

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).

```
In [2]: 1  ## Import and setups
2
3  import time
4  import numpy as np
5  import matplotlib.pyplot as plt
6  from nndl.fc_net import *
7  from nndl.layers import *
8  from utils.data_utils import get_CIFAR10_data
9  from utils.gradient_check import eval_numerical_gradient, eval_numerical_gradient_array
10 from utils.solver import Solver
11
12 %matplotlib inline
13 plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
14 plt.rcParams['image.interpolation'] = 'nearest'
15 plt.rcParams['image.cmap'] = 'gray'
16
17 # for auto-reloading external modules
18 # see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-ipython
19 %load_ext autoreload
20 %autoreload 2
21
22 def rel_error(x, y):
23     """ returns relative error """
24     return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

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

```
In [3]: 1  # Load the (preprocessed) CIFAR10 data.
2
3  data = get_CIFAR10_data()
4  for k in data.keys():
5      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,)
```

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.

```
In [4]: 1 # Check the training-time forward pass by checking means and variances
2 # of features both before and after batch normalization
3
4 # Simulate the forward pass for a two-layer network
5 N, D1, D2, D3 = 200, 50, 60, 3
6 X = np.random.randn(N, D1)
7 W1 = np.random.randn(D1, D2)
8 W2 = np.random.randn(D2, D3)
9 a = np.maximum(0, X.dot(W1)).dot(W2)
10
11 print('Before batch normalization:')
12 print(' means: ', a.mean(axis=0))
13 print(' stds: ', a.std(axis=0))
14
15 # Means should be close to zero and stds close to one
16 print('After batch normalization (gamma=1, beta=0)')
17 a_norm, _ = batchnorm_forward(a, np.ones(D3), np.zeros(D3), {'mode': 'train'})
18 print(' mean: ', a_norm.mean(axis=0))
19 print(' std: ', a_norm.std(axis=0))
20
21 # Now means should be close to beta and stds close to gamma
22 gamma = np.asarray([1.0, 2.0, 3.0])
23 beta = np.asarray([11.0, 12.0, 13.0])
24 a_norm, _ = batchnorm_forward(a, gamma, beta, {'mode': 'train'})
25 print('After batch normalization (nontrivial gamma, beta)')
26 print(' means: ', a_norm.mean(axis=0))
27 print(' stds: ', a_norm.std(axis=0))
```

```
Before batch normalization:
 means: [-25.27203987 -8.10778565  5.39254912]
 stds: [29.71997558 38.51999143 37.20360818]
After batch normalization (gamma=1, beta=0)
 mean: [-2.50910404e-16  1.40703421e-16 -3.48332474e-17]
 std: [0.99999999 1.         1.         ]
After batch normalization (nontrivial gamma, beta)
 means: [11. 12. 13.]
 stds: [0.99999999 1.99999999 2.99999999]
```

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

```
In [5]: 1 # Check the test-time forward pass by running the training-time
2 # forward pass many times to warm up the running averages, and then
3 # checking the means and variances of activations after a test-time
4 # forward pass.
5
6 N, D1, D2, D3 = 200, 50, 60, 3
7 W1 = np.random.randn(D1, D2)
8 W2 = np.random.randn(D2, D3)
9
10 bn_param = {'mode': 'train'}
11 gamma = np.ones(D3)
12 beta = np.zeros(D3)
13 for t in np.arange(50):
14     X = np.random.randn(N, D1)
15     a = np.maximum(0, X.dot(W1)).dot(W2)
16     batchnorm_forward(a, gamma, beta, bn_param)
17 bn_param['mode'] = 'test'
18 X = np.random.randn(N, D1)
19 a = np.maximum(0, X.dot(W1)).dot(W2)
20 a_norm, _ = batchnorm_forward(a, gamma, beta, bn_param)
21
22 # Means should be close to zero and stds close to one, but will be
23 # noisier than training-time forward passes.
24 print('After batch normalization (test-time):')
25 print(' means: ', a_norm.mean(axis=0))
26 print(' stds: ', a_norm.std(axis=0))
```

```
After batch normalization (test-time):
 means: [-0.0523224  0.15029738  0.09722565]
 stds: [0.94254433 1.01253859 1.07017573]
```

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.

