

# Batch Normalization

---

One way to make deep networks easier to train is to use more sophisticated optimization procedures such as SGD+momentum, RMSProp, or Adam. 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 [3] 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.

[3] [Sergey Ioffe and Christian Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift", ICML 2015.](<https://arxiv.org/abs/1502.03167>)

```

# As usual, a bit of setup
import time
import numpy as np
import matplotlib.pyplot as plt
from cs231n.classifiers.fc_net 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))))

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
```

```

# 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,)

```

## Batch normalization: forward

In the file `cs231n/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!

## 解释

- BN的简单理解：机器学习领域有个很重要的假设：IID独立同分布假设，就是假设训练数据和测试数据是满足相同分布的，这是通过训练数据获得的模型能够在测试集获得好的效果的一个基本保障。而BatchNorm就是在深度神经网络训练过程中使得每一层神经网络的输入保持相同分布的。那么为什么不用BN分布有可能会变？因为网络的各层参数在不断变化导致隐藏层的输入 $x$ 的分布在不断变化。对于BN的基本思想就是：对于每个隐层神经元，把逐渐向非线性函数映射后向取值区间极限饱和区靠拢的输入分布强制拉回到均值为0方差为1的比较标准的正态分布，使得非线性变换函数的输入值落入对输入比较敏感的区域，以此避免梯度消失问题。因为梯度一直都能保持比较大的状态，所以很明显对神经网络的参数调整效率比较高，就是变动大，就是说向损失函数最优值迈动的步子大，也就是说收敛地快。
- BN的前向传播算法如下：

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots x_m\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

根据图中算法得到代码如下：

```
batch_mean = np.mean(x, axis=0) # mini-batch mean
batch_var = np.var(x, axis=0) # mini-batch variance
x_normmm = (x - batch_mean) / np.sqrt(batch_var + eps) # normalize
out = gamma * x_normmm + beta # scale and shift
cache = (x, batch_mean, batch_var, x_normmm, gamma, beta, eps)
running_mean = momentum * running_mean + (1 - momentum) * batch_mean
running_var = momentum * running_var + (1 - momentum) * batch_var
```

- 以上是训练过程中的BN前向传播算法，因为算法是基于mini-batch的，即上面训练过程中的均值和方差都是用了一个mini-batch的数据算出来的，但是实际上我们在测试的时候希望的是用全局的均值和方差。要得到全局的我们只需要将每次batch算出来的mean和var记下来然后求一次数学期望即可，如下（代码对应上面的

running\_mean和running\_var) :

$$\begin{aligned} E[x] &\leftarrow E_{\mathcal{B}}[\mu_{\mathcal{B}}] \\ \text{Var}[x] &\leftarrow \frac{m}{m-1} E_{\mathcal{B}}[\sigma_{\mathcal{B}}^2] \end{aligned}$$

由上，在

预测过程中时候的BN算法为：

$$y = \frac{\gamma}{\sqrt{\text{Var}[x] + \epsilon}} \cdot x + \left( \beta - \frac{\gamma E[x]}{\sqrt{\text{Var}[x] + \epsilon}} \right)$$

可得代码

如下：

```
x_normm = (x - running_mean) / np.sqrt(running_var + eps) # normalize
out = gamma * x_normm + beta # scale and shift
```

```
# 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  4.22578639e-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]
```

```
# 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]
```

## Batch normalization: backward

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.

## 解释

- BN的后向传播算法如下：

$$\frac{\partial \ell}{\partial \hat{x}_i} = \frac{\partial \ell}{\partial y_i} \cdot \gamma$$

$$\frac{\partial \ell}{\partial \sigma_B^2} = \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot (x_i - \mu_B) \cdot \frac{-1}{2} (\sigma_B^2 + \epsilon)^{-3/2}$$

$$\frac{\partial \ell}{\partial \mu_B} = \left( \sum_{i=1}^m \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \frac{\partial \ell}{\partial \sigma_B^2} \cdot \frac{\sum_{i=1}^m -2(x_i - \mu_B)}{m}$$

$$\frac{\partial \ell}{\partial x_i} = \frac{\partial \ell}{\partial \hat{x}_i} \cdot \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial \ell}{\partial \sigma_B^2} \cdot \frac{2(x_i - \mu_B)}{m} + \frac{\partial \ell}{\partial \mu_B} \cdot \frac{1}{m}$$

$$\frac{\partial \ell}{\partial \gamma} = \sum_{i=1}^m \frac{\partial \ell}{\partial y_i} \cdot \hat{x}_i$$

$$\frac{\partial \ell}{\partial \beta} = \sum_{i=1}^m \frac{\partial \ell}{\partial y_i}$$

