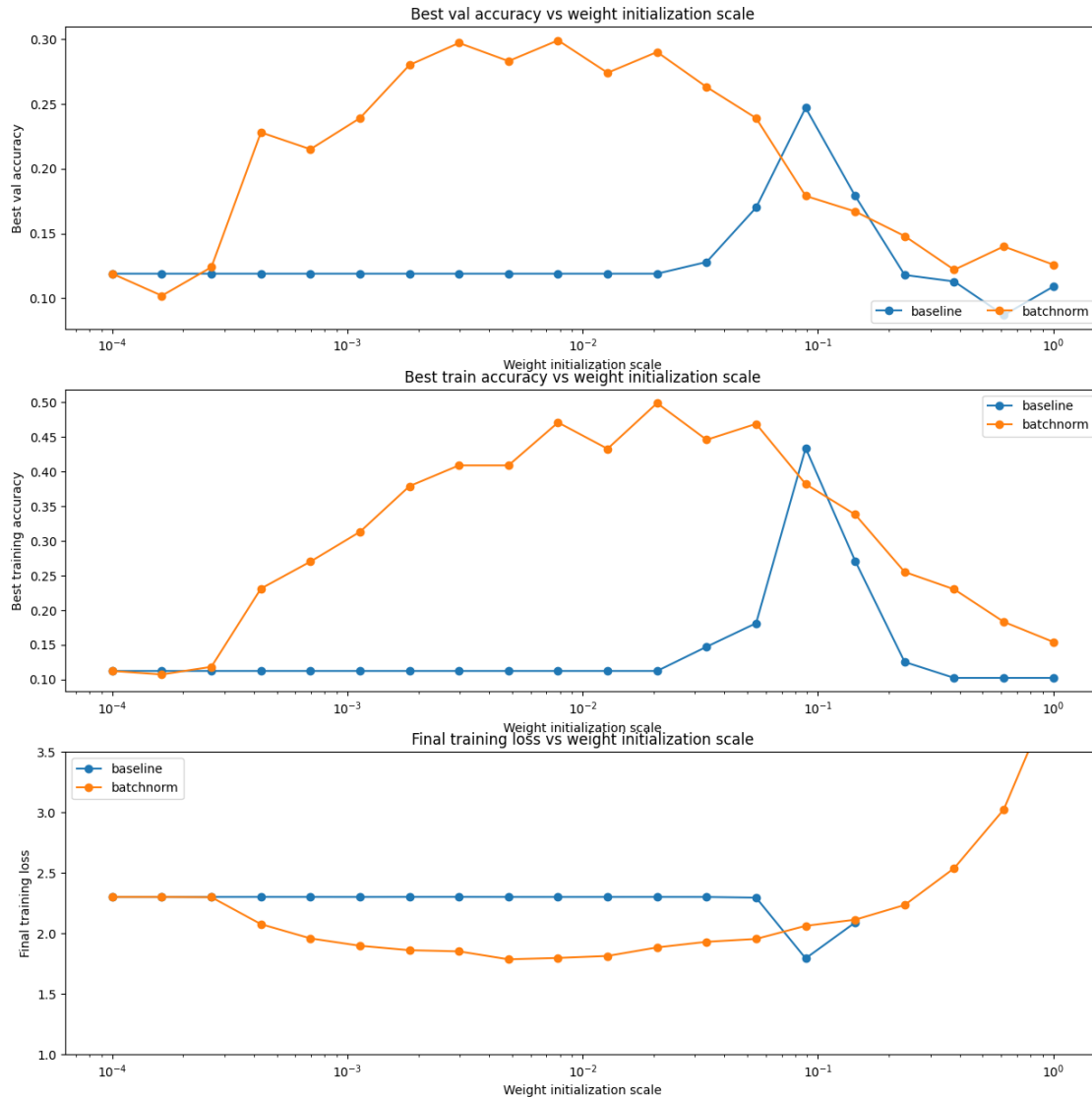
plt.title('Best val accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Best val accuracy')
plt.semilogx(weight_scales, best_val_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_val_accs, '-o', label='batchnorm')
plt.legend(ncol=2, loc='lower 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.gca().set_ylim(1.0, 3.5)