```
In [6]: 1 # Gradient check batchnorm backward pass
2
3 N, D = 4, 5
4 x = 5 * np.random.randn(N, D) + 12
5 gamma = np.random.randn(D)
6 beta = np.random.randn(D)
7 dout = np.random.randn(N, D)
8
9 bn_param = {'mode': 'train'}
10 fx = lambda x: batchnorm_forward(x, gamma, beta, bn_param)[0]
11 fg = lambda gamma: batchnorm_forward(x, gamma, beta, bn_param)[0]
12 fb = lambda beta: batchnorm_forward(x, gamma, beta, bn_param)[0]
13
14 dx_num = eval_numerical_gradient_array(fx, x, dout)
15 da_num = eval_numerical_gradient_array(fg, gamma, dout)
16 db_num = eval_numerical_gradient_array(fb, beta, dout)
17
18 _, cache = batchnorm_forward(x, gamma, beta, bn_param)
19 dx, dgamma, dbeta = batchnorm_backward(dout, cache)
20 print('dx error: ', rel_error(dx_num, dx))
21 print('dgamma error: ', rel_error(da_num, dgamma))
22 print('dbeta error: ', rel_error(db_num, dbeta))
```

```
dx error: 2.1040090821756552e-08
dgamma error: 1.0792566814777373e-11
dbeta error: 3.1514364236988523e-12
```

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 $1e-4$.

```
In [7]: 1 N, D, H1, H2, C = 2, 15, 20, 30, 10
2 X = np.random.randn(N, D)
3 y = np.random.randint(C, size=(N,))
4
5 for reg in [0, 3.14]:
6     print('Running check with reg = ', reg)
7     model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
8                               reg=reg, weight_scale=5e-2, dtype=np.float64,
9                               use_batchnorm=True)
10
11     loss, grads = model.loss(X, y)
12     print('Initial loss: ', loss)
13
14     for name in sorted(grads):
15         f = lambda _: model.loss(X, y)[0]
16         grad_num = eval_numerical_gradient(f, model.params[name], verbose=False, h=1e-5)
17         print('{} relative error: {}'.format(name, rel_error(grad_num, grads[name])))
18     if reg == 0: print('\n')
```

```
Running check with reg = 0
Initial loss: 2.353457733616361
W1 relative error: 7.320331999975745e-06
W2 relative error: 1.669363329093752e-05
W3 relative error: 3.1732216224979363e-09
b1 relative error: 0.0022204460492503126
b2 relative error: 0.00222042384478982
b3 relative error: 1.3774575839430264e-10
beta1 relative error: 1.4335235401576925e-08
beta2 relative error: 1.752101382445816e-07
gamma1 relative error: 1.4239179511780286e-08
gamma2 relative error: 6.028171267693151e-09
```

```
Running check with reg = 3.14
Initial loss: 6.7810199172672
W1 relative error: 1.4181005436854099e-05
W2 relative error: 1.2865560232585118e-06
W3 relative error: 3.1851622215920295e-06
b1 relative error: 2.7755575615628914e-09
b2 relative error: 5.551115123125783e-09
b3 relative error: 9.817564423472065e-11
beta1 relative error: 3.5782495407921424e-07
beta2 relative error: 6.21788848725212e-09
gamma1 relative error: 7.734868323961019e-08
gamma2 relative error: 6.0129117102249746e-09
```

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.

```

In [8]: 1 # Try training a very deep net with batchnorm
2 hidden_dims = [100, 100, 100, 100, 100]
3
4 num_train = 1000
5 small_data = {
6     'X_train': data['X_train'][:num_train],
7     'y_train': data['y_train'][:num_train],
8     'X_val': data['X_val'],
9     'y_val': data['y_val'],
10 }
11
12 weight_scale = 2e-2
13 bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=True)
14 model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=False)
15
16 bn_solver = Solver(bn_model, small_data,
17                     num_epochs=10, batch_size=50,
18                     update_rule='adam',
19                     optim_config={
20                         'learning_rate': 1e-3,
21                     },
22                     verbose=True, print_every=200)
23 bn_solver.train()
24
25 solver = Solver(model, small_data,
26                 num_epochs=10, batch_size=50,
27                 update_rule='adam',
28                 optim_config={
29                     'learning_rate': 1e-3,
30                 },
31                 verbose=True, print_every=200)
32 solver.train()

```