根据图中算法得到代码如下：

```
(x, batch_mean, batch_var, x_normmm, gamma, beta, eps) = cache
N = x.shape[0]
dbeta = np.sum(dout, axis=0)
dgamma = np.sum(x_normmm*dout, axis = 0)
dx_normmm = gamma* dout
dbatch_var = np.sum(-1.0/2*dx_normmm*(x-batch_mean)/(batch_var+eps)**(3.0/2), axis =0)
dbatch_mean = np.sum(-1/np.sqrt(batch_var+eps)* dx_normmm, axis = 0) + 1.0/N*dbatch_var
* np.sum(-2*(x-batch_mean), axis = 0)
dx = 1/np.sqrt(batch_var+eps)*dx_normmm + dbatch_var*2.0/N*(x-batch_mean) + 1.0/N*dbatch_mean
```

```
# 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
```

## 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  $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_N \end{bmatrix}$ , 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  $\mu$  and  $v$  calculated, we can calculate the standard deviation  $\sigma = \sqrt{v + \epsilon}$  and normalized data  $\mathbf{Y}$  with  $y_i = \frac{x_i - \mu}{\sigma}$ .

The meat of our problem is to get  $\frac{\partial L}{\partial \mathbf{X}}$  from the upstream gradient  $\frac{\partial L}{\partial \mathbf{Y}}$ . It might be challenging to directly reason about the gradients over  $\mathbf{X}$  and  $\mathbf{Y}$  - 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.

## 解释

这题要求实现的是加速版的BN的前向传播过程。我们仔细看一下要求，实际上就要要门化简表达式，方法就是将一些要重复计算的结果先算出来保存下来，后面的计算过程中直接调用这个结果来计算，也就是说减少了不必要的重复计算过程，以达到提速的效果。

代码如下：

```
(x, batch_mean, batch_var, x_norm, gamma, beta, eps) = cache
N = x.shape[0]

zmean_x = x - batch_mean
sqrt_var = (batch_var+eps) ** -0.5

dx_norm = gamma * dout
var_dx = dx_norm * sqrt_var

dbatch_var = np.sum(-0.5*dx_norm*zmean_x*sqrt_var**3, axis=0)
dbatch_mean = np.sum(-var_dx, axis=0) + dbatch_var * np.sum(-2*zmean_x, axis=0) / N

dx = var_dx + 2.0*dbatch_var*zmean_x/N + dbatch_mean/N
dgamma = np.sum(x_norm*dout, axis=0)
dbeta = np.sum(dout, axis=0)
```



```

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: 9.200043712246197e-13
dgamma difference: 0.0
dbeta difference: 0.0
speedup: 0.66x

```

## 解释:

我们看到通过优化之后，BN的计算速度提升了0.66倍

# Fully Connected Nets with Batch Normalization

Now that you have a working implementation for batch normalization, go back to your `FullyConnectedNet` in the file `cs231n/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 `cs231n/layer_utils.py`. If you decide to do so, do it in the file `cs231n/classifiers/fc_net.py`.

## 解释:

为了方便起见，我们模仿之前的`affine_relu_forward`和`affine_relu_backward`两个函数，添加在`relu`之后加入了BN层的新的前向传播接口`affine_bn_relu_forward`和`affine_bn_relu_backward`。代码如下：

```

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_bn_relu_forward(x, w, b, gamma, beta, bn_param):
    a, fc_cache = affine_forward(x, w, b)
    bn, bn_cache = batchnorm_forward(a, gamma, beta, bn_param)
    out, relu_cache = relu_forward(bn)
    cache = (fc_cache, bn_cache, relu_cache)
    return out, cache

```

然后我们只需要在fc\_net.py中 FullyConnetedNet类里判断是否使用batchnorm，使得话就调用这两个接口进行前向传播和反向传播。代码暂时不贴，在后面加上layer normalization之后一起放出。

