

Optimization for Fully Connected Networks

In this notebook, we will implement different optimization rules for gradient descent. We have provided starter code; however, you will need to copy and paste your code from your implementation of the modular fully connected nets in HW #3 to build upon this.

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

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

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

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

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

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


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

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

X_train: (49000, 3, 32, 32)
y_train: (49000,)
X_val: (1000, 3, 32, 32)
y_val: (1000,)
X_test: (1000, 3, 32, 32)
y_test: (1000,)
```

Building upon your HW #3 implementation

Copy and paste the following functions from your HW #3 implementation of a modular FC net:

- `affine_forward` in `nndl/layers.py`
- `affine_backward` in `nndl/layers.py`

- `relu_forward` in `nndl/layers.py`
- `relu_backward` in `nndl/layers.py`
- `affine_relu_forward` in `nndl/layer_utils.py`
- `affine_relu_backward` in `nndl/layer_utils.py`
- The `FullyConnectedNet` class in `nndl/fc_net.py`

Test all functions you copy and pasted

```
In [ ]: from nndl.layer_tests import *

affine_forward_test(); print('\n')
affine_backward_test(); print('\n')
relu_forward_test(); print('\n')
relu_backward_test(); print('\n')
affine_relu_test(); print('\n')
fc_net_test()
```

If `affine_forward` function is working, difference should be less than 1e-9:
difference: 9.7698500479884e-10

If `affine_backward` is working, error should be less than 1e-9::
dx error: 3.6393407187524124e-10
dw error: 9.309579801860401e-10
db error: 3.761852417915519e-11

If `relu_forward` function is working, difference should be around 1e-8:
difference: 4.999999798022158e-08

If `relu_forward` function is working, error should be less than 1e-9:
dx error: 3.2756264683652314e-12

If `affine_relu_forward` and `affine_relu_backward` are working, error should be less than 1e-9::
dx error: 1.1768547642195208e-10
dw error: 1.1334531161941464e-09
db error: 3.2755880437749807e-12

Running check with reg = 0
Initial loss: 2.3052013069776676
W1 relative error: 5.692714102328828e-07
W2 relative error: 1.0681747892166097e-05
W3 relative error: 3.122515137335495e-06
b1 relative error: 7.197229320847002e-08
b2 relative error: 2.0633537168458619e-07
b3 relative error: 1.0896914116372043e-10
Running check with reg = 3.14
Initial loss: 7.188867799558902
W1 relative error: 6.486678513821688e-07
W2 relative error: 9.801203239989081e-08
W3 relative error: 3.706387821098089e-08
b1 relative error: 2.186394840150542e-08
b2 relative error: 1.8510824966348172e-08
b3 relative error: 1.4545751157246795e-10

Training a larger model

In general, proceeding with vanilla stochastic gradient descent to optimize models may be fraught with problems and limitations, as discussed in class. Thus, we implement optimizers that improve on SGD.

SGD + momentum

In the following section, implement SGD with momentum. Read the `nndl/optim.py` API, which is provided by CS231n, and be sure you understand it. After, implement `sgd_momentum` in `nndl/optim.py`. Test your implementation of `sgd_momentum` by running the cell below.

```
In [ ]: from nndl.optim import sgd_momentum

N, D = 4, 5
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
v = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-3, 'velocity': v}
next_w, _ = sgd_momentum(w, dw, config=config)

expected_next_w = np.asarray([
    [ 0.1406,       0.20738947,   0.27417895,   0.34096842,   0.40775789],
    [ 0.47454737,   0.54133684,   0.60812632,   0.67491579,   0.74170526],
    [ 0.80849474,   0.87528421,   0.94207368,   1.00886316,   1.07565263],
    [ 1.14244211,   1.20923158,   1.27602105,   1.34281053,   1.4096      ]])
expected_velocity = np.asarray([
    [ 0.5406,       0.55475789,   0.56891579,   0.58307368,   0.59723158],
    [ 0.61138947,   0.62554737,   0.63970526,   0.65386316,   0.66802105],
    [ 0.68217895,   0.69633684,   0.71049474,   0.72465263,   0.73881053],
    [ 0.75296842,   0.76712632,   0.78128421,   0.79544211,   0.8096      ]])

print('next_w error: {}'.format(rel_error(next_w, expected_next_w)))
print('velocity error: {}'.format(rel_error(expected_velocity, config['velocity'])))

next_w error: 8.882347033505819e-09
velocity error: 4.269287743278663e-09
```

SGD + Nesterov momentum

Implement `sgd_nesterov_momentum` in `ndl/optim.py`.

```
In [ ]: from nndl.optim import sgd_nesterov_momentum

N, D = 4, 5
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
v = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-3, 'velocity': v}
next_w, _ = sgd_nesterov_momentum(w, dw, config=config)

expected_next_w = np.asarray([
    [0.08714,       0.15246105,   0.21778211,   0.28310316,   0.34842421],
```

```

[ 0.41374526,  0.47906632,  0.54438737,  0.60970842,  0.67502947],
[ 0.74035053,  0.80567158,  0.87099263,  0.93631368,  1.00163474],
[ 1.06695579,  1.13227684,  1.19759789,  1.26291895,  1.32824   ]])
expected_velocity = np.asarray([
[ 0.5406,      0.55475789,  0.56891579,  0.58307368,  0.59723158],
[ 0.61138947,  0.62554737,  0.63970526,  0.65386316,  0.66802105],
[ 0.68217895,  0.69633684,  0.71049474,  0.72465263,  0.73881053],
[ 0.75296842,  0.76712632,  0.78128421,  0.79544211,  0.8096   ]])

print('next_w error: {}'.format(rel_error(next_w, expected_next_w)))
print('velocity error: {}'.format(rel_error(expected_velocity, config['velocity'])))

```

```

next_w error: 1.0875186845081027e-08
velocity error: 4.269287743278663e-09

```

Evaluating SGD, SGD+Momentum, and SGD+NesterovMomentum

Run the following cell to train a 6 layer FC net with SGD, SGD+momentum, and SGD+Nesterov momentum. You should see that SGD+momentum achieves a better loss than SGD, and that SGD+Nesterov momentum achieves a slightly better loss (and training accuracy) than SGD+momentum.

```

In [ ]: num_train = 4000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

solvers = {}

for update_rule in ['sgd', 'sgd_momentum', 'sgd_nesterov_momentum']:
    print('Optimizing with {}'.format(update_rule))
    model = FullyConnectedNet([100, 100, 100, 100, 100], weight_scale=5e-2)

    solver = Solver(model, small_data,
                    num_epochs=5, batch_size=100,
                    update_rule=update_rule,
                    optim_config={
                        'learning_rate': 1e-2,
                    },
                    verbose=False)
    solvers[update_rule] = solver
    solver.train()
    print

fig, axes = plt.subplots(3, 1)

ax = axes[0]
ax.set_title('Training loss')
ax.set_xlabel('Iteration')

ax = axes[1]
ax.set_title('Training accuracy')
ax.set_xlabel('Epoch')

ax = axes[2]

```

```
ax.set_title('Validation accuracy')
ax.set_xlabel('Epoch')

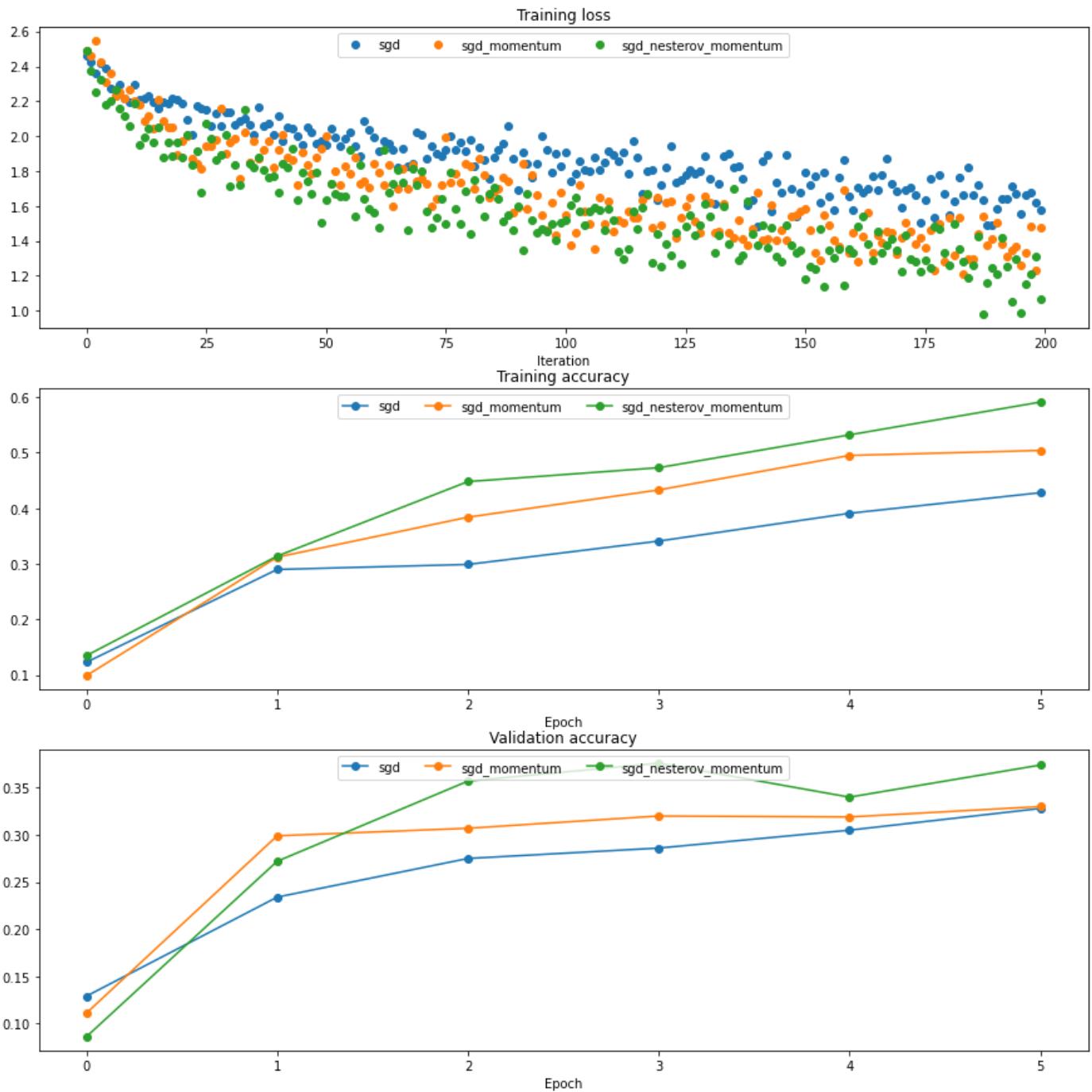
for update_rule, solver in solvers.items():
    ax = axes[0]
    ax.plot(solver.loss_history, 'o', label=update_rule)

    ax = axes[1]
    ax.plot(solver.train_acc_history, '-o', label=update_rule)

    ax = axes[2]
    ax.plot(solver.val_acc_history, '-o', label=update_rule)

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

Optimizing with sgd
Optimizing with sgd_momentum
Optimizing with sgd_nesterov_momentum



RMSProp

Now we go to techniques that adapt the gradient. Implement `rmsprop` in `nndl/optim.py`. Test your implementation by running the cell below.