```

(Iteration 1 / 200) loss: 2.287024
(Epoch 0 / 10) train acc: 0.089000; val_acc: 0.094000
(Epoch 1 / 10) train acc: 0.337000; val_acc: 0.238000
(Epoch 2 / 10) train acc: 0.437000; val_acc: 0.285000
(Epoch 3 / 10) train acc: 0.475000; val_acc: 0.302000
(Epoch 4 / 10) train acc: 0.572000; val_acc: 0.318000
(Epoch 5 / 10) train acc: 0.630000; val_acc: 0.338000
(Epoch 6 / 10) train acc: 0.688000; val_acc: 0.331000
(Epoch 7 / 10) train acc: 0.716000; val_acc: 0.323000
(Epoch 8 / 10) train acc: 0.715000; val_acc: 0.333000
(Epoch 9 / 10) train acc: 0.770000; val_acc: 0.343000
(Epoch 10 / 10) train acc: 0.815000; val_acc: 0.320000
(Iteration 1 / 200) loss: 2.302971
(Epoch 0 / 10) train acc: 0.140000; val_acc: 0.157000
(Epoch 1 / 10) train acc: 0.273000; val_acc: 0.242000
(Epoch 2 / 10) train acc: 0.340000; val_acc: 0.290000
(Epoch 3 / 10) train acc: 0.379000; val_acc: 0.296000
(Epoch 4 / 10) train acc: 0.396000; val_acc: 0.301000
(Epoch 5 / 10) train acc: 0.465000; val_acc: 0.318000
(Epoch 6 / 10) train acc: 0.478000; val_acc: 0.304000
(Epoch 7 / 10) train acc: 0.539000; val_acc: 0.320000
(Epoch 8 / 10) train acc: 0.588000; val_acc: 0.314000
(Epoch 9 / 10) train acc: 0.578000; val_acc: 0.306000
(Epoch 10 / 10) train acc: 0.616000; val_acc: 0.331000

```

```
In [9]: 1 plt.subplot(3, 1, 1)
2 plt.title('Training loss')
3 plt.xlabel('Iteration')
4
5 plt.subplot(3, 1, 2)
6 plt.title('Training accuracy')
7 plt.xlabel('Epoch')
8
9 plt.subplot(3, 1, 3)
10 plt.title('Validation accuracy')
11 plt.xlabel('Epoch')
12
13 plt.subplot(3, 1, 1)
14 plt.plot(solver.loss_history, 'o', label='baseline')
15 plt.plot(bn_solver.loss_history, 'o', label='batchnorm')
16
17 plt.subplot(3, 1, 2)
18 plt.plot(solver.train_acc_history, '-o', label='baseline')
19 plt.plot(bn_solver.train_acc_history, '-o', label='batchnorm')
20
21 plt.subplot(3, 1, 3)
22 plt.plot(solver.val_acc_history, '-o', label='baseline')
23 plt.plot(bn_solver.val_acc_history, '-o', label='batchnorm')
24
25 for i in [1, 2, 3]:
26     plt.subplot(3, 1, i)
27     plt.legend(loc='upper center', ncol=4)
28 plt.gcf().set_size_inches(15, 15)
29 plt.show()
```

/var/folders/gx/lzqxm5853svcnyjc53mcryxc0000gn/T/ipykernel_5818/38527370.py: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)
```

/var/folders/gx/lzqxm5853svcnyjc53mcryxc0000gn/T/ipykernel_5818/38527370.py: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)
```

/var/folders/gx/lzqxm5853svcnyjc53mcryxc0000gn/T/ipykernel_5818/38527370.py: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)
```

/var/folders/gx/lzqxm5853svcnyjc53mcryxc0000gn/T/ipykernel_5818/38527370.py: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)
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "orientation" which is no longer supported as of 3.3 and will become an error two minor releases later

```
fig.canvas.print_figure(bytes_io, **kw)
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "dpi" which is no longer supported as of 3.3 and will become an error two minor releases later

```
fig.canvas.print_figure(bytes_io, **kw)
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "facecolor" which is no longer supported as of 3.3 and will become an error two minor releases later