另外下面测试结果中b3的精度总是太小，可能是电脑的问题。

```

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: 3.92e-10
b1 relative error: 4.44e-08
b2 relative error: 2.22e-08
b3 relative error: 4.78e-11
beta1 relative error: 7.33e-09
beta2 relative error: 1.89e-09
gamma1 relative error: 7.57e-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: 6.94e-10
b2 relative error: 2.22e-08
b3 relative error: 2.23e-10
beta1 relative error: 6.65e-09
beta2 relative error: 3.48e-09
gamma1 relative error: 5.94e-09
gamma2 relative error: 4.14e-09
```

## Batchnorm 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.

```

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='adam',
    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='adam',
    optim_config={
        'learning_rate': 1e-3,
    },
    verbose=True, print_every=20)
solver.train()

```

```
(Iteration 1 / 200) loss: 2.340974
(Epoch 0 / 10) train acc: 0.107000; val_acc: 0.115000
(Epoch 1 / 10) train acc: 0.314000; val_acc: 0.266000
(Iteration 21 / 200) loss: 2.039345
(Epoch 2 / 10) train acc: 0.393000; val_acc: 0.285000
(Iteration 41 / 200) loss: 2.045770
(Epoch 3 / 10) train acc: 0.480000; val_acc: 0.322000
(Iteration 61 / 200) loss: 1.772057
(Epoch 4 / 10) train acc: 0.529000; val_acc: 0.316000
(Iteration 81 / 200) loss: 1.238848
(Epoch 5 / 10) train acc: 0.603000; val_acc: 0.328000
(Iteration 101 / 200) loss: 1.356776
(Epoch 6 / 10) train acc: 0.640000; val_acc: 0.317000
(Iteration 121 / 200) loss: 1.111586
(Epoch 7 / 10) train acc: 0.675000; val_acc: 0.316000
(Iteration 141 / 200) loss: 1.258607
(Epoch 8 / 10) train acc: 0.711000; val_acc: 0.308000
(Iteration 161 / 200) loss: 0.776065
(Epoch 9 / 10) train acc: 0.786000; val_acc: 0.339000
(Iteration 181 / 200) loss: 0.749002
(Epoch 10 / 10) train acc: 0.779000; val_acc: 0.336000
(Iteration 1 / 200) loss: 2.302332
(Epoch 0 / 10) train acc: 0.129000; val_acc: 0.131000
(Epoch 1 / 10) train acc: 0.283000; val_acc: 0.250000
(Iteration 21 / 200) loss: 2.041970
(Epoch 2 / 10) train acc: 0.316000; val_acc: 0.277000
(Iteration 41 / 200) loss: 1.900473
(Epoch 3 / 10) train acc: 0.373000; val_acc: 0.282000
(Iteration 61 / 200) loss: 1.713156
(Epoch 4 / 10) train acc: 0.390000; val_acc: 0.310000
(Iteration 81 / 200) loss: 1.662208
(Epoch 5 / 10) train acc: 0.434000; val_acc: 0.300000
(Iteration 101 / 200) loss: 1.696063
(Epoch 6 / 10) train acc: 0.536000; val_acc: 0.346000
(Iteration 121 / 200) loss: 1.550786
(Epoch 7 / 10) train acc: 0.530000; val_acc: 0.310000
(Iteration 141 / 200) loss: 1.436340
(Epoch 8 / 10) train acc: 0.624000; val_acc: 0.348000
(Iteration 161 / 200) loss: 0.998066
(Epoch 9 / 10) train acc: 0.660000; val_acc: 0.332000
(Iteration 181 / 200) loss: 0.943247
(Epoch 10 / 10) train acc: 0.740000; val_acc: 0.345000
```

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.