In []: `from nndl.optim import rmsprop`

```
N, D = 4, 5
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
a = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-2, 'a': a}
next_w, _ = rmsprop(w, dw, config=config)

expected_next_w = np.asarray([
    [ 0.4945,  0.4835, -0.4375, -0.4445],
    [ 0.3675,  0.3565, -0.2995, -0.3065],
    [-0.0285, -0.0175, -0.4035, -0.4105],
    [-0.4475, -0.4365, -0.0995, -0.1065]
])
```

```

[ -0.39223849, -0.34037513, -0.28849239, -0.23659121, -0.18467247],
[ -0.132737, -0.08078555, -0.02881884, 0.02316247, 0.07515774],
[ 0.12716641, 0.17918792, 0.23122175, 0.28326742, 0.33532447],
[ 0.38739248, 0.43947102, 0.49155973, 0.54365823, 0.59576619]]))
expected_cache = np.asarray([
[ 0.5976, 0.6126277, 0.6277108, 0.64284931, 0.65804321],
[ 0.67329252, 0.68859723, 0.70395734, 0.71937285, 0.73484377],
[ 0.75037008, 0.7659518, 0.78158892, 0.79728144, 0.81302936],
[ 0.82883269, 0.84469141, 0.86060554, 0.87657507, 0.8926 ]])

print('next_w error: {}'.format(rel_error(expected_next_w, next_w)))
print('cache error: {}'.format(rel_error(expected_cache, config['a'])))

```

next_w error: 9.502645229894295e-08
cache error: 2.6477955807156126e-09

Adaptive moments

Now, implement `adam` in `nndl/optim.py`. Test your implementation by running the cell below.

```
In [ ]: # Test Adam implementation; you should see errors around 1e-7 or less
from nndl.optim import adam

N, D = 4, 5
w = np.linspace(-0.4, 0.6, num=N*D).reshape(N, D)
dw = np.linspace(-0.6, 0.4, num=N*D).reshape(N, D)
v = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)
a = np.linspace(0.7, 0.5, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-2, 'v': v, 'a': a, 't': 5}
next_w, _ = adam(w, dw, config=config)

expected_next_w = np.asarray([
[ -0.40094747, -0.34836187, -0.29577703, -0.24319299, -0.19060977],
[ -0.1380274, -0.08544591, -0.03286534, 0.01971428, 0.0722929],
[ 0.1248705, 0.17744702, 0.23002243, 0.28259667, 0.33516969],
[ 0.38774145, 0.44031188, 0.49288093, 0.54544852, 0.59801459]]))
expected_a = np.asarray([
[ 0.69966, 0.68908382, 0.67851319, 0.66794809, 0.65738853, ],
[ 0.64683452, 0.63628604, 0.6257431, 0.61520571, 0.60467385, ],
[ 0.59414753, 0.58362676, 0.57311152, 0.56260183, 0.55209767, ],
[ 0.54159906, 0.53110598, 0.52061845, 0.51013645, 0.49966, ]])
expected_v = np.asarray([
[ 0.48, 0.49947368, 0.51894737, 0.53842105, 0.55789474],
[ 0.57736842, 0.59684211, 0.61631579, 0.63578947, 0.65526316],
[ 0.67473684, 0.69421053, 0.71368421, 0.73315789, 0.75263158],
[ 0.77210526, 0.79157895, 0.81105263, 0.83052632, 0.85 ]])

print('next_w error: {}'.format(rel_error(expected_next_w, next_w)))
print('a error: {}'.format(rel_error(expected_a, config['a'])))
print('v error: {}'.format(rel_error(expected_v, config['v'])))
```

next_w error: 1.1395691798535431e-07
a error: 4.208314038113071e-09
v error: 4.214963193114416e-09

Comparing SGD, SGD+NesterovMomentum, RMSProp, and Adam

The following code will compare optimization with SGD, Momentum, Nesterov Momentum, RMSProp and Adam. In our code, we find that RMSProp, Adam, and SGD + Nesterov Momentum achieve approximately the same training error after a few training epochs.

```
In [ ]: learning_rates = {'rmsprop': 2e-4, 'adam': 1e-3}

for update_rule in ['adam', 'rmsprop']:
    print('Optimizing with {}'.format(update_rule))
    model = FullyConnectedNet([100, 100, 100, 100, 100], weight_scale=5e-2)

    solver = Solver(model, small_data,
                    num_epochs=5, batch_size=100,
                    update_rule=update_rule,
                    optim_config={
                        'learning_rate': learning_rates[update_rule]
                    },
                    verbose=False)
    solvers[update_rule] = solver
    solver.train()
    print

fig, axes = plt.subplots(3, 1)

ax = axes[0]
ax.set_title('Training loss')
ax.set_xlabel('Iteration')

ax = axes[1]
ax.set_title('Training accuracy')
ax.set_xlabel('Epoch')

ax = axes[2]
ax.set_title('Validation accuracy')
ax.set_xlabel('Epoch')

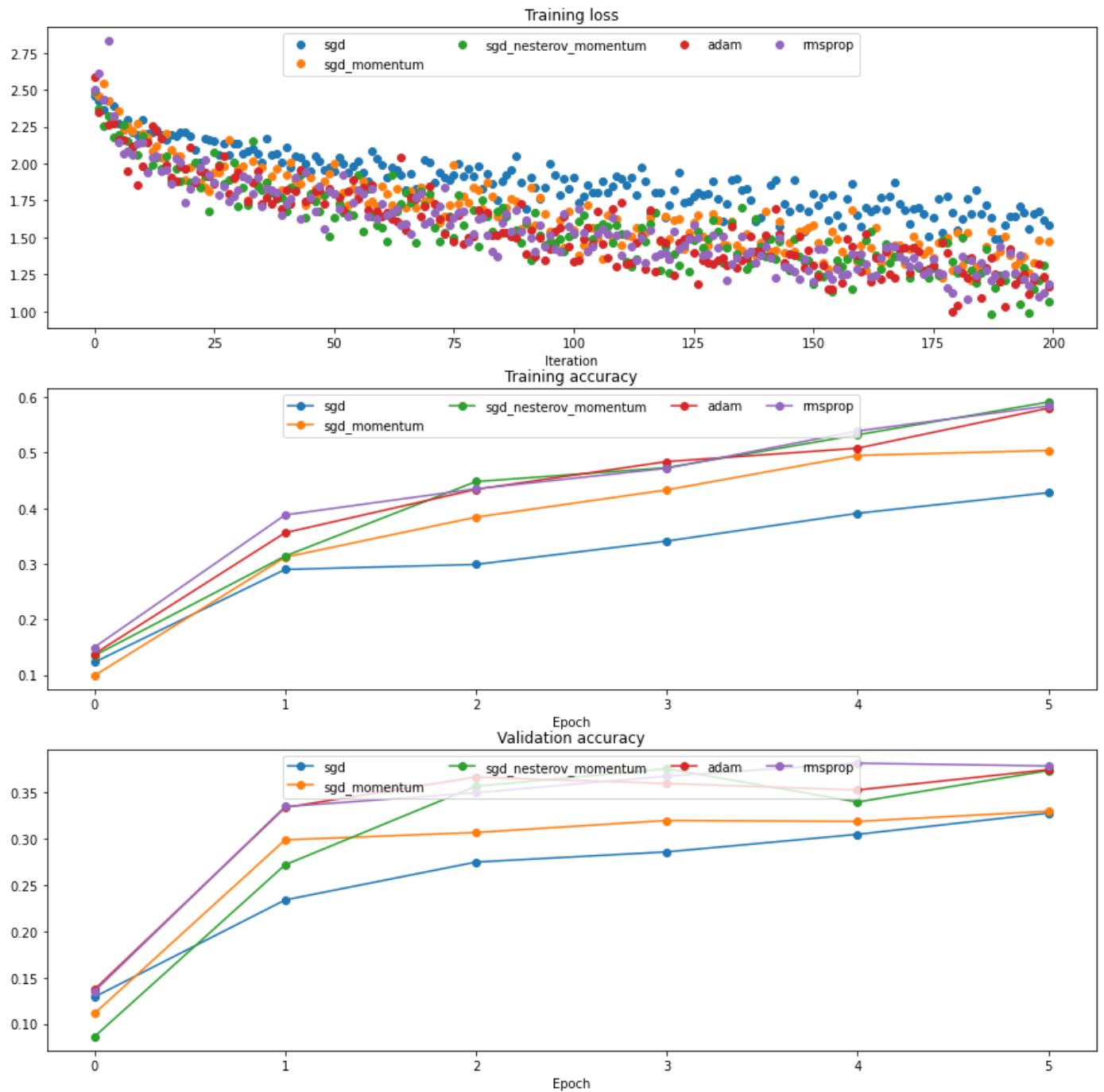
for update_rule, solver in solvers.items():
    ax = axes[0]
    ax.plot(solver.loss_history, 'o', label=update_rule)

    ax = axes[1]
    ax.plot(solver.train_acc_history, '-o', label=update_rule)

    ax = axes[2]
    ax.plot(solver.val_acc_history, '-o', label=update_rule)

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

Optimizing with adam
Optimizing with rmsprop
```



Easier optimization

In the following cell, we'll train a 4 layer neural network having 500 units in each hidden layer with the different optimizers, and find that it is far easier to get up to 50+% performance on CIFAR-10. After we implement batchnorm and dropout, we'll ask you to get 55+% on CIFAR-10.