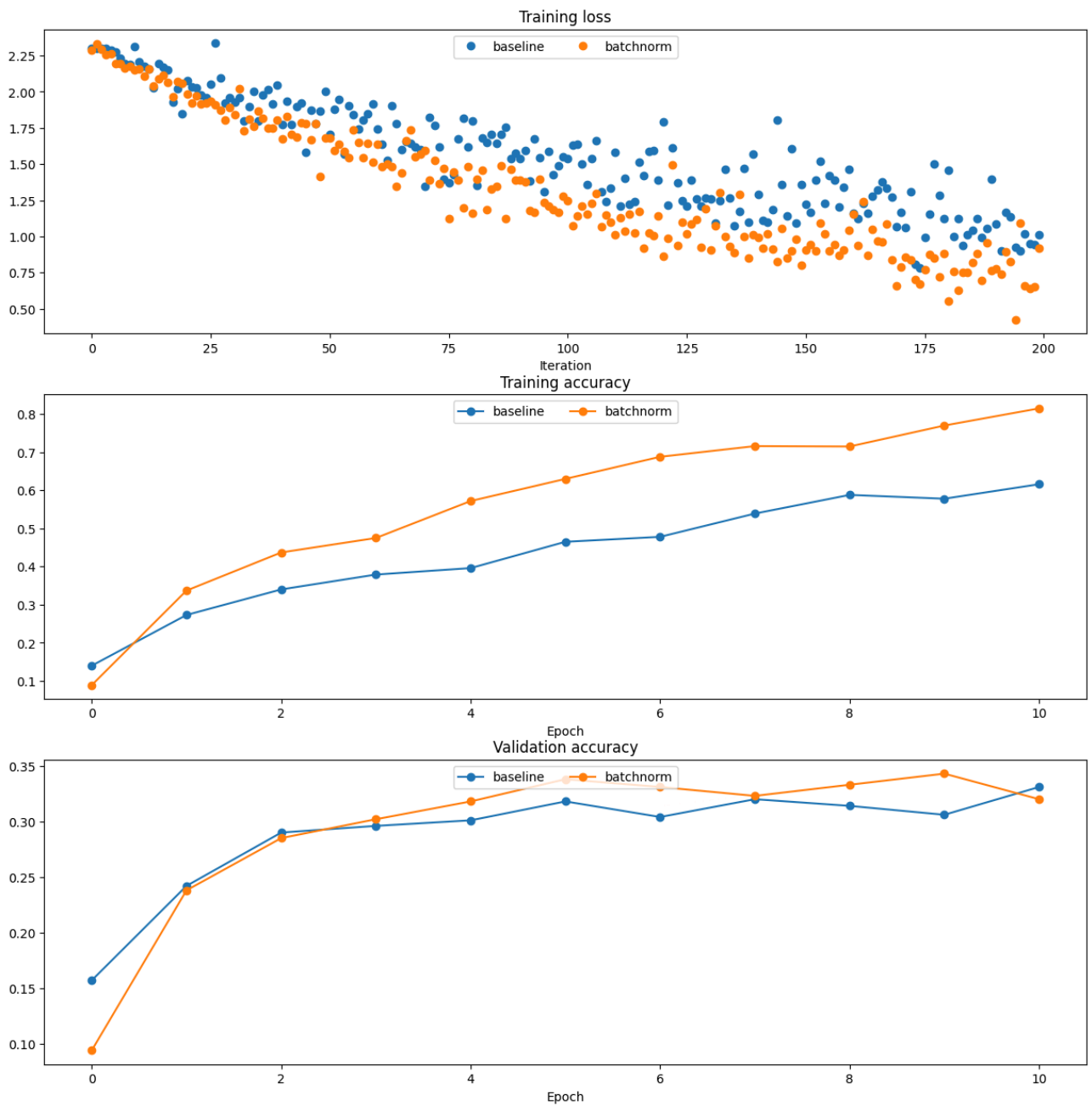
```
fig.canvas.print_figure(bytes_io, **kw)
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "edgecolor" which is no longer supported as of 3.3 and will become an error two minor releases later

```
fig.canvas.print_figure(bytes_io, **kw)
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "bbox_inches_restore" which is no longer supported as of 3.3 and will become an error two minor releases later

```
fig.canvas.print_figure(bytes_io, **kw)
```