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

```

6.1 Inline Question 1:

Describe the results of this experiment. How does the scale of weight initialization affect models with/without batch normalization differently, and why?

6.2 Answer:

With no batch normalization, the network's performance depends heavily on the initial weight scale: too small or too large weights prevent meaningful learning, making the network nearly stuck at random accuracy. With batch normalization, the model becomes far more robust to different weight scales. Even when initialized poorly, batch normalization stabilizes the activations and gradients, allowing the network to train effectively and achieve better accuracy. In this experiment we can conclude that batch normalization reduces the network's sensitivity to weight initialization,

making it easier to train deep models across a wide range of weight scales.

7 Batch normalization and batch size

We will now run a small experiment to study the interaction of batch normalization and batch size.

The first cell will train 6-layer networks both with and without batch normalization using different batch sizes. The second layer will plot training accuracy and validation set accuracy over time.

```
[51]: import numpy as np

def run_batchsize_experiments(normalization_mode):
    np.random.seed(231)
    # 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'],
    }
    n_epochs=10
    weight_scale = 2e-2
    batch_sizes = [5,10,50]
    lr = 10**(-3.5)
    solver_bsize = batch_sizes[0]

    print('No normalization: batch size = ',solver_bsize)
    model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale,
    ↪normalization=None)
    solver = Solver(model, small_data,
                    num_epochs=n_epochs, batch_size=solver_bsize,
                    update_rule='sgd_momentum',
                    optim_config={
                        'learning_rate': lr,
                    },
                    verbose=False)
    solver.train()

    bn_solvers = []
    for i in range(len(batch_sizes)):
        b_size=batch_sizes[i]
        print('Normalization: batch size = ',b_size)
        bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale,
    ↪normalization=normalization_mode)
```

```

        bn_solver = Solver(bn_model, small_data,
                            num_epochs=n_epochs, batch_size=b_size,
                            update_rule='sgd_momentum',
                            optim_config={
                                'learning_rate': lr,
                            },
                            verbose=False)
        bn_solver.train()
        bn_solvers.append(bn_solver)

    return bn_solvers, solver, batch_sizes

batch_sizes = [5,10,50]
bn_solvers_bsize, solver_bsize, batch_sizes = \
    ↪run_batchsize_experiments('batchnorm')

```

```

No normalization: batch size = 5
Normalization: batch size = 5
Normalization: batch size = 10
Normalization: batch size = 50

```

```

[52]: from google.colab import drive
      drive.mount('/content/drive')

```

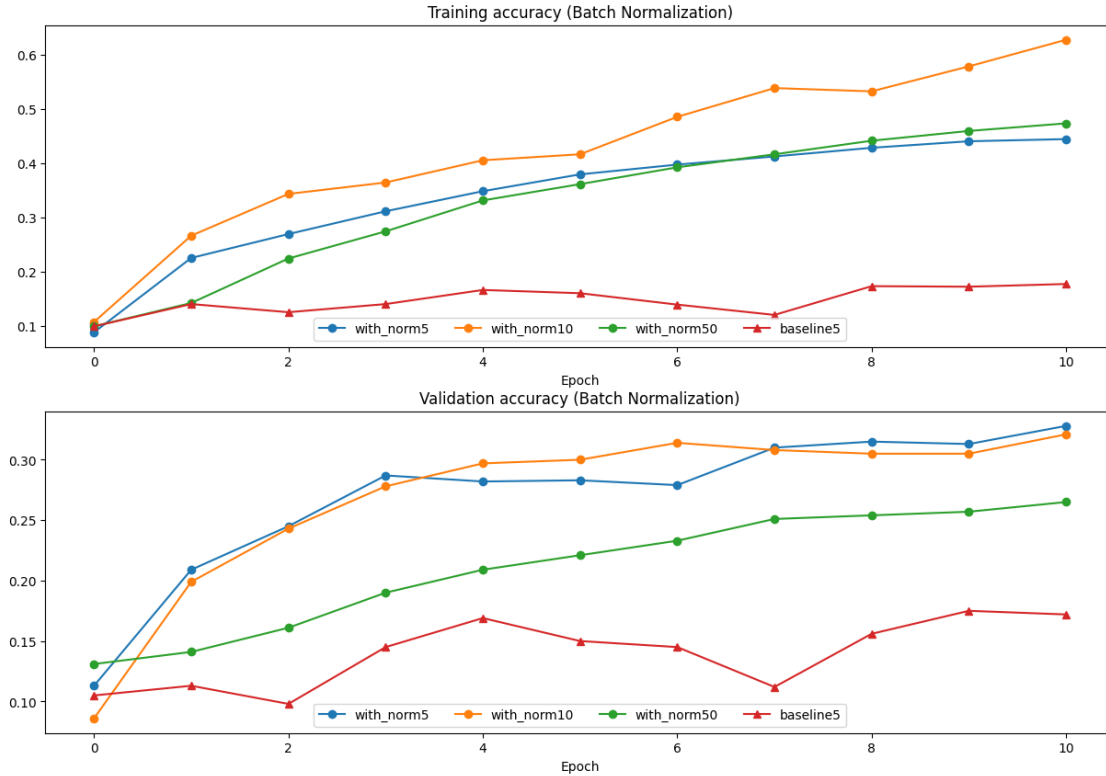
Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