In []:

```
optimizer = 'adam'
best_model = None

layer_dims = [500, 500, 500]
weight_scale = 0.01
learning_rate = 1e-3
lr_decay = 0.9

model = FullyConnectedNet(layer_dims, weight_scale=weight_scale,
                         use_batchnorm=True)
```

```
solver = Solver(model, data,
                num_epochs=10, batch_size=100,
                update_rule=optimizer,
                optim_config={
                    'learning_rate': learning_rate,
                },
                lr_decay=lr_decay,
                verbose=True, print_every=50)
solver.train()
```

```
(Iteration 1 / 4900) loss: 2.332614
(Epoch 0 / 10) train acc: 0.120000; val_acc: 0.093000
(Iteration 51 / 4900) loss: 2.069911
(Iteration 101 / 4900) loss: 1.748196
(Iteration 151 / 4900) loss: 1.673191
(Iteration 201 / 4900) loss: 1.706806
(Iteration 251 / 4900) loss: 1.568937
(Iteration 301 / 4900) loss: 1.562796
(Iteration 351 / 4900) loss: 1.719840
(Iteration 401 / 4900) loss: 1.615221
(Iteration 451 / 4900) loss: 1.385631
(Epoch 1 / 10) train acc: 0.438000; val_acc: 0.440000
(Iteration 501 / 4900) loss: 1.357284
(Iteration 551 / 4900) loss: 1.602889
(Iteration 601 / 4900) loss: 1.527623
(Iteration 651 / 4900) loss: 1.498156
(Iteration 701 / 4900) loss: 1.541310
(Iteration 751 / 4900) loss: 1.417537
(Iteration 801 / 4900) loss: 1.567270
(Iteration 851 / 4900) loss: 1.458091
(Iteration 901 / 4900) loss: 1.404160
(Iteration 951 / 4900) loss: 1.284027
(Epoch 2 / 10) train acc: 0.528000; val_acc: 0.492000
(Iteration 1001 / 4900) loss: 1.275148
(Iteration 1051 / 4900) loss: 1.216280
(Iteration 1101 / 4900) loss: 1.369018
(Iteration 1151 / 4900) loss: 1.245002
(Iteration 1201 / 4900) loss: 1.304045
(Iteration 1251 / 4900) loss: 1.187996
(Iteration 1301 / 4900) loss: 1.378438
(Iteration 1351 / 4900) loss: 1.257191
(Iteration 1401 / 4900) loss: 1.263489
(Iteration 1451 / 4900) loss: 1.503795
(Epoch 3 / 10) train acc: 0.548000; val_acc: 0.514000
(Iteration 1501 / 4900) loss: 1.284813
(Iteration 1551 / 4900) loss: 1.231590
(Iteration 1601 / 4900) loss: 1.346853
(Iteration 1651 / 4900) loss: 1.133964
(Iteration 1701 / 4900) loss: 1.118430
(Iteration 1751 / 4900) loss: 1.250812
(Iteration 1801 / 4900) loss: 1.387983
(Iteration 1851 / 4900) loss: 1.135359
(Iteration 1901 / 4900) loss: 1.390784
(Iteration 1951 / 4900) loss: 1.344122
(Epoch 4 / 10) train acc: 0.598000; val_acc: 0.518000
(Iteration 2001 / 4900) loss: 1.197539
(Iteration 2051 / 4900) loss: 1.203903
(Iteration 2101 / 4900) loss: 1.212395
(Iteration 2151 / 4900) loss: 1.128789
(Iteration 2201 / 4900) loss: 1.102054
(Iteration 2251 / 4900) loss: 0.881218
(Iteration 2301 / 4900) loss: 1.170962
(Iteration 2351 / 4900) loss: 1.494178
(Iteration 2401 / 4900) loss: 1.093521
(Epoch 5 / 10) train acc: 0.617000; val_acc: 0.518000
(Iteration 2451 / 4900) loss: 1.205427
(Iteration 2501 / 4900) loss: 1.154193
(Iteration 2551 / 4900) loss: 1.290390
(Iteration 2601 / 4900) loss: 1.162692
(Iteration 2651 / 4900) loss: 1.302425
(Iteration 2701 / 4900) loss: 1.166759
(Iteration 2751 / 4900) loss: 1.148194
```

```
(Iteration 2801 / 4900) loss: 1.097889
(Iteration 2851 / 4900) loss: 1.001535
(Iteration 2901 / 4900) loss: 1.125101
(Epoch 6 / 10) train acc: 0.616000; val_acc: 0.527000
(Iteration 2951 / 4900) loss: 0.950356
(Iteration 3001 / 4900) loss: 0.956059
(Iteration 3051 / 4900) loss: 1.053148
(Iteration 3101 / 4900) loss: 1.278789
(Iteration 3151 / 4900) loss: 1.074644
(Iteration 3201 / 4900) loss: 0.987589
(Iteration 3251 / 4900) loss: 1.018784
(Iteration 3301 / 4900) loss: 1.037282
(Iteration 3351 / 4900) loss: 0.955591
(Iteration 3401 / 4900) loss: 0.725424
(Epoch 7 / 10) train acc: 0.639000; val_acc: 0.528000
(Iteration 3451 / 4900) loss: 1.179690
(Iteration 3501 / 4900) loss: 0.814921
(Iteration 3551 / 4900) loss: 0.957090
(Iteration 3601 / 4900) loss: 0.842116
(Iteration 3651 / 4900) loss: 1.089793
(Iteration 3701 / 4900) loss: 0.791398
(Iteration 3751 / 4900) loss: 1.047696
(Iteration 3801 / 4900) loss: 0.771098
(Iteration 3851 / 4900) loss: 1.053682
(Iteration 3901 / 4900) loss: 0.818712
(Epoch 8 / 10) train acc: 0.703000; val_acc: 0.511000
(Iteration 3951 / 4900) loss: 1.060274
(Iteration 4001 / 4900) loss: 0.932583
(Iteration 4051 / 4900) loss: 0.951160
(Iteration 4101 / 4900) loss: 0.969140
(Iteration 4151 / 4900) loss: 0.993867
(Iteration 4201 / 4900) loss: 1.104124
(Iteration 4251 / 4900) loss: 0.682630
(Iteration 4301 / 4900) loss: 1.023785
(Iteration 4351 / 4900) loss: 0.815212
(Iteration 4401 / 4900) loss: 0.956450
(Epoch 9 / 10) train acc: 0.685000; val_acc: 0.532000
(Iteration 4451 / 4900) loss: 0.950370
(Iteration 4501 / 4900) loss: 0.744691
(Iteration 4551 / 4900) loss: 0.926480
(Iteration 4601 / 4900) loss: 0.827626
(Iteration 4651 / 4900) loss: 0.673550
(Iteration 4701 / 4900) loss: 0.962829
(Iteration 4751 / 4900) loss: 0.836745
(Iteration 4801 / 4900) loss: 0.856362
(Iteration 4851 / 4900) loss: 0.838408
(Epoch 10 / 10) train acc: 0.707000; val_acc: 0.539000
```

```
In [ ]: y_test_pred = np.argmax(model.loss(data['X_test']), axis=1)
y_val_pred = np.argmax(model.loss(data['X_val']), axis=1)
print('Validation set accuracy: {}'.format(np.mean(y_val_pred == data['y_val'])))
print('Test set accuracy: {}'.format(np.mean(y_test_pred == data['y_test'])))

Validation set accuracy: 0.539
Test set accuracy: 0.541
```

```
In [ ]:
```

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

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

Inputs:

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

Returns:

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

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

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

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

```
"""

```

Performs vanilla stochastic gradient descent.

config format:

- learning_rate: Scalar learning rate.

```
"""

```

```
if config is None: config = {}
```

```

config.setdefault('learning_rate', 1e-2)

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

def sgd_momentum(w, dw, config=None):
    """
    Performs stochastic gradient descent with momentum.

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

    # =====#
    # YOUR CODE HERE:
    # Implement the momentum update formula. Return the updated
    weights
    # as next_w, and the updated velocity as v.
    # =====#
    v = config['momentum'] * v - config['learning_rate'] * dw
    next_w = w + v
    # next_w = np.linalg.norm(w - v[-1])

    # =====#
    # END YOUR CODE HERE
    # =====#
    config['velocity'] = v

    return next_w, config

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

```

Performs stochastic gradient descent with Nesterov momentum.

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

def rmsprop(w, dw, config=None):
 """
 Uses the RMSProp update rule, which uses a moving average of
 squared gradient
 values to set adaptive per-parameter learning rates.

 config format:
 - learning_rate: Scalar learning rate.
 - decay_rate: Scalar between 0 and 1 giving the decay rate for the

```

squared
    gradient cache.
    - epsilon: Small scalar used for smoothing to avoid dividing by
zero.
    - beta: Moving average of second moments of gradients.
"""
if config is None: config = {}
config.setdefault('learning_rate', 1e-2)
config.setdefault('decay_rate', 0.99)
config.setdefault('epsilon', 1e-8)
config.setdefault('a', np.zeros_like(w))

next_w = None

# =====#
# YOUR CODE HERE:
# Implement RMSProp. Store the next value of w as next_w. You
need
# to also store in config['a'] the moving average of the second
# moment gradients, so they can be used for future gradients.
Concretely,
# config['a'] corresponds to "a" in the lecture notes.
# =====#
# beta = 0.99
# config['a'] = beta * config['a'] + (1 - beta) * dw * dw
# next_w = w - config['learning_rate'] * dw / (np.sqrt(config['a'] +
# config['epsilon']))

# =====#
# END YOUR CODE HERE
# =====#
# return next_w, config

def adam(w, dw, config=None):
"""
    Uses the Adam update rule, which incorporates moving averages of
both the
gradient and its square and a bias correction term.

    config format:
    - learning_rate: Scalar learning rate.
    - beta1: Decay rate for moving average of first moment of
gradient.

```

```

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

    next_w = None

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

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

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

    return next_w, config

```

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.

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

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

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

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

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

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


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

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

```
X_train: (49000, 3, 32, 32)
y_train: (49000,)
X_val: (1000, 3, 32, 32)
y_val: (1000,)
X_test: (1000, 3, 32, 32)
y_test: (1000,)
```

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

# Simulate the forward pass for a two-layer network
N, D1, D2, D3 = 200, 50, 60, 3
X = np.random.randn(N, D1)
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)
a = np.maximum(0, X.dot(W1)).dot(W2)

print('Before batch normalization:')
print(' means: ', a.mean(axis=0))
print(' stds: ', a.std(axis=0))

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

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

```
Before batch normalization:
means: [-0.6742318  0.80407948  3.95415552]
stds: [30.51327184 26.81786846 36.13805553]
After batch normalization (gamma=1, beta=0)
mean: [-4.32986980e-17 -4.20843915e-17 -3.38618023e-17]
std: [0.99999999 0.99999999 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 [ ]: # Check the test-time forward pass by running the training-time
# forward pass many times to warm up the running averages, and then
# checking the means and variances of activations after a test-time
# forward pass.

N, D1, D2, D3 = 200, 50, 60, 3
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)

bn_param = {'mode': 'train'}
gamma = np.ones(D3)
beta = np.zeros(D3)
for t in np.arange(50):
    X = np.random.randn(N, D1)
```

```

a = np.maximum(0, X.dot(W1)).dot(W2)
batchnorm_forward(a, gamma, beta, bn_param)
bn_param['mode'] = 'test'
X = np.random.randn(N, D1)
a = np.maximum(0, X.dot(W1)).dot(W2)
a_norm, _ = batchnorm_forward(a, gamma, beta, bn_param)

# Means should be close to zero and stds close to one, but will be
# noisier than training-time forward passes.
print('After batch normalization (test-time):')
print(' means: ', a_norm.mean(axis=0))
print(' stds: ', a_norm.std(axis=0))

```

```

After batch normalization (test-time):
means: [-0.02747245 -0.01851219 -0.07292147]
stds: [1.07141508 0.916886 0.92060113]

```

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 [ ]: # Gradient check batchnorm backward pass

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

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

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

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

# print("--")
# print("dout", dout)
# print("da", da_num)
# print("dgamma", dgamma)
# print("--")
# print("db", db_num)
# print("dbeta", dbeta)
# print("--")
# print("dx check", dx_num)
# print("dx", dx)

```

```
dx error: 3.4477536046741704e-08
dgamma error: 8.978965955969402e-12
dbeta error: 3.3104710507890927e-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 [ ]: N, D, H1, H2, C = 2, 15, 20, 30, 10
X = np.random.randn(N, D)
y = np.random.randint(C, size=(N,))

for reg in [0, 3.14]:
    print('Running check with reg = ', reg)
    model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
                             reg=reg, weight_scale=5e-2, dtype=np.float64,
                             use_batchnorm=True)

    loss, grads = model.loss(X, y)
    print('Initial loss: ', loss)

    for name in sorted(grads):
        f = lambda _: model.loss(X, y)[0]
        grad_num = eval_numerical_gradient(f, model.params[name], verbose=False)
        print('{} relative error: {}'.format(name, rel_error(grad_num, grads[name])))
    if reg == 0: print('\n')
```

```
Running check with reg = 0
Initial loss: 2.3710490024211275
W1 relative error: 0.00015101419349347454
W2 relative error: 3.672777870998256e-05
W3 relative error: 2.443181439414118e-09
b1 relative error: 0.0022203572314083426
b2 relative error: 1.1102230246251565e-08
b3 relative error: 8.523086614559495e-11
beta1 relative error: 5.446086592476153e-09
beta2 relative error: 8.75544283248942e-08
gamma1 relative error: 5.436763679722301e-09
gamma2 relative error: 5.354138682022913e-08
```

```
Running check with reg = 3.14
Initial loss: 7.192403210139613
W1 relative error: 8.459596789798366e-07
W2 relative error: 8.081855100065242e-06
W3 relative error: 1.6006226075338707e-08
b1 relative error: 2.7755575615628914e-09
b2 relative error: 2.220446049250313e-08
b3 relative error: 2.726096590945677e-10
beta1 relative error: 7.750597039229821e-09
beta2 relative error: 2.395750124612298e-08
gamma1 relative error: 2.1389851900663752e-08
gamma2 relative error: 9.288900409058029e-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 [ ]: # Try training a very deep net with batchnorm
hidden_dims = [100, 100, 100, 100, 100]

num_train = 1000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

weight_scale = 2e-2
bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=False)
model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=True)

bn_solver = Solver(bn_model, small_data,
                  num_epochs=10, batch_size=50,
                  update_rule='adam',
                  optim_config={
                      'learning_rate': 1e-3,
                  },
                  verbose=True, print_every=200)
bn_solver.train()
```