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.

```

In [12]: 1 # Try training a very deep net with batchnorm
2 hidden_dims = [50, 50, 50, 50, 50, 50, 50]
3
4 num_train = 1000
5 small_data = {
6     'X_train': data['X_train'][:num_train],
7     'y_train': data['y_train'][:num_train],
8     'X_val': data['X_val'],
9     'y_val': data['y_val'],
10 }
11
12 bn_solvers = {}
13 solvers = {}
14 weight_scales = np.logspace(-4, 0, num=20)
15 for i, weight_scale in enumerate(weight_scales):
16     print('Running weight scale {} / {}'.format(i + 1, len(weight_scales)))
17     bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=True)
18     model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=False)
19
20     bn_solver = Solver(bn_model, small_data,
21                       num_epochs=10, batch_size=50,
22                       update_rule='adam',
23                       optim_config={
24                           'learning_rate': 1e-3,
25                       },
26                       verbose=False, print_every=200)
27     bn_solver.train()
28     bn_solvers[weight_scale] = bn_solver
29
30     solver = Solver(model, small_data,
31                     num_epochs=10, batch_size=50,
32                     update_rule='adam',
33                     optim_config={
34                         'learning_rate': 1e-3,
35                     },
36                     verbose=False, print_every=200)
37     solver.train()
38     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

```

```

/Users/sujitsilas/Desktop/UCLA/Winter 2025/EE ENGR 247/Homeworks/HW4/HW4_code/nndl/layers.py:429: 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

```



```

In [11]: 1 # Plot results of weight scale experiment
2 best_train_accs, bn_best_train_accs = [], []
3 best_val_accs, bn_best_val_accs = [], []
4 final_train_loss, bn_final_train_loss = [], []
5
6 for ws in weight_scales:
7     best_train_accs.append(max(solvers[ws].train_acc_history))
8     bn_best_train_accs.append(max(bn_solvers[ws].train_acc_history))
9
10    best_val_accs.append(max(solvers[ws].val_acc_history))
11    bn_best_val_accs.append(max(bn_solvers[ws].val_acc_history))
12
13    final_train_loss.append(np.mean(solvers[ws].loss_history[-100:]))
14    bn_final_train_loss.append(np.mean(bn_solvers[ws].loss_history[-100:]))
15
16 plt.subplot(3, 1, 1)
17 plt.title('Best val accuracy vs weight initialization scale')
18 plt.xlabel('Weight initialization scale')
19 plt.ylabel('Best val accuracy')
20 plt.semilogx(weight_scales, best_val_accs, '-o', label='baseline')
21 plt.semilogx(weight_scales, bn_best_val_accs, '-o', label='batchnorm')
22 plt.legend(ncol=2, loc='lower right')
23
24 plt.subplot(3, 1, 2)
25 plt.title('Best train accuracy vs weight initialization scale')
26 plt.xlabel('Weight initialization scale')
27 plt.ylabel('Best training accuracy')
28 plt.semilogx(weight_scales, best_train_accs, '-o', label='baseline')
29 plt.semilogx(weight_scales, bn_best_train_accs, '-o', label='batchnorm')
30 plt.legend()
31
32 plt.subplot(3, 1, 3)
33 plt.title('Final training loss vs weight initialization scale')
34 plt.xlabel('Weight initialization scale')
35 plt.ylabel('Final training loss')
36 plt.semilogx(weight_scales, final_train_loss, '-o', label='baseline')
37 plt.semilogx(weight_scales, bn_final_train_loss, '-o', label='batchnorm')
38 plt.legend()
39
40 plt.gcf().set_size_inches(10, 15)
41 plt.show()

```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "orientation" which is no longer supported as of 3.3 and will become an error two minor releases later

```
fig.canvas.print_figure(bytes_io, **kw)
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "dpi" which is no longer supported as of 3.3 and will become an error two minor releases later

```
fig.canvas.print_figure(bytes_io, **kw)
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "facecolor" which is no longer supported as of 3.3 and will become an error two minor releases later