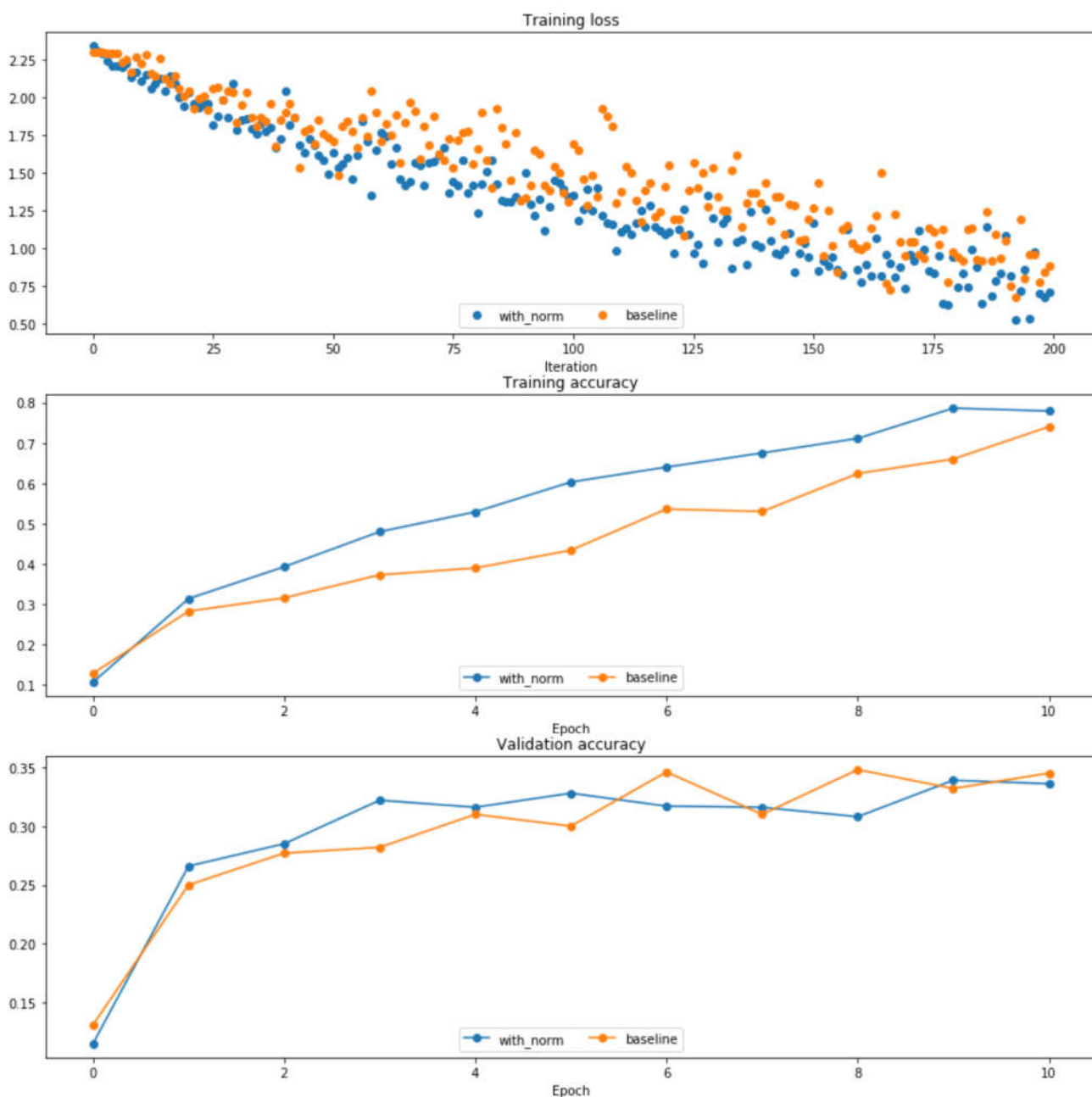
```

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()

```



注解:

由图我们可以看到，在训练集上，使用了BN层的网络学习地更快，正确率更高，增长地更快。在测试集上，使用了BN的测试结果也比较快趋于稳定。

## 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.

```

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='adam',
                        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='adam',
                    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  
Running weight scale 18 / 20  
Running weight scale 19 / 20  
Running weight scale 20 / 20

```

# Plot results of weight scale experiment
best_train_accs, bn_best_train_accs = [], []
best_val_accs, bn_best_val_accs = [], []
final_train_loss, bn_final_train_loss = [], []

for ws in weight_scales:
    best_train_accs.append(max(solvers_ws[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:]))