```

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

(Iteration 1 / 200) loss: 2.295101
(Epoch 0 / 10) train acc: 0.116000; val_acc: 0.123000
(Epoch 1 / 10) train acc: 0.278000; val_acc: 0.240000
(Epoch 2 / 10) train acc: 0.408000; val_acc: 0.328000
(Epoch 3 / 10) train acc: 0.483000; val_acc: 0.308000
(Epoch 4 / 10) train acc: 0.543000; val_acc: 0.339000
(Epoch 5 / 10) train acc: 0.586000; val_acc: 0.320000
(Epoch 6 / 10) train acc: 0.644000; val_acc: 0.342000
(Epoch 7 / 10) train acc: 0.733000; val_acc: 0.335000
(Epoch 8 / 10) train acc: 0.780000; val_acc: 0.353000
(Epoch 9 / 10) train acc: 0.827000; val_acc: 0.343000
(Epoch 10 / 10) train acc: 0.829000; val_acc: 0.331000
(Iteration 1 / 200) loss: 2.302552
(Epoch 0 / 10) train acc: 0.108000; val_acc: 0.083000
(Epoch 1 / 10) train acc: 0.241000; val_acc: 0.215000
(Epoch 2 / 10) train acc: 0.278000; val_acc: 0.243000
(Epoch 3 / 10) train acc: 0.333000; val_acc: 0.282000
(Epoch 4 / 10) train acc: 0.354000; val_acc: 0.276000
(Epoch 5 / 10) train acc: 0.369000; val_acc: 0.270000
(Epoch 6 / 10) train acc: 0.368000; val_acc: 0.267000
(Epoch 7 / 10) train acc: 0.430000; val_acc: 0.307000
(Epoch 8 / 10) train acc: 0.443000; val_acc: 0.312000
(Epoch 9 / 10) train acc: 0.552000; val_acc: 0.309000
(Epoch 10 / 10) train acc: 0.583000; val_acc: 0.299000

```

```

In [ ]: fig, axes = plt.subplots(3, 1)

ax = axes[0]
ax.set_title('Training loss')
ax.set_xlabel('Iteration')

ax = axes[1]
ax.set_title('Training accuracy')
ax.set_xlabel('Epoch')

ax = axes[2]
ax.set_title('Validation accuracy')
ax.set_xlabel('Epoch')

ax = axes[0]
ax.plot(solver.loss_history, 'o', label='baseline')
ax.plot(bn_solver.loss_history, 'o', label='batchnorm')

ax = axes[1]
ax.plot(solver.train_acc_history, '-o', label='baseline')
ax.plot(bn_solver.train_acc_history, '-o', label='batchnorm')

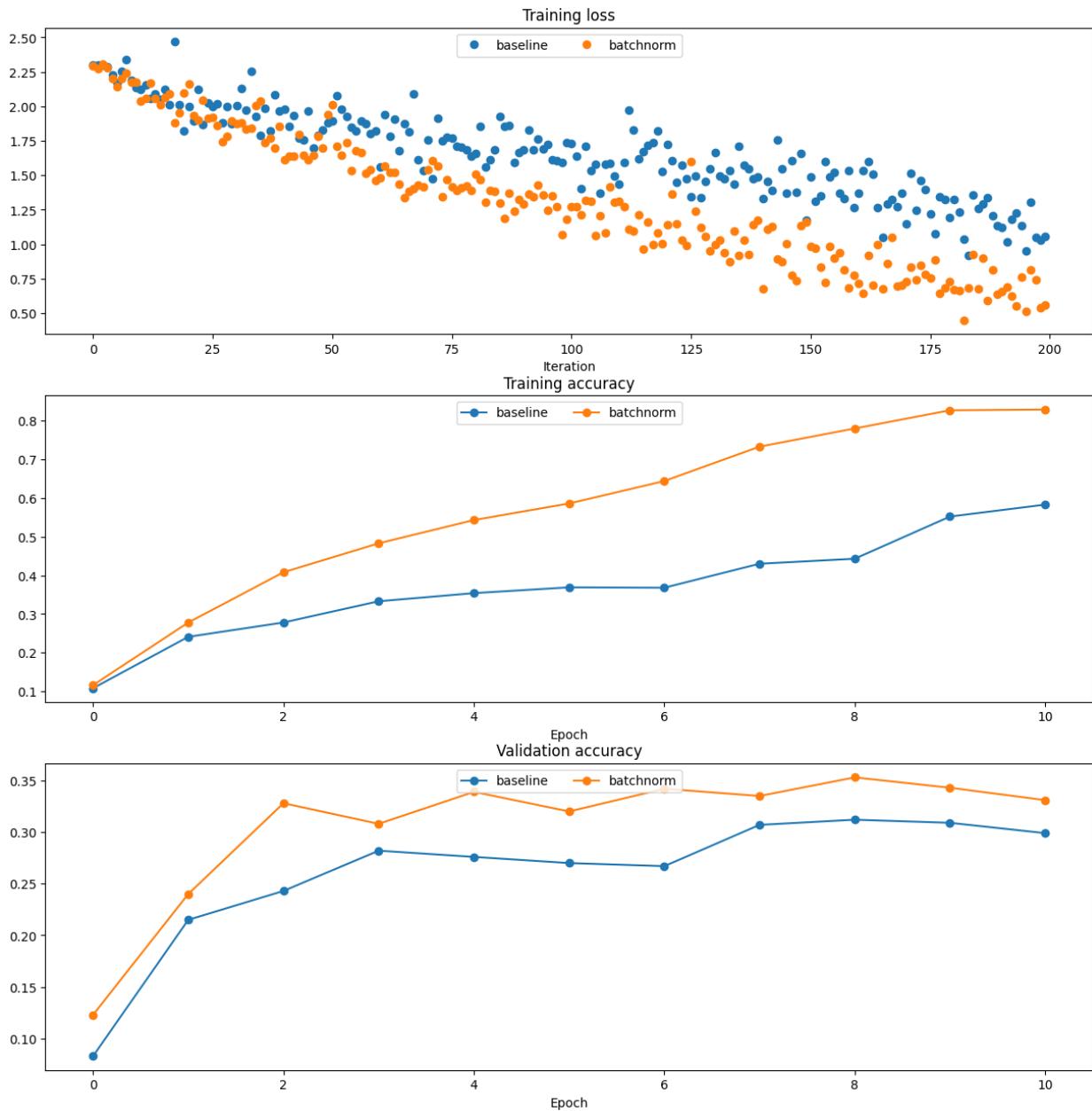
ax = axes[2]
ax.plot(solver.val_acc_history, '-o', label='baseline')
ax.plot(bn_solver.val_acc_history, '-o', label='batchnorm')

```

```

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

```



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 [ ]: # Try training a very deep net with batchnorm
hidden_dims = [50, 50, 50, 50, 50, 50, 50]

num_train = 1000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
}

```

```

'x_val': data['x_val'],
'y_val': data['y_val'],
}

bn_solvers = {}
solvers = {}
weight_scales = np.logspace(-4, 0, num=20)
for i, weight_scale in enumerate(weight_scales):
    print('Running weight scale {} / {}'.format(i + 1, len(weight_scales)))
    bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batch_norm=False)
    model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batch_norm=True)

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

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

```

Running weight scale 1 / 20
 Running weight scale 2 / 20
 Running weight scale 3 / 20
 Running weight scale 4 / 20
 Running weight scale 5 / 20
 Running weight scale 6 / 20
 Running weight scale 7 / 20
 Running weight scale 8 / 20
 Running weight scale 9 / 20
 Running weight scale 10 / 20
 Running weight scale 11 / 20
 Running weight scale 12 / 20
 Running weight scale 13 / 20
 Running weight scale 14 / 20
 Running weight scale 15 / 20
 Running weight scale 16 / 20
 Running weight scale 17 / 20
 Running weight scale 18 / 20
 Running weight scale 19 / 20
 Running weight scale 20 / 20

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

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

```
bn_best_train_accs.append(max(bn_solvers[ws].train_acc_history))

best_val_accs.append(max(solvers[ws].val_acc_history))
bn_best_val_accs.append(max(bn_solvers[ws].val_acc_history))

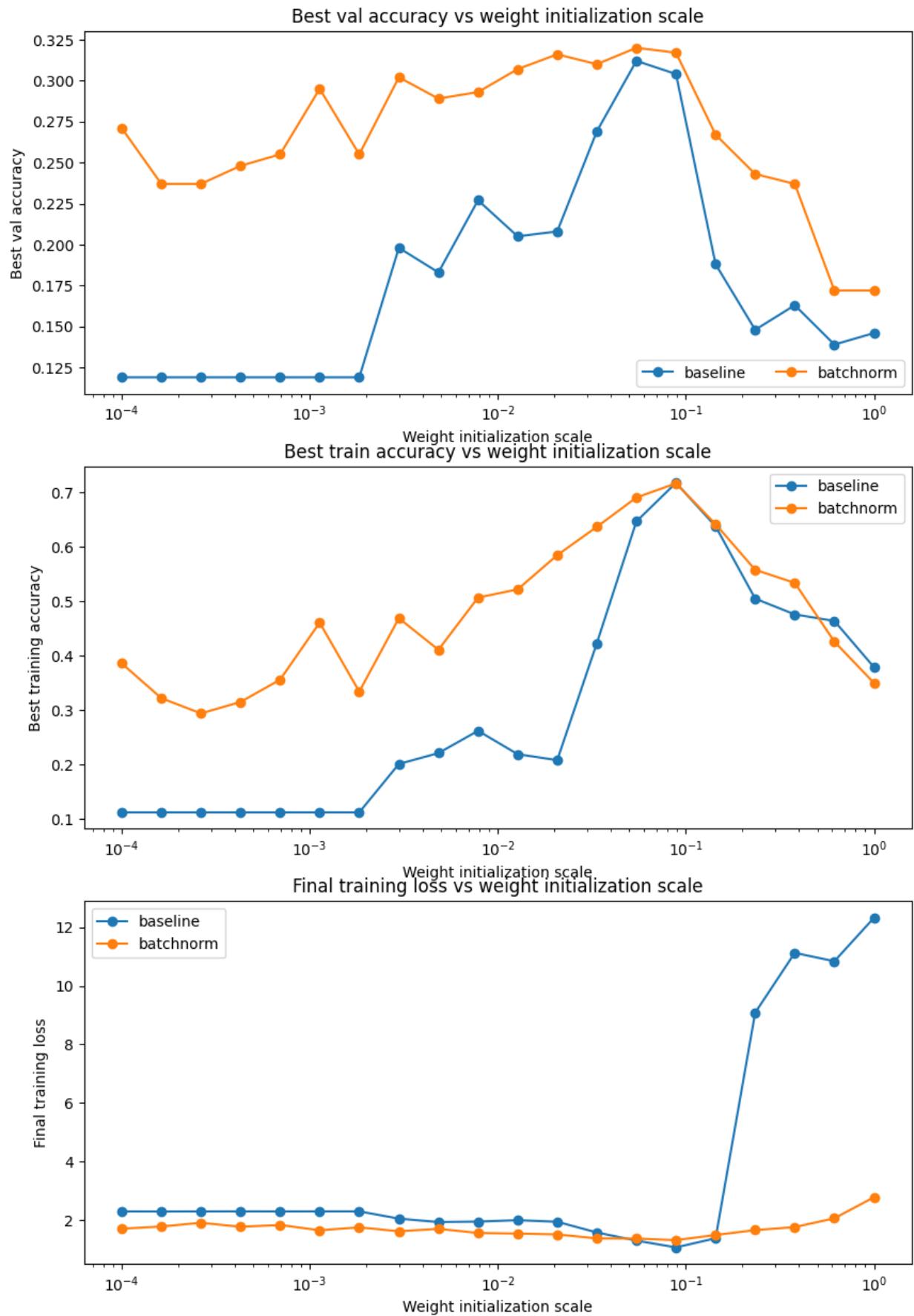
final_train_loss.append(np.mean(solvers[ws].loss_history[-100:]))
bn_final_train_loss.append(np.mean(bn_solvers[ws].loss_history[-100:]))

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

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

plt.subplot(3, 1, 3)
plt.title('Final training loss vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Final training loss')
plt.semilogx(weight_scales, final_train_loss, '-o', label='baseline')
plt.semilogx(weight_scales, bn_final_train_loss, '-o', label='batchnorm')
plt.legend()

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



Question:

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

Answer:

For both batchnorm and non-batchnorm, the best train and validation accuracy are at weight initialization scale 10e-1.

For the final training, there is a slight decrease along the graph and 10e-1 is the best for both batchnorm and non-batchnorm. For the non-batchnorm (baseline) in particular, loss skyrockets up for weight initialization scale > 10e-1.

These results make sense because when the weight scale is 10e-1, this reduces the problem of disappearing gradients. This is similar to the Xavier weight initialization of $-(1/\sqrt{n})$ and $1/\sqrt{n}$, which makes sense because this is the optimal initialization strategy. This strategy helps mitigate the exploding/vanishing gradients problem.

In []:

```

import numpy as np
import pdb

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

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

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

    Returns a tuple of:
    - out: output, of shape (N, M)
    - cache: (x, w, b)
    """
    out = None
    # =====#
    #
    # YOUR CODE HERE:
    #   Calculate the output of the forward pass. Notice the
    dimensions
    #   of w are D x M, which is the transpose of what we did in
earlier
    #   assignments.
    # =====#
    #

    xr = x.reshape(x.shape[0], -1)
    out = xr.dot(w) + b

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

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

```

```

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

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

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

    # =====#
    # YOUR CODE HERE:
    # Calculate the gradients for the backward pass.
    # Notice:
    #   dout is N x M
    #   dx should be N x d1 x ... x dk; it relates to dout through
    #   multiplication with w, which is D x M
    #   dw should be D x M; it relates to dout through multiplication
    #   with x, which is N x D after reshaping
    #   db should be M; it is just the sum over dout examples
    # =====#
    # dx = np.dot(dout, w.T)
    # dx = dx.reshape(x.shape)
    # xr = x.reshape(x.shape[0], -1)
    # dw = np.dot(xr.T, dout)
    # db = np.sum(dout, axis=0)

    # =====#
    # END YOUR CODE HERE
    # =====#
    # return dx, dw, db

```

```

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

    Input:
    - x: Inputs, of any shape

    Returns a tuple of:
    - out: Output, of the same shape as x
    - cache: x
    """
    # =====
    #
    # YOUR CODE HERE:
    # Implement the ReLU forward pass.
    # =====
    #

    f = lambda x: x * (x > 0)
    out = f(x)
    # =====
    #
    # END YOUR CODE HERE
    # =====
    #

    cache = x
    return out, cache


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

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

    Returns:
    - dx: Gradient with respect to x
    """
    x = cache

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

```

```

#
# dx = dout * (x > 0)
#
# =====
#
# END YOUR CODE HERE
# =====
#
return dx

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

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

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

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

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

Input:
- x: Data of shape (N, D)
- gamma: Scale parameter of shape (D,)
- beta: Shift parameter of shape (D,)
- bn_param: Dictionary with the following keys:

```

```

    - mode: 'train' or 'test'; required
    - eps: Constant for numeric stability
    - momentum: Constant for running mean / variance.
    - running_mean: Array of shape (D,) giving running mean of
features
        - running_var Array of shape (D,) giving running variance of
features

    Returns a tuple of:
    - out: of shape (N, D)
    - cache: A tuple of values needed in the backward pass
    """
mode = bn_param['mode']
eps = bn_param.get('eps', 1e-5)
momentum = bn_param.get('momentum', 0.9)

N, D = x.shape
running_mean = bn_param.get('running_mean', np.zeros(D,
dtype=x.dtype))
running_var = bn_param.get('running_var', np.zeros(D,
dtype=x.dtype))

out, cache = None, None
if mode == 'train':

    #
===== # YOUR CODE HERE:
    # A few steps here:
    # (1) Calculate the running mean and variance of the
minibatch.
    # (2) Normalize the activations with the running mean and
variance.
    # (3) Scale and shift the normalized activations. Store
this
    # as the variable 'out'
    # (4) Store any variables you may need for the backward
pass in
    # the 'cache' variable.
    #
===== #

    # 1 calculate running mean and variance
    # calculate sample mean and var
    mun = np.mean(x, axis=0)
    var = np.var(x, axis=0)

    # update running
    running_mean = momentum * running_mean + (1 - momentum) * mun
    running_var = momentum * running_var + (1 - momentum) * var

```

```

# 2 normalize activations
xc = x - mun
x_hat = (x - mun)/(np.sqrt(var + eps))

# 3 scale and shift the normalized activations
out = gamma * x_hat + beta

# 4 store variables for backward pass in cache var
cache = eps, var, gamma, beta, x, x_hat, mun, xc

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

elif mode == 'test':
#
===== #
# YOUR CODE HERE:
# Calculate the testing time normalized activation.
Normalize using
# the running mean and variance, and then scale and shift
appropriately.
# Store the output as 'out'.
#
===== #

# training time you use batches, testing time you use all the
data
xc = x - running_mean
var = running_var
mun = running_mean
x_hat = (x - running_mean)/(np.sqrt(running_var + eps))
out = gamma * x_hat + beta
cache = eps, var, gamma, beta, x, x_hat, mun, xc

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

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

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

```

```

    return out, cache

def batchnorm_backward(dout, cache):
    """
    Backward pass for batch normalization.

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

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

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

    # =====#
    # YOUR CODE HERE:
    # Implement the batchnorm backward pass, calculating dx, dgamma,
    and dbeta.
    # =====#
    #

    # x_hat, gamma, beta, running_mean, running_var, eps = cache
    # m = dout.shape[0]

    # dbeta = np.sum(dout, axis=0)
    # dgamma = np.sum(x_hat*dout, axis=0)

    # dx_hat = dout * gamma

    # inv_var = 1/np.sqrt(running_var + eps)

    # # print("inside dout", dout)

    # dx = (1/m) * inv_var * (m*dx_hat - np.sum(dx_hat, axis=0) -
    x_hat*np.sum(dx_hat*x_hat, axis=0))

    # try again
    eps, var, gamma, beta, x, x_hat, mun, xc = cache

```

```

m = dout.shape[0]

dbeta = np.sum(dout, axis=0)
dgamma = np.sum(dout * (x - mun) / np.sqrt(var + eps), axis=0)

dxhat = dout * gamma
dsiginv = np.sum(dxhat * xc, axis=0)
dsig = dsiginv * -1/(var + eps)
dvar = dsig / 2 * 1/np.sqrt(var + eps)

dxc = dxhat / np.sqrt(var + eps)
dxc += 2.0/m*xc*dvar

dmun = -np.sum(dxc/m, axis=0)

dx = dmum + dxc

# # print("xhat shape", dx_hat.shape)
# # print("gamma shape", gamma.shape)
# # print("x shape", x.shape)
# da = 1/(np.sqrt(running_var + eps)) * dx_hat
# dnu = -1/(np.sqrt(running_var + eps)) * np.sum(dx_hat, axis=0)
# db = (x - running_mean).dot(dx_hat.T)
# dc = -1/(running_var + eps) * db
# de = -1/2 * (running_var + eps)**(-1/2) * dc
# dvar = np.sum(de, axis=0)

# d1 = da
# d2 = 2 * (x_hat - running_mean)/m * dvar
# d3 = 1/m * dnu

# dx = d1 + d2 + d3

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

return dx, dgamma, dbeta

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

    Inputs:
    - x: Input data, of any shape
    - dropout_param: A dictionary with the following keys:
        - p: Dropout parameter. We keep each neuron output with
    probability p.
    """

```

```

        - mode: 'test' or 'train'. If the mode is train, then perform
dropout;
            if the mode is test, then just return the input.
        - seed: Seed for the random number generator. Passing seed makes
this
            function deterministic, which is needed for gradient checking
but not in
            real networks.

    Outputs:
        - out: Array of the same shape as x.
        - cache: A tuple (dropout_param, mask). In training mode, mask is
the dropout
            mask that was used to multiply the input; in test mode, mask is
None.
        """
        p, mode = dropout_param['p'], dropout_param['mode']
        assert (0 < p <= 1), "Dropout probability is not in (0,1]"
        if 'seed' in dropout_param:
            np.random.seed(dropout_param['seed'])

        mask = None
        out = None

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

            # mask
            mask = (np.random.rand(*x.shape) < p) / p
            out = x * mask

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

        elif mode == 'test':
            #
            ===== #

```

```

# YOUR CODE HERE:
#   Implement the inverted dropout forward pass during test
time.
#
=====
# for test time, just return the input
out = x

#
=====
# END YOUR CODE HERE
#
=====
# cache = (dropout_param, mask)
out = out.astype(x.dtype, copy=False)

return out, cache

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

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

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

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

```

```

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

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

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

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

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

```

Inputs:

- x: Input data, of shape (N, C) where $x[i, j]$ is the score for the j th class for the i th input.
- y: Vector of labels, of shape (N,) where $y[i]$ is the label for $x[i]$ and $0 \leq y[i] < C$

Returns a tuple of:

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

```
probs = np.exp(x - np.max(x, axis=1, keepdims=True))
probs /= np.sum(probs, axis=1, keepdims=True)
N = x.shape[0]
loss = -np.sum(np.log(np.maximum(probs[np.arange(N), y], 1e-8))) /
N
dx = probs.copy()
dx[np.arange(N), y] -= 1
dx /= N
return loss, dx
```

```

import numpy as np
import pdb

from .layers import *
from .layer_utils import *

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

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

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

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

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

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

    #
===== #

```

```

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

        self.params['W1'] = weight_scale * np.random.randn(input_dim,
hidden_dims)
        self.params['b1'] = np.zeros(hidden_dims)
        self.params['W2'] = weight_scale *
np.random.randn(hidden_dims, num_classes)
        self.params['b2'] = np.zeros(num_classes)

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

```

```

def loss(self, X, y=None):
    """
    Compute loss and gradient for a minibatch of data.

    Inputs:
    - X: Array of input data of shape (N, d_1, ..., d_k)
    - y: Array of labels, of shape (N,). y[i] gives the label for
X[i]. 

    Returns:
    If y is None, then run a test-time forward pass of the model
and return:
    - scores: Array of shape (N, C) giving classification scores,
where
        scores[i, c] is the classification score for X[i] and class
c.