```

[53]: plt.subplot(2, 1, 1)
      plot_training_history('Training accuracy (Batch Normalization)', 'Epoch', \
        ↪solver_bsize, bn_solvers_bsize, \
        ↪lambda x: x.train_acc_history, bl_marker='^-', \
        ↪bn_marker='-o', labels=batch_sizes)
      plt.subplot(2, 1, 2)
      plot_training_history('Validation accuracy (Batch Normalization)', 'Epoch', \
        ↪solver_bsize, bn_solvers_bsize, \
        ↪lambda x: x.val_acc_history, bl_marker='^-', \
        ↪bn_marker='-o', labels=batch_sizes)

      plt.gcf().set_size_inches(15, 10)
      plt.show()

```



7.1 Inline Question 2:

Describe the results of this experiment. What does this imply about the relationship between batch normalization and batch size? Why is this relationship observed?

7.2 Answer:

The experiment highlights how batch normalization helps models perform better, especially when using smaller batch sizes like 5 or 10. With batch normalization, models learn faster and achieve higher accuracy compared to those without it, which often struggle to train effectively.

This happens because batch normalization makes training more stable by managing the variations in activations during training (internal covariate shift). This stability is particularly important for small batches, where the gradients can be noisy. By normalizing layer activations, batch normalization allows the model to train efficiently even with fewer samples in each batch, making it a great tool for small batch sizes.

8 Layer Normalization

Batch normalization has proved to be effective in making networks easier to train, but the dependency on batch size makes it less useful in complex networks which have a cap on the input batch size due to hardware limitations.

Several alternatives to batch normalization have been proposed to mitigate this problem; one such technique is Layer Normalization [2]. Instead of normalizing over the batch, we normalize over the features. In other words, when using Layer Normalization, each feature vector corresponding to a single datapoint is normalized based on the sum of all terms within that feature vector.

[2] Ba, Jimmy Lei, Jamie Ryan Kiros, and Geoffrey E. Hinton. “Layer Normalization.” *stat* 1050 (2016): 21.

8.1 Inline Question 3:

Which of these data preprocessing steps is analogous to batch normalization, and which is analogous to layer normalization?

1. Scaling each image in the dataset, so that the RGB channels for each row of pixels within an image sums up to 1.
2. Scaling each image in the dataset, so that the RGB channels for all pixels within an image sums up to 1.
3. Subtracting the mean image of the dataset from each image in the dataset.
4. Setting all RGB values to either 0 or 1 depending on a given threshold.

8.2 Answer:

1. Analogous to layer normalization
2. Analogous to layer normalization
3. Analogous to batch normalization
4. Neither ??

Explanation:

1. Scaling each image so that the RGB channels for each row of pixels sum to 1: This is like Layer Normalization because the adjustment happens within a single image, focusing on each row of pixels individually. It doesn't depend on other images in the batch, just like layer normalization works independently on individual samples.
2. Scaling each image so that the RGB channels for all pixels in the image sum to 1: This is also like Layer Normalization, as it adjusts each image on its own by considering all its pixels together, similar to how layer normalization operates on a single data point.