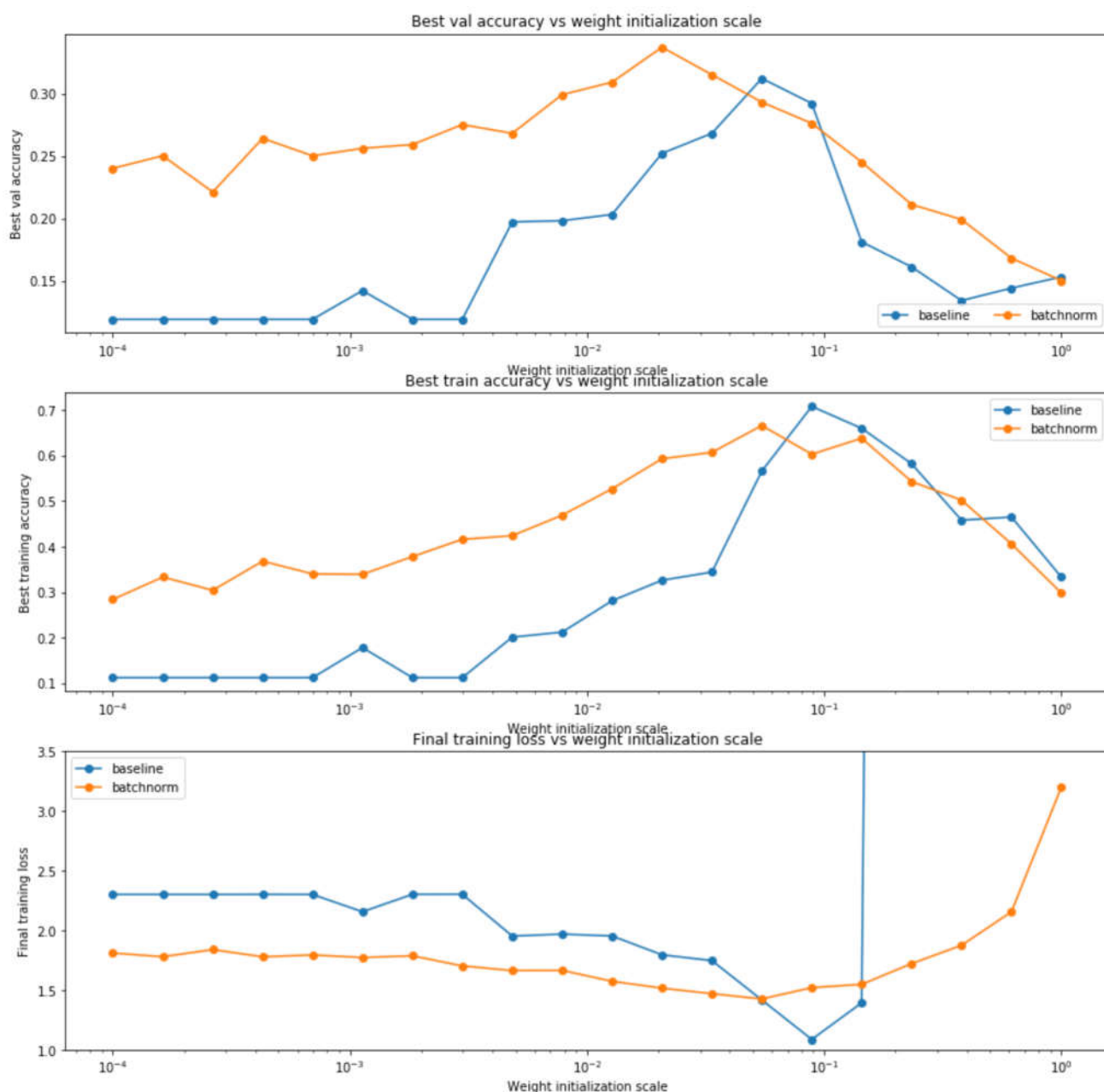
plt.subplot(3, 1, 1)
plt.title('Best val accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Best val accuracy')
plt.semilogx(weight_scales, best_val_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_val_accs, '-o', label='batchnorm')
plt.legend(ncol=2, 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()

```



## 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?

### Answer:

在没有使用batch normalization时, `weight_scale`逐层传递, 层层变化, 对模型的影响会比较大, 甚至当 `weight_scale`过大还会造成梯度爆炸的现象。但是做了归一化之后, 就不会有这个问题。

## 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.