    If y is not None, then run a training-time forward and
backward pass and
    return a tuple of:
    - loss: Scalar value giving the loss
    - grads: Dictionary with the same keys as self.params, mapping

```

```

parameter
    names to gradients of the loss with respect to those
parameters.
"""
scores = None

#
===== #
# YOUR CODE HERE:
# Implement the forward pass of the two-layer neural
network. Store
# the class scores as the variable 'scores'. Be sure to use
the layers
# you prior implemented.
#
===== #

# Unpack variables from the params dictionary
W1, b1 = self.params['W1'], self.params['b1']
W2, b2 = self.params['W2'], self.params['b2']

h1 = affine_forward(X, W1, b1)
relu_1 = relu_forward(h1[0])
h2 = affine_forward(relu_1[0], W2, b2)
scores = h2[0]

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

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

loss, grads = 0, {}
#
===== #
# YOUR CODE HERE:
# Implement the backward pass of the two-layer neural net.
Store
# the loss as the variable 'loss' and store the gradients in
the
# 'grads' dictionary. For the grads dictionary, grads['W1']
holds
# the gradient for W1, grads['b1'] holds the gradient for
b1, etc.
# i.e., grads[k] holds the gradient for self.params[k].

```

```

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

softmax = softmax_loss(h2[0], y)
loss = softmax[0] + 0.5 * self.reg * (np.sum(W1 ** 2) +
np.sum(W2 ** 2))

(h2_grad, grads['W2'], grads['b2']) =
affine_backward(softmax[1], (relu_1[0], W2, b2))

relu_grad = relu_backward(h2_grad, h1[0])

(h1_grad, grads['W1'], grads['b1']) =
affine_backward(relu_grad, (X, W1, b1))

grads['W2'] += self.reg*W2
grads['W1'] += self.reg*W1

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

return loss, grads

```

```

class FullyConnectedNet(object):
"""
A fully-connected neural network with an arbitrary number of
hidden layers,
ReLU nonlinearities, and a softmax loss function. This will also
implement
dropout and batch normalization as options. For a network with L
layers,
the architecture will be

{affine - [batch norm] - relu - [dropout]} x (L - 1) - affine -
softmax

where batch normalization and dropout are optional, and the {...}
block is

```

repeated $L - 1$ times.

Similar to the TwoLayerNet above, learnable parameters are stored in the self.params dictionary and will be learned using the Solver class.
"""

```
def __init__(self, hidden_dims, input_dim=3*32*32, num_classes=10,
            dropout=1, use_batchnorm=False, reg=0.0,
            weight_scale=1e-2, dtype=np.float32, seed=None):
    """
```

Initialize a new FullyConnectedNet.

Inputs:

- hidden_dims: A list of integers giving the size of each hidden layer.
- input_dim: An integer giving the size of the input.
- num_classes: An integer giving the number of classes to classify.
- dropout: Scalar between 0 and 1 giving dropout strength. If dropout=1 then the network should not use dropout at all.
- use_batchnorm: Whether or not the network should use batch normalization.
- reg: Scalar giving L2 regularization strength.
- weight_scale: Scalar giving the standard deviation for random initialization of the weights.
- dtype: A numpy datatype object; all computations will be performed using this datatype. float32 is faster but less accurate, so you should use float64 for numeric gradient checking.
- seed: If not None, then pass this random seed to the dropout layers. This will make the dropout layers deterministic so we can gradient check the model.
"""

```
        self.use_batchnorm = use_batchnorm
        self.use_dropout = dropout < 1
        self.reg = reg
        self.num_layers = 1 + len(hidden_dims)
        self.dtype = dtype
        self.params = {}
```

=====

```
# YOUR CODE HERE:
# Initialize all parameters of the network in the
```

```

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

self.params['W1'] = weight_scale * np.random.randn(input_dim,
hidden_dims[0])
self.params['b1'] = np.zeros(hidden_dims[0])

for i in range(1, self.num_layers - 1):
    w = str("W" + str(i + 1))
    b = str("b" + str(i + 1))
    self.params[w] = weight_scale *
np.random.randn(hidden_dims[i-1], hidden_dims[i])
    self.params[b] = np.zeros(hidden_dims[i])

w = str("W" + str(self.num_layers))
b = str("b" + str(self.num_layers))
self.params[w] = weight_scale *
np.random.randn(hidden_dims[-1], num_classes)
self.params[b] = np.zeros(num_classes)

# batch norm
if self.use_batchnorm:
    for i in range(self.num_layers - 1):
        self.params['gamma'+str(i+1)] = np.ones(hidden_dims[i])
        self.params['beta'+str(i+1)] = np.zeros(hidden_dims[i])

#
===== #

```

```

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

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

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

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

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

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

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

```

```

        if self.use_batchnorm:
            for bn_param in self.bn_params:
                bn_param['mode'] = mode

        scores = None

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

# try again
layer_input = X
relu_cache = []
affine_cache = []
bn_cache = []
drop_cache = []

for i in range(self.num_layers):
    layer_out, acache = affine_forward(layer_input,
self.params['W'+str(i+1)], self.params['b'+str(i+1)])
    affine_cache.append(acache)

    # batchnorm/dropout maybe
    if i != self.num_layers - 1:
        if self.use_batchnorm:
            layer_out, bcache = batchnorm_forward(layer_out,
self.params['gamma'+str(i+1)], self.params['beta'+str(i+1)],
self.bn_params[i])
            bn_cache.append(bcache)

        # relu
        layer_out, rcache = relu_forward(layer_out)
        relu_cache.append(rcache)

```

```

        # potentially dropout
        if self.use_dropout:
            layer_out, dcache = dropout_forward(layer_out,
self.dropout_param)
            drop_cache.append(dcache)

    layer_input = layer_out

    # last layer scores
    scores = layer_out

    # if self.use_batchnorm:
    #     # use batchnorm
    #     print("using batchnorm")

    #     Hs = [X]
    #     caches = []

    #     print("num layers", self.num_layers)

    #     # { affine - batchnorm - relu } * ( L-1 )
    #     for i in range(self.num_layers - 1):
    #         # print(i)
    #         W = self.params[str('W' + str(i+1))]
    #         b = self.params[str('b' + str(i+1))]
    #         gamma = self.params[str('gamma' + str(i+1))]
    #         beta = self.params[str('beta' + str(i+1))]
    #         H = Hs[-1]
    #         out, cache = affine_batchnorm_relu_forward(H, W, b,
gamma, beta, self.bn_params[i])
    #         # print("w shape", W.shape)
    #         # print("b shape", b.shape)
    #         # print("x shape", H.shape)

    #         Hs.append(out)
    #         caches.append(cache)

    #     # affine - (softmax is in the next part)

    #     # print(i+1)
    #     W = self.params[str('W' + str(i+2))]
    #     b = self.params[str('b' + str(i+2))]
    #     H = Hs[-1]
    #     # print("w shape", W.shape)
    #     # print("b shape", b.shape)
    #     # print("x shape", H.shape)
    #     out, cache = affine_forward(H, W, b)
    #     Hs.append(out)

```

```

#     caches.append(cache)

# # forward prop w/o batch norm
# else:
#     Hs = [X]
#     Zs = [X]
#     Ws = []
#     bs = []

#     W = self.params['W1']
#     b = self.params['b1']
#     aff_fwd = affine_forward(X, W, b)
#     Z = aff_fwd[0]

#     for i in range(1, self.num_layers):
#         relu_h = relu_forward(Z)
#         H = relu_h[0]
#         Hs.append(H)
#         Ws.append(W)
#         Zs.append(Z)
#         bs.append(b)

#         H = Hs[-1]
#         W = self.params[str('W' + str(i+1))]
#         b = self.params[str('b' + str(i+1))]

#         aff_fwd = affine_forward(H, W, b)
#         Z = aff_fwd[0]
#         scores = Z

#         Zs.append(Z)
#         Ws.append(W)
#         bs.append(b)

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

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

loss, grads = 0.0, {}
#

```

```

=====
# YOUR CODE HERE:
# Implement the backwards pass of the FC net and store the
gradients
# in the grads dict, so that grads[k] is the gradient of
self.params[k]
# Be sure your L2 regularization includes a 0.5 factor.
#
# BATCHNORM: Incorporate the backward pass of the batchnorm.
#
# DROPOUT: Incorporate the backward pass of dropout.
#
===== #

# print('num layers', self.num_layers)

loss, grad = softmax_loss(scores, y)
up_grad = grad

for i in reversed(np.arange(self.num_layers)):
    # print(i)
    # add regularization to loss
    loss += 0.5 * self.reg * np.sum((self.params['W'+str(i+1)] *
self.params['W'+str(i+1)]))

    # backward pass
    if i == 0:
        da, dw, db = affine_backward(up_grad, affine_cache.pop())

    else:
        da, dw, db = affine_backward(up_grad, affine_cache.pop())

        if self.use_dropout:
            da = dropout_backward(da, drop_cache.pop())

    dz = relu_backward(da, relu_cache.pop())
    if self.use_batchnorm:
        dz, dgamma, dbeta = batchnorm_backward(dz,
bn_cache.pop())
        grads['gamma'+str(i)] = dgamma
        grads['beta'+str(i)] = dbeta

    grads['W'+str(i+1)] = dw + self.reg *
self.params['W'+str(i+1)]
    grads['b'+str(i+1)] = db

    up_grad = dz

# print(grads.keys())

```

```

# print(grads.keys())

# softmax = softmax_loss(scores, y)

# # regularization
# agg_sum = 0
# for i in range(self.num_layers):
#     agg_sum += np.sum(self.params[str('W' + str(i+1))]) ** 2

# loss = softmax[0] + 0.5 * self.reg * agg_sum

# print('i got here')

# loss_grads = [softmax[1]]
# loss_grad = loss_grads[-1]

# print("num layers", self.num_layers)

# # first do the affine layer
# print(self.num_layers-1)
# dx, grads['W' + str(self.num_layers)], grads['b' +
str(self.num_layers)] = affine_backward(loss_grad,
affine_cache[self.num_layers-1])
# # print("x shape", affine_cache[self.num_layers][0].shape)
# # print("w shape", affine_cache[self.num_layers][1].shape)
# # print("b shape", affine_cache[self.num_layers-1][2].shape)

# # print("dx shape", dx.shape)
# # print("dw shape", grads['W' +
str(self.num_layers-1)].shape)
# # print("db shape", grads['b' +
str(self.num_layers-1)].shape)

# loss_grads.append(dx)

# # then the remaining layers
# for i in range(self.num_layers-1, 0, -1):
#     print(i-1)
#     loss_grad = loss_grads[-1]
#     relu_grad = relu_backward(loss_grad, relu_cache[i-1])

#     if self.use_batchnorm:
#         d_batchnorm, grads['gamma' + str(i)], grads['beta' +
str(i)] = batchnorm_backward(relu_grad, bn_cache[i-1])
#         dx, grads['W' + str(i)], grads['b' + str(i)] =
affine_backward(d_batchnorm, affine_cache[i-1])

#     else:
#         dx, grads['W' + str(i)], grads['b' + str(i)] =

```

```

affine_backward(loss_grad, affine_cache[i-1])

# # print("x shape", affine_cache[i-1][0].shape)
# # print("w shape", affine_cache[i-1][1].shape)
# # print("b shape", affine_cache[i-1][2].shape)

# # print("dx shape", dx.shape)
# # print("dw shape", grads['W' + str(i-1)].shape)
# # print("db shape", grads['b' + str(i-1)].shape)

# loss_grads.append(dx)

# # back prop
# # if self.use_batchnorm: # batch norm
# #     loss_grads = [softmax[1]]
# #     loss_grad = loss_grads[-1]
# #     cache = caches[-1]

# #     dx, grads['W' + str(self.num_layers-1)], grads['b' +
str(self.num_layers-1)] = affine_backward(loss_grad, cache)
# #     loss_grads.append(dx)

# #     for i in range(self.num_layers-1, 1, -1):
# #         print(i)
# #         loss_grad = loss_grads[-1]
# #         dx, grads[str('W' + str(i))], grads[str('b' + str(i))] =
affine_batchnorm_relu_backward(loss_grad, caches[i-1])
# #         loss_grads.append(dx)

# # else: # no batchnorm
# #     loss_grads = [softmax[1]]

# #     for i in range(self.num_layers, 1, -1):
# #         loss_grad = loss_grads[-1]
# #         (h_grad, grads[str('W' + str(i))], grads[str('b' +
str(i))]) = affine_backward(loss_grad, (Hs[i-1], Ws[i-1], bs[i-1]))
# #         relu_grad = relu_backward(h_grad, Zs[i-1])
# #         loss_grads.append(relu_grad)