3. Subtracting the mean image of the dataset from each image: This is like Batch Normalization because the adjustment involves the entire dataset (or batch) to compute the mean, making it similar to batch normalization, which depends on batch-wide statistics.
4. Setting RGB values to 0 or 1 based on a threshold: This is not like either Batch Normalization or Layer Normalization. It's a binarization process and doesn't involve calculations like mean or variance, which are key to normalization methods.

9 Layer Normalization: Implementation

Now you'll implement layer normalization. This step should be relatively straightforward, as conceptually the implementation is almost identical to that of batch normalization. One significant difference though is that for layer normalization, we do not keep track of the moving moments, and the testing phase is identical to the training phase, where the mean and variance are directly calculated per datapoint.

Here's what you need to do:

- In `cs6353/layers.py`, implement the forward pass for layer normalization in the function `layernorm_backward`.

Run the cell below to check your results. * In `cs6353/layers.py`, implement the backward pass for layer normalization in the function `layernorm_backward`.

Run the second cell below to check your results. * Modify `cs6353/classifiers/fc_net.py` to add layer normalization to the `FullyConnectedNet`. When the `normalization` flag is set to `"layernorm"` in the constructor, you should insert a layer normalization layer before each ReLU nonlinearity.

Run the third cell below to run the batch size experiment on layer normalization.

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

      # Simulate the forward pass for a two-layer network
      np.random.seed(231)
      N, D1, D2, D3 = 4, 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 layer normalization:')
      print_mean_std(a,axis=1)

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

      gamma = np.asarray([3.0,3.0,3.0])
      beta = np.asarray([5.0,5.0,5.0])
      # Now means should be close to beta and stds close to gamma
      print('After layer normalization (gamma=', gamma, ', beta=', beta, ')')
      a_norm, _ = layernorm_forward(a, gamma, beta, {'mode': 'train'})
      print_mean_std(a_norm,axis=1)
```

Before layer normalization:

```
means: [-59.06673243 -47.60782686 -43.31137368 -26.40991744]
stds:  [10.07429373 28.39478981 35.28360729  4.01831507]
```

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

```
means: [ 4.81096644e-16 -7.40148683e-17  2.22044605e-16 -5.92118946e-16]
stds:  [0.99999995 0.99999999 1.          0.99999969]
```

After layer normalization (gamma= [3. 3. 3.] , beta= [5. 5. 5.])

```
means: [5. 5. 5. 5.]
stds:  [2.99999985 2.99999998 2.99999999 2.99999907]
```

```
[60]: # Gradient check batchnorm backward pass
np.random.seed(231)
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)

ln_param = {}
fx = lambda x: layernorm_forward(x, gamma, beta, ln_param)[0]
fg = lambda a: layernorm_forward(x, a, beta, ln_param)[0]
fb = lambda b: layernorm_forward(x, gamma, b, ln_param)[0]

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

_, cache = layernorm_forward(x, gamma, beta, ln_param)
dx, dgamma, dbeta = layernorm_backward(dout, cache)

#You should expect to see relative errors between 1e-12 and 1e-8
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:  1.433615657860454e-09
dgamma error:  4.519489546032799e-12
dbeta error:  2.276445013433725e-12
```