```
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='adam',
                    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='adam',
                        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

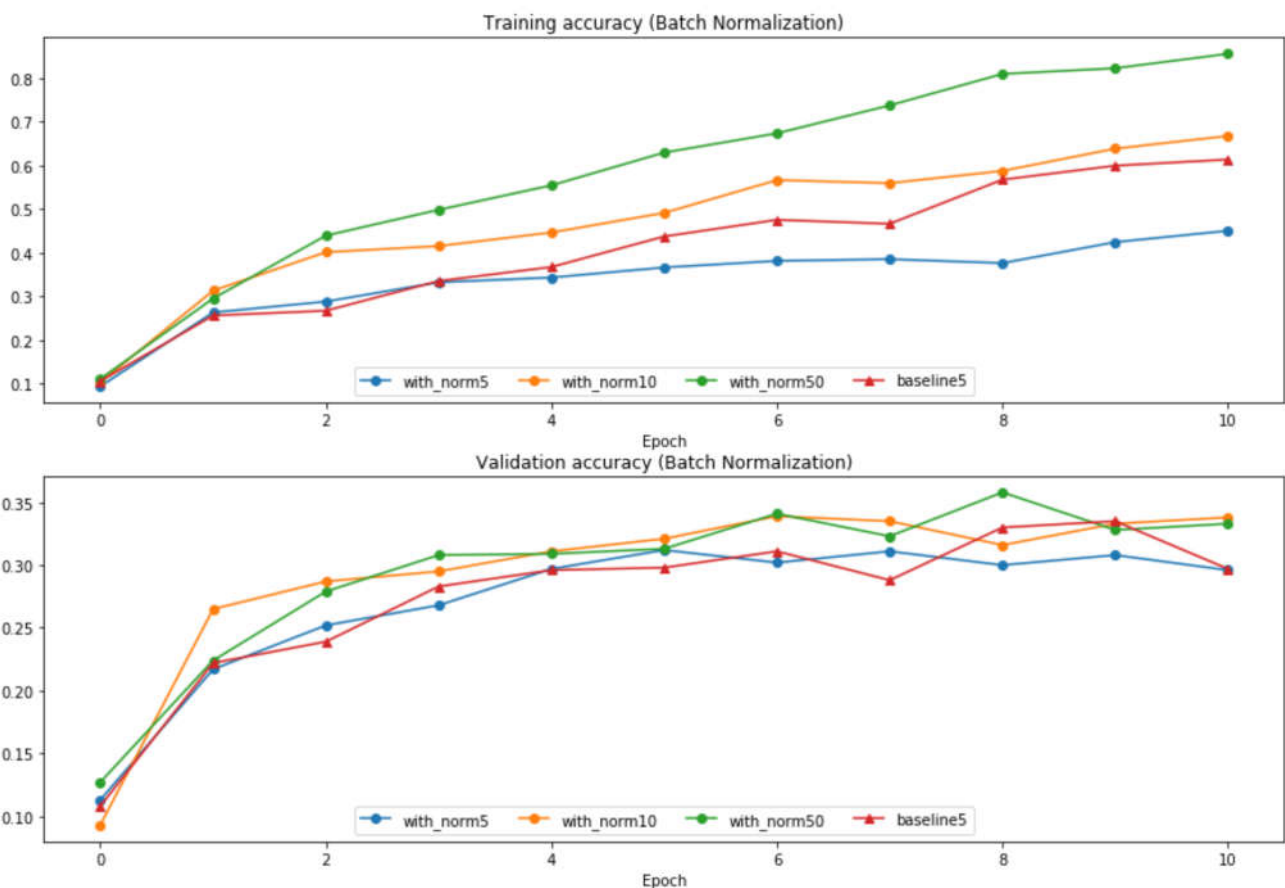
```

```

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()

```



## 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?

**Answer:**

当batch\_size较小时，由于样本数量太小，归一化受噪音影响会比较大，结果反而会不如不用归一化的baseline的结果。当时当batch\_size变大时，表示的样本的分布会更好，让结果也更加好。

## 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 [4]. 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.

[4][Ba, Jimmy Lei, Jamie Ryan Kiros, and Geoffrey E. Hinton. "Layer Normalization." stat 1050 (2016): 21.] (<https://arxiv.org/pdf/1607.06450.pdf>)

### 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.

### Answer:

---

BN: 3 LN: 2

## 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 `cs231n/layers.py`, implement the forward pass for layer normalization in the function `layernorm_forward`.

Run the cell below to check your results.

- In `cs231n/layers.py`, implement the backward pass for layer normalization in the function `layernorm_backward`.

Run the second cell below to check your results.

- Modify `cs231n/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.

## 注解：

Layer Normalization 是根据 Batch Normalization 转变过来的。我们假设一个batch的数据有N个样本，每个样本有D维特征，则我们可以将一个batch样本的数据表示成一个 $N \times D$ 的矩阵。

- 那么对于Batch Normalization，我们是对这个矩阵的**每一列**进行归一化，但BN有两个明显不足：
  1. 高度依赖于mini-batch的大小，实际使用中会对mini-Batch大小进行约束，不适合类似在线学习（mini-batch为1）情况；