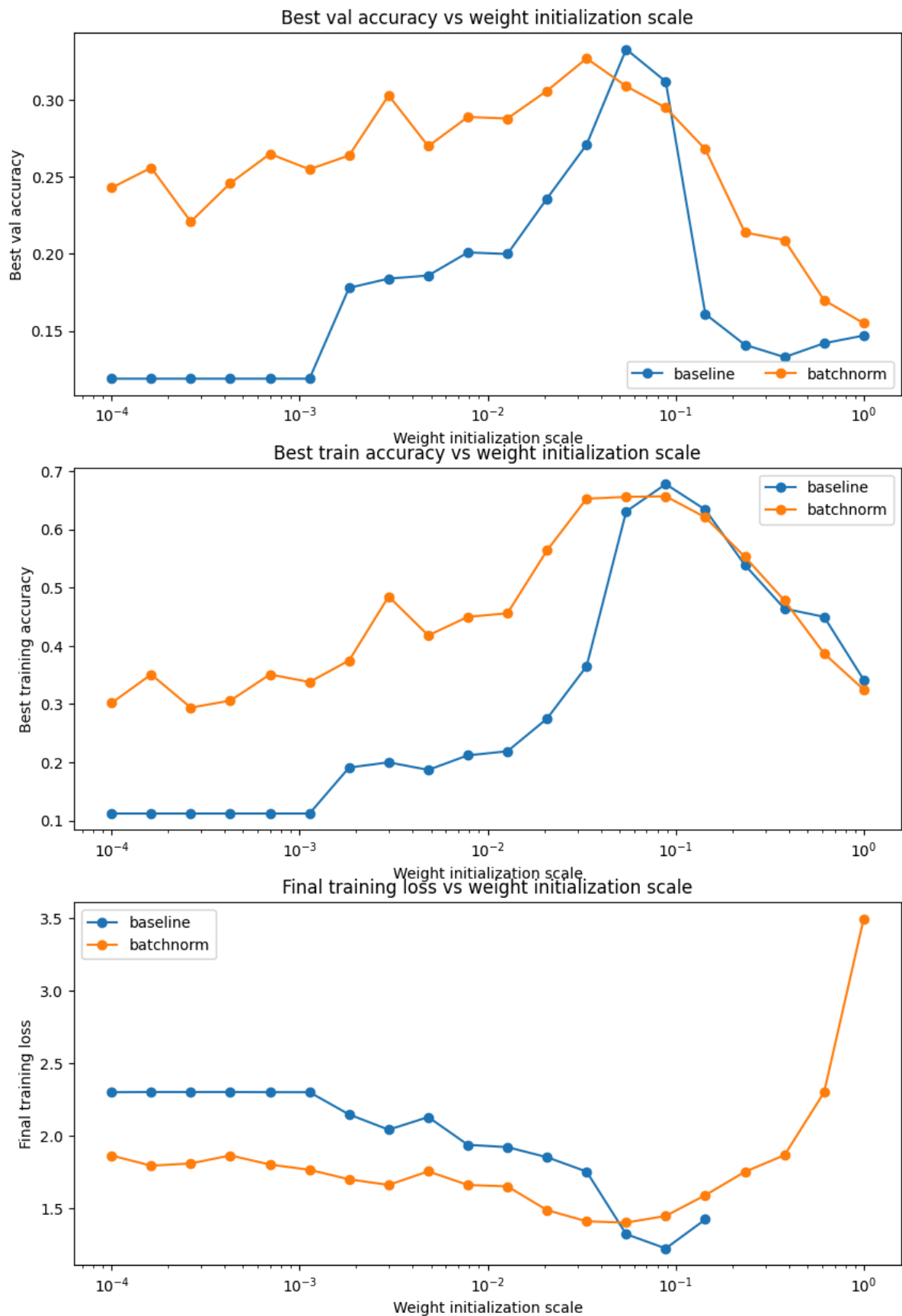
```
fig.canvas.print_figure(bytes_io, **kw)
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "edgecolor" which is no longer supported as of 3.3 and will become an error two minor releases later

```
fig.canvas.print_figure(bytes_io, **kw)
```

/opt/homebrew/lib/python3.11/site-packages/IPython/core/pylabtools.py:152: MatplotlibDeprecationWarning: save fig() got unexpected keyword argument "bbox_inches_restore" which is no longer supported as of 3.3 and will become an error two minor releases later

```
fig.canvas.print_figure(bytes_io, **kw)
```



Question:

In the cell below, summarize the findings of this experiment, and WHY these results make sense.

Answer:

- Upon observing the plots of training and testing error versus weight initialization, we see that the model with batch normalization generally achieves higher accuracy compared to the model without it across a wide range of weight initializations.
- This improvement occurs because batch normalization helps mitigate the internal covariate shift, leading to more stable activations during training
- Additionally, batch normalization acts as a form of regularization, reducing overfitting and enhancing generalization, which results in better test performance.
- Batch normalization also introduces learnable parameters, β (beta) and γ (gamma), allowing the model to scale and shift the inputs to each layer. This contributes to improved optimization and convergence, thereby increasing accuracy in both training and testing.
- In the third graph, we observe that batch normalization minimizes the final training loss across a wide range of weight initializations.
- This suggests that the model is less sensitive to weight initialization and exhibits greater stability during training.

Type *Markdown* and LaTeX: α^2