# #     loss_grad = loss_grads[-1]
# #     (h_grad, grads['W1'], grads['b1']) =
affine_backward(loss_grad, (X, Ws[0], bs[0]))

# print(grads.keys())

# for i in range(self.num_layers):
#     print("i: ", i)
#     print("self params shape", self.params[str('W' +

```

```
str(i+1)).shape)

# for i in range(self.num_layers):
#   print("i: ", i)
#   print("w grad shape", grads[str('W' + str(i+1))].shape)

# # regularization
# for i in range(self.num_layers):
#   print("i: ", i)
#   print("grad shape", grads[str('W' + str(i+1))].shape)
#   print("w shape", self.params[str('W' + str(i+1))].shape)
#   grads[str('W' + str(i+1))] += self.reg *
self.params[str('W' + str(i+1))]

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

return loss, grads
```

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.

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

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

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

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

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

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


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

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

```
X_train: (49000, 3, 32, 32)
y_train: (49000,)
X_val: (1000, 3, 32, 32)
y_val: (1000,)
X_test: (1000, 3, 32, 32)
y_test: (1000,)
```

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

# Simulate the forward pass for a two-layer network
N, D1, D2, D3 = 200, 50, 60, 3
X = np.random.randn(N, D1)
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)
a = np.maximum(0, X.dot(W1)).dot(W2)

print('Before batch normalization:')
print(' means: ', a.mean(axis=0))
print(' stds: ', a.std(axis=0))

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

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

```
Before batch normalization:
means: [-0.6742318  0.80407948  3.95415552]
stds: [30.51327184 26.81786846 36.13805553]
After batch normalization (gamma=1, beta=0)
mean: [-4.32986980e-17 -4.20843915e-17 -3.38618023e-17]
std: [0.99999999 0.99999999 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 [ ]: # Check the test-time forward pass by running the training-time
# forward pass many times to warm up the running averages, and then
# checking the means and variances of activations after a test-time
# forward pass.

N, D1, D2, D3 = 200, 50, 60, 3
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)

bn_param = {'mode': 'train'}
gamma = np.ones(D3)
beta = np.zeros(D3)
for t in np.arange(50):
    X = np.random.randn(N, D1)
```

```

a = np.maximum(0, X.dot(W1)).dot(W2)
batchnorm_forward(a, gamma, beta, bn_param)
bn_param['mode'] = 'test'
X = np.random.randn(N, D1)
a = np.maximum(0, X.dot(W1)).dot(W2)
a_norm, _ = batchnorm_forward(a, gamma, beta, bn_param)

# Means should be close to zero and stds close to one, but will be
# noisier than training-time forward passes.
print('After batch normalization (test-time):')
print(' means: ', a_norm.mean(axis=0))
print(' stds: ', a_norm.std(axis=0))

```

```

After batch normalization (test-time):
means: [-0.02747245 -0.01851219 -0.07292147]
stds: [1.07141508 0.916886 0.92060113]

```

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 [ ]: # Gradient check batchnorm backward pass

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

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

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

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

# print("--")
# print("dout", dout)
# print("da", da_num)
# print("dgamma", dgamma)
# print("--")
# print("db", db_num)
# print("dbeta", dbeta)
# print("--")
# print("dx check", dx_num)
# print("dx", dx)

```

```
dx error: 3.4477536046741704e-08
dgamma error: 8.978965955969402e-12
dbeta error: 3.3104710507890927e-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 [ ]: N, D, H1, H2, C = 2, 15, 20, 30, 10
X = np.random.randn(N, D)
y = np.random.randint(C, size=(N,))

for reg in [0, 3.14]:
    print('Running check with reg = ', reg)
    model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
                             reg=reg, weight_scale=5e-2, dtype=np.float64,
                             use_batchnorm=True)

    loss, grads = model.loss(X, y)
    print('Initial loss: ', loss)

    for name in sorted(grads):
        f = lambda _: model.loss(X, y)[0]
        grad_num = eval_numerical_gradient(f, model.params[name], verbose=False)
        print('{} relative error: {}'.format(name, rel_error(grad_num, grads[name])))
    if reg == 0: print('\n')
```

```
Running check with reg = 0
Initial loss: 2.3710490024211275
W1 relative error: 0.00015101419349347454
W2 relative error: 3.672777870998256e-05
W3 relative error: 2.443181439414118e-09
b1 relative error: 0.0022203572314083426
b2 relative error: 1.1102230246251565e-08
b3 relative error: 8.523086614559495e-11
beta1 relative error: 5.446086592476153e-09
beta2 relative error: 8.75544283248942e-08
gamma1 relative error: 5.436763679722301e-09
gamma2 relative error: 5.354138682022913e-08
```

```
Running check with reg = 3.14
Initial loss: 7.192403210139613
W1 relative error: 8.459596789798366e-07
W2 relative error: 8.081855100065242e-06
W3 relative error: 1.6006226075338707e-08
b1 relative error: 2.7755575615628914e-09
b2 relative error: 2.220446049250313e-08
b3 relative error: 2.726096590945677e-10
beta1 relative error: 7.750597039229821e-09
beta2 relative error: 2.395750124612298e-08
gamma1 relative error: 2.1389851900663752e-08
gamma2 relative error: 9.288900409058029e-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 [ ]: # Try training a very deep net with batchnorm
hidden_dims = [100, 100, 100, 100, 100]

num_train = 1000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

weight_scale = 2e-2
bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=False)
model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batchnorm=True)

bn_solver = Solver(bn_model, small_data,
                  num_epochs=10, batch_size=50,
                  update_rule='adam',
                  optim_config={
                      'learning_rate': 1e-3,
                  },
                  verbose=True, print_every=200)
bn_solver.train()
```

```

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

(Iteration 1 / 200) loss: 2.295101
(Epoch 0 / 10) train acc: 0.116000; val_acc: 0.123000
(Epoch 1 / 10) train acc: 0.278000; val_acc: 0.240000
(Epoch 2 / 10) train acc: 0.408000; val_acc: 0.328000
(Epoch 3 / 10) train acc: 0.483000; val_acc: 0.308000
(Epoch 4 / 10) train acc: 0.543000; val_acc: 0.339000
(Epoch 5 / 10) train acc: 0.586000; val_acc: 0.320000
(Epoch 6 / 10) train acc: 0.644000; val_acc: 0.342000
(Epoch 7 / 10) train acc: 0.733000; val_acc: 0.335000
(Epoch 8 / 10) train acc: 0.780000; val_acc: 0.353000
(Epoch 9 / 10) train acc: 0.827000; val_acc: 0.343000
(Epoch 10 / 10) train acc: 0.829000; val_acc: 0.331000
(Iteration 1 / 200) loss: 2.302552
(Epoch 0 / 10) train acc: 0.108000; val_acc: 0.083000
(Epoch 1 / 10) train acc: 0.241000; val_acc: 0.215000
(Epoch 2 / 10) train acc: 0.278000; val_acc: 0.243000
(Epoch 3 / 10) train acc: 0.333000; val_acc: 0.282000
(Epoch 4 / 10) train acc: 0.354000; val_acc: 0.276000
(Epoch 5 / 10) train acc: 0.369000; val_acc: 0.270000
(Epoch 6 / 10) train acc: 0.368000; val_acc: 0.267000
(Epoch 7 / 10) train acc: 0.430000; val_acc: 0.307000
(Epoch 8 / 10) train acc: 0.443000; val_acc: 0.312000
(Epoch 9 / 10) train acc: 0.552000; val_acc: 0.309000
(Epoch 10 / 10) train acc: 0.583000; val_acc: 0.299000

```

```

In [ ]: fig, axes = plt.subplots(3, 1)

ax = axes[0]
ax.set_title('Training loss')
ax.set_xlabel('Iteration')

ax = axes[1]
ax.set_title('Training accuracy')
ax.set_xlabel('Epoch')

ax = axes[2]
ax.set_title('Validation accuracy')
ax.set_xlabel('Epoch')

ax = axes[0]
ax.plot(solver.loss_history, 'o', label='baseline')
ax.plot(bn_solver.loss_history, 'o', label='batchnorm')

ax = axes[1]
ax.plot(solver.train_acc_history, '-o', label='baseline')
ax.plot(bn_solver.train_acc_history, '-o', label='batchnorm')

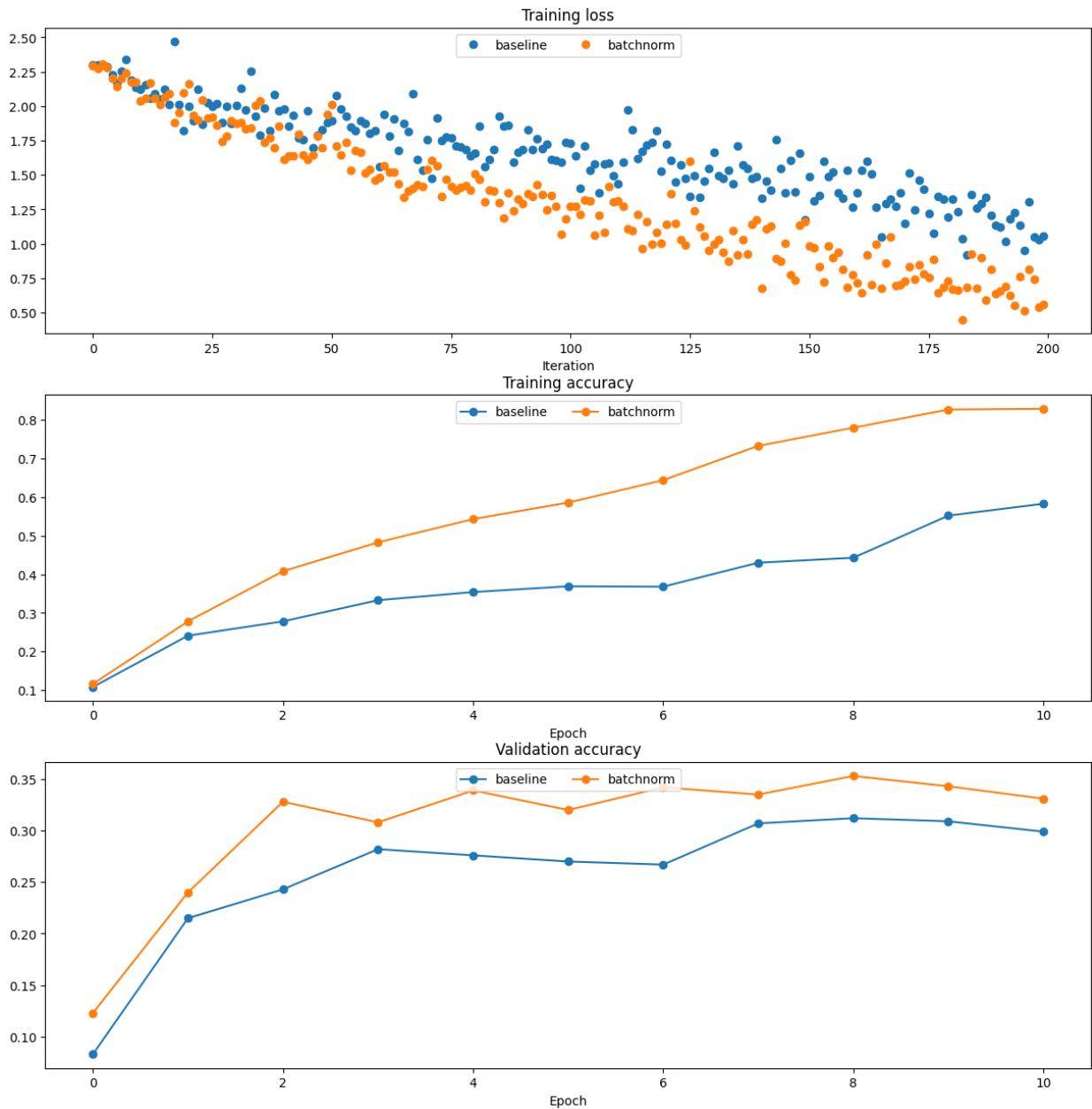
ax = axes[2]
ax.plot(solver.val_acc_history, '-o', label='baseline')
ax.plot(bn_solver.val_acc_history, '-o', label='batchnorm')