  2. 不适用于RNN网络中normalize操作：BN实际使用时需要计算并且保存某一层神经网络mini-batch的均值和方差等统计信息，对于对一个固定深度的前向神经网络（DNN，CNN）使用BN，很方便；但对于RNN来说，sequence的长度是不一致的，换句话说RNN的深度不是固定的，不同的time-step需要保存不同的statics特征，可能存在一个特殊sequence比其他的sequence长很多，这样training时，计算很麻烦。
- 针对BN的问题，LN就被提出了，LN是对这个矩阵的**每一行**进行归一化。也就是说对于网络中的某一层所有神经元进行归一化。由此可以看出LN中同层神经元输入拥有相同的均值和方差，不同的输入样本有不同的均值和方差；而BN中则针对不同神经元输入计算均值和方差，同一个minibatch中的输入拥有相同的均值和方差。因此，LN不依赖于mini-batch的大小和输入sequence的深度，因此可以用于batch-size为1和RNN中对边长的输入sequence的normalize操作。

LN用于RNN进行Normalization时，取得了比BN更好的效果。但用于CNN时，效果并不如BN明显

根据上面的描述，我们知道了LN代码实际上就是将 $N \times D$ 的矩阵转置成 $D \times N$ 的矩阵然后做跟BN一样的归一化计算，注意由于LN不依赖于batch\_size,所以我们不需要保存每一次batch的均值和方差最后再求期望，即不要求前面BN中的running\_mean和running\_var。训练过程和测试过程用同样的一段归一化代码即可。

前向传播代码：

```
x = x.T
mean = np.mean(x, axis=0)
var = np.var(x, axis=0)
x_norm = (x - mean) / np.sqrt(var + eps)
x_norm = x_norm.T

out = x_norm * gamma + beta
cache = (x, mean, var, x_norm, gamma, eps)
```

反向传播代码：

```

(x, mean, var, x_norm, gamma, eps) = cache
N = dout.shape[1]

zmean_x = x - mean
sqrt_var = (var+eps) ** -0.5

dx_norm = (gamma * dout).T
var_dx = dx_norm * sqrt_var

dvar = np.sum(-0.5*dx_norm*zmean_x*sqrt_var**3, axis=0)
dmean = np.sum(-var_dx, axis=0) + dvar * np.sum(-2*zmean_x, axis=0) / N

dx = (var_dx + 2.0*dvar*zmean_x/N + dmean/N).T
dgamma = np.sum(x_norm*dout, axis=0)
dbeta = np.sum(dout, axis=0)

```

```

# 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  0.00000000e+00  7.40148683e-17 -5.55111512e-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]
```

```
# 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:  1.4336166798862177e-09
```

```
dgamma error:  4.519489546032799e-12
```

```
dbeta error:  2.276445013433725e-12
```

## 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!