10 Layer Normalization and batch size

We will now run the previous batch size experiment with layer normalization instead of batch normalization. Compared to the previous experiment, you should see a markedly smaller influence of batch size on the training history!

```
[64]: ln_solvers_bsize, solver_bsize, batch_sizes = \
    ↪run_batchsize_experiments('layernorm')

plt.subplot(2, 1, 1)
plot_training_history('Training accuracy (Layer Normalization)', 'Epoch', \
    ↪solver_bsize, ln_solvers_bsize, \
    ↪lambda x: x.train_acc_history, bl_marker='-^', \
    ↪bn_marker='-o', labels=batch_sizes)
plt.subplot(2, 1, 2)
plot_training_history('Validation accuracy (Layer Normalization)', 'Epoch', \
    ↪solver_bsize, ln_solvers_bsize, \
    ↪lambda x: x.val_acc_history, bl_marker='-^', \
    ↪bn_marker='-o', labels=batch_sizes)

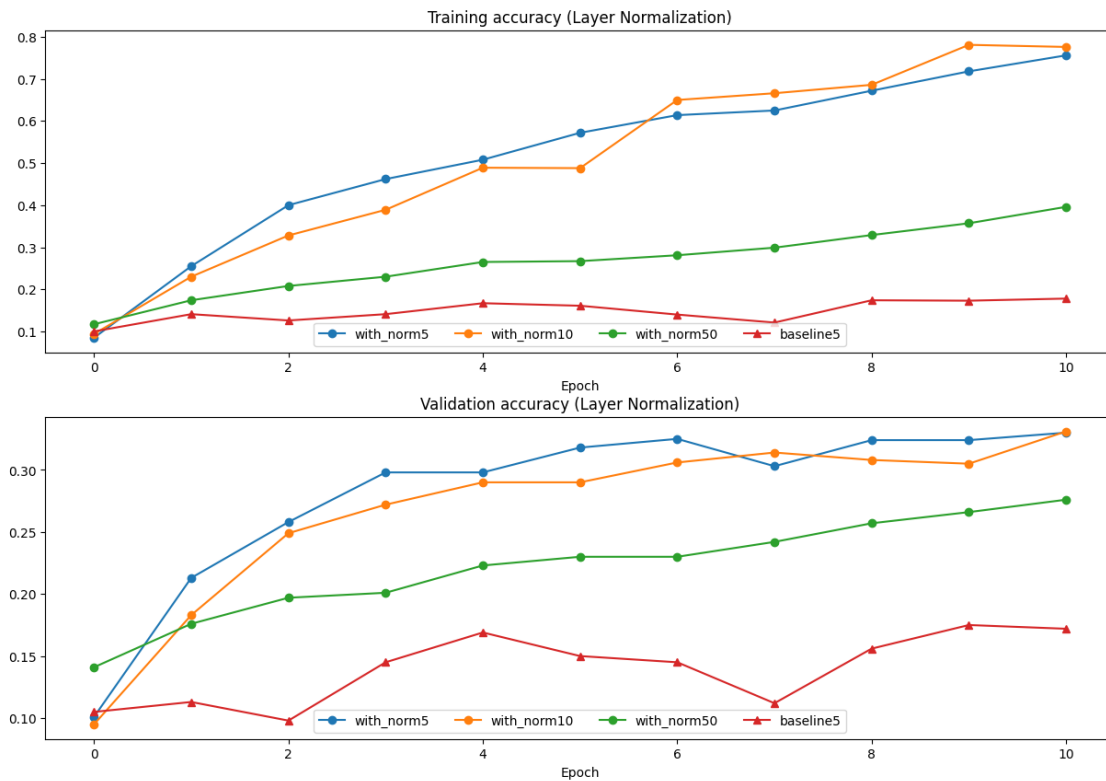
plt.gcf().set_size_inches(15, 10)
plt.show()
```

No normalization: batch size = 5

Normalization: batch size = 5

Normalization: batch size = 10

Normalization: batch size = 50



10.1 Inline Question 4:

When is layer normalization likely to not work well, and why?

1. Using it in a very deep network
2. Having a very small dimension of features
3. Having a high regularization term

10.2 Answer:

Layer normalization is more likely to not work well when the feature dimension is too small. This is because layer normalization computes the mean and variance across the feature dimension for each data point. If the number of features is too small, the computed mean and variance can become unstable or less representative of the underlying data distribution. This can lead to ineffective normalization, making the network prone to overfitting or underfitting. Thus, the lack of meaningful statistics can result in poor gradient flow, slowing down or destabilizing training.