```

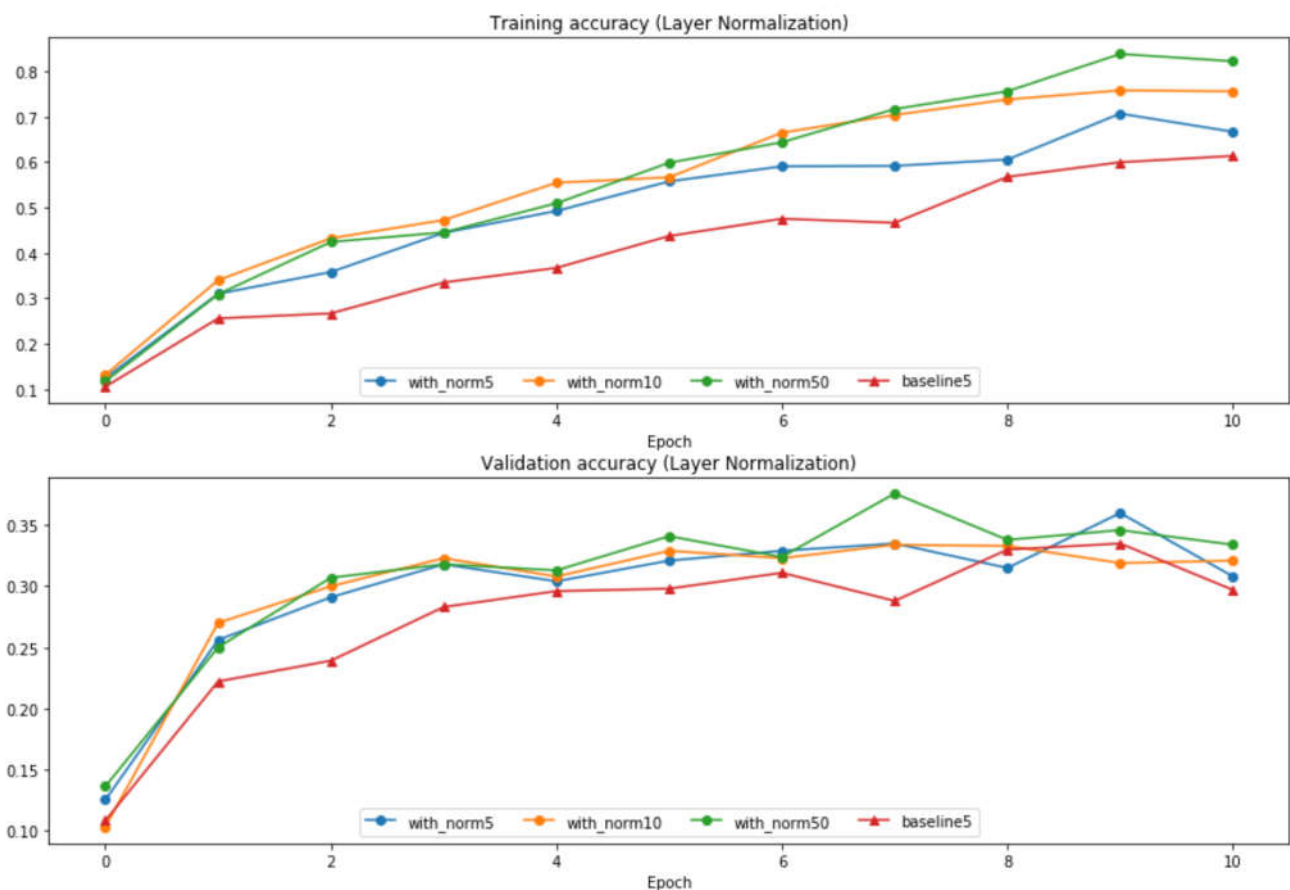
ln_solvers_bsize, solver_bsize, batch_sizes = run_batchsize_experiments('layernorm')

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



## 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

## Answer:

---

1. 网络很深的时候可以work well， 因为如果网络太深的时候batch\_size 不能太大，一个是内存可能吃不消，另一个原因是batch\_size太大的话速度会非常慢。这样用BN效果就不好，这时用LN可以取得很好的速度和效果。
2. 如果特征向量维数太少，用LN不会取得很好的效果，因为这是受噪声影响比较大，道理同batch\_size较小时用BN效果不好一样。
3. 很高正则项相当于将网络的模型结果简单化了，可能会影响LN的效果