```

```

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

```



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 [ ]: # Try training a very deep net with batchnorm
hidden_dims = [50, 50, 50, 50, 50, 50, 50]

num_train = 1000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
}

```

```

'x_val': data['x_val'],
'y_val': data['y_val'],
}

bn_solvers = {}
solvers = {}
weight_scales = np.logspace(-4, 0, num=20)
for i, weight_scale in enumerate(weight_scales):
    print('Running weight scale {} / {}'.format(i + 1, len(weight_scales)))
    bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batch_norm=False)
    model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale, use_batch_norm=True)

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

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

```

Running weight scale 1 / 20
 Running weight scale 2 / 20
 Running weight scale 3 / 20
 Running weight scale 4 / 20
 Running weight scale 5 / 20
 Running weight scale 6 / 20
 Running weight scale 7 / 20
 Running weight scale 8 / 20
 Running weight scale 9 / 20
 Running weight scale 10 / 20
 Running weight scale 11 / 20
 Running weight scale 12 / 20
 Running weight scale 13 / 20
 Running weight scale 14 / 20
 Running weight scale 15 / 20
 Running weight scale 16 / 20
 Running weight scale 17 / 20
 Running weight scale 18 / 20
 Running weight scale 19 / 20
 Running weight scale 20 / 20

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

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

```
bn_best_train_accs.append(max(bn_solvers[ws].train_acc_history))

best_val_accs.append(max(solvers[ws].val_acc_history))
bn_best_val_accs.append(max(bn_solvers[ws].val_acc_history))

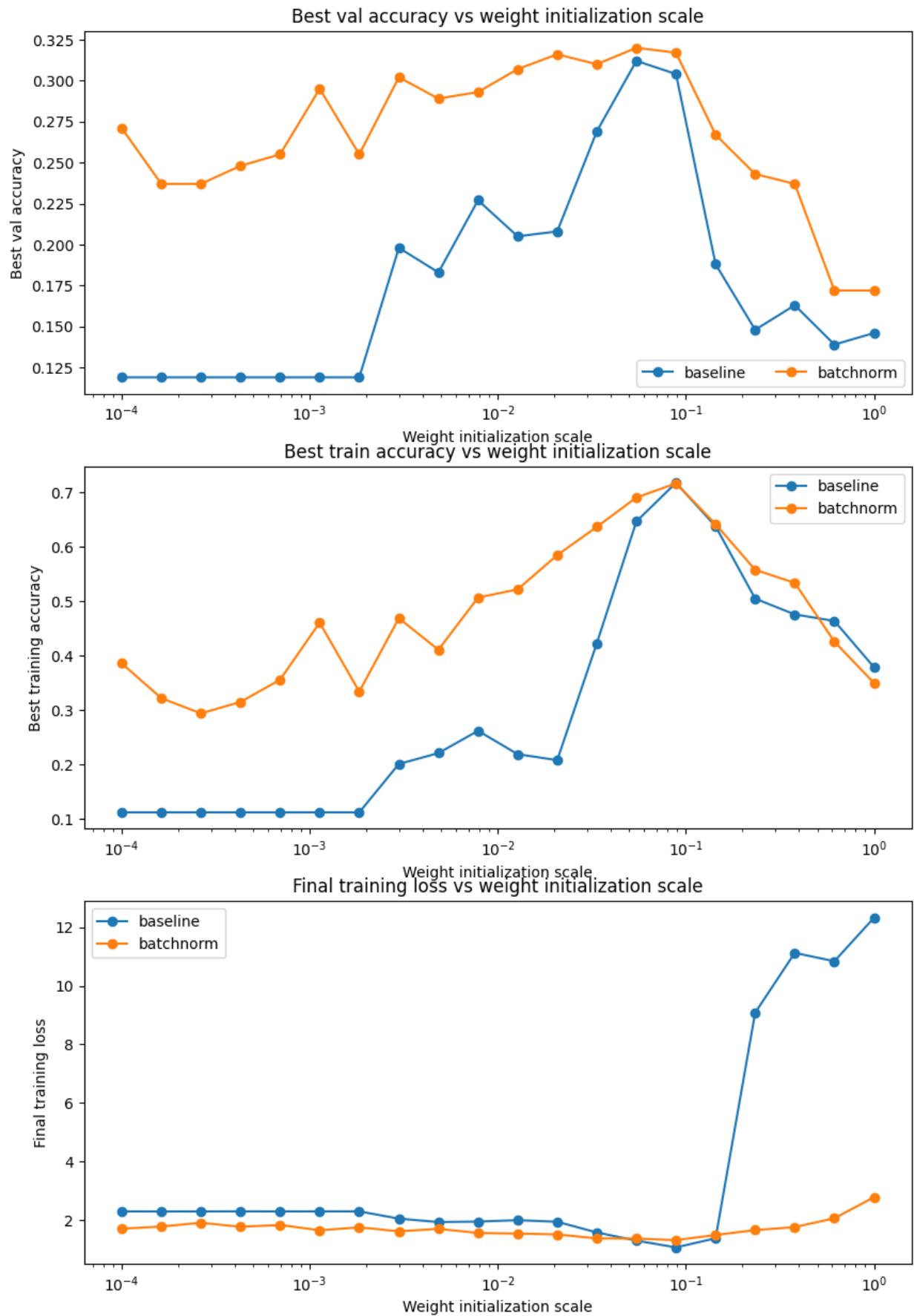
final_train_loss.append(np.mean(solvers[ws].loss_history[-100:]))
bn_final_train_loss.append(np.mean(bn_solvers[ws].loss_history[-100:]))

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

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

plt.subplot(3, 1, 3)
plt.title('Final training loss vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Final training loss')
plt.semilogx(weight_scales, final_train_loss, '-o', label='baseline')
plt.semilogx(weight_scales, bn_final_train_loss, '-o', label='batchnorm')
plt.legend()

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



Question:

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

Answer:

For both batchnorm and non-batchnorm, the best train and validation accuracy are at weight initialization scale 10e-1.

For the final training, there is a slight decrease along the graph and 10e-1 is the best for both batchnorm and non-batchnorm. For the non-batchnorm (baseline) in particular, loss skyrockets up for weight initialization scale > 10e-1.

These results make sense because when the weight scale is 10e-1, this reduces the problem of disappearing gradients. This is similar to the Xavier weight initialization of $-(1/\sqrt{n})$ and $1/\sqrt{n}$, which makes sense because this is the optimal initialization strategy. This strategy helps mitigate the exploding/vanishing gradients problem.

In []:

Dropout

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

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

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

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

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

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

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


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

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

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

x_train: (49000, 3, 32, 32)
y_train: (49000,)
x_val: (1000, 3, 32, 32)
y_val: (1000,)
x_test: (1000, 3, 32, 32)
y_test: (1000,)
```

Dropout forward pass

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

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

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

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

```
Running tests with p =  0.3
Mean of input:  10.001199068997368
Mean of train-time output:  10.0700616634419
Mean of test-time output:  10.001199068997368
Fraction of train-time output set to zero:  0.697984
Fraction of test-time output set to zero:  0.0
Running tests with p =  0.6
Mean of input:  10.001199068997368
Mean of train-time output:  10.017545453767633
Mean of test-time output:  10.001199068997368
Fraction of train-time output set to zero:  0.398944
Fraction of test-time output set to zero:  0.0
Running tests with p =  0.75
Mean of input:  10.001199068997368
Mean of train-time output:  9.999852007144597
Mean of test-time output:  10.001199068997368
Fraction of train-time output set to zero:  0.250056
Fraction of test-time output set to zero:  0.0
```

Dropout backward pass

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

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

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

```
print('dx relative error: ', rel_error(dx, dx_num))
```

```
dx relative error:  5.445610892205213e-11
```

Implement a fully connected neural network with dropout layers

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

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

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

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

In []:

```
N, D, H1, H2, C = 2, 15, 20, 30, 10
X = np.random.randn(N, D)
y = np.random.randint(C, size=(N,))

for dropout in [0.5, 0.75, 1.0]:
    print('Running check with dropout = ', dropout)
    model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
                             weight_scale=5e-2, dtype=np.float64,
                             dropout=dropout, seed=123)

    loss, grads = model.loss(X, y)
    print('Initial loss: ', loss)

    for name in sorted(grads):
        f = lambda _: model.loss(X, y)[0]
        grad_num = eval_numerical_gradient(f, model.params[name], verbose=False, h=1e-5)
        print('{} relative error: {}'.format(name, rel_error(grad_num, grads[name])))
    print('\n')
```

```
Running check with dropout =  0.5
Initial loss:  2.309771209610118
W1 relative error: 2.694274363733021e-07
W2 relative error: 7.439246147919978e-08
W3 relative error: 1.910371122296728e-08
b1 relative error: 4.112891126518e-09
b2 relative error: 5.756217724722137e-10
b3 relative error: 1.3204470857080166e-10
```

```
Running check with dropout =  0.75
Initial loss:  2.306133548427975
W1 relative error: 8.72986097970181e-08
W2 relative error: 2.9777307885797295e-07
W3 relative error: 1.8832780806174298e-08
b1 relative error: 5.379486003985169e-08
b2 relative error: 3.6529949080385546e-09
b3 relative error: 9.987242764516995e-11
```

```
Running check with dropout =  1.0
Initial loss:  2.3053332250963194
W1 relative error: 1.2744095365229032e-06
W2 relative error: 4.678743300473988e-07
W3 relative error: 4.331673892536035e-08
b1 relative error: 4.0853539035931665e-08
b2 relative error: 1.951342257912746e-09
b3 relative error: 9.387142701440351e-11
```

Dropout as a regularizer

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

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

num_train = 500
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

solvers = {}
dropout_choices = [0.6, 1.0]
for dropout in dropout_choices:
    model = FullyConnectedNet([100, 100, 100], dropout=dropout)

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

```
(Iteration 1 / 125) loss: 2.300199
(Epoch 0 / 25) train acc: 0.158000; val_acc: 0.127000
(Epoch 1 / 25) train acc: 0.130000; val_acc: 0.119000
(Epoch 2 / 25) train acc: 0.226000; val_acc: 0.180000
(Epoch 3 / 25) train acc: 0.312000; val_acc: 0.253000
(Epoch 4 / 25) train acc: 0.360000; val_acc: 0.279000
(Epoch 5 / 25) train acc: 0.376000; val_acc: 0.283000
(Epoch 6 / 25) train acc: 0.392000; val_acc: 0.262000
(Epoch 7 / 25) train acc: 0.422000; val_acc: 0.280000
(Epoch 8 / 25) train acc: 0.436000; val_acc: 0.294000
(Epoch 9 / 25) train acc: 0.446000; val_acc: 0.299000
(Epoch 10 / 25) train acc: 0.498000; val_acc: 0.316000
(Epoch 11 / 25) train acc: 0.470000; val_acc: 0.283000
(Epoch 12 / 25) train acc: 0.520000; val_acc: 0.311000
(Epoch 13 / 25) train acc: 0.558000; val_acc: 0.322000
(Epoch 14 / 25) train acc: 0.552000; val_acc: 0.320000
(Epoch 15 / 25) train acc: 0.614000; val_acc: 0.326000
(Epoch 16 / 25) train acc: 0.626000; val_acc: 0.314000
(Epoch 17 / 25) train acc: 0.628000; val_acc: 0.314000
(Epoch 18 / 25) train acc: 0.674000; val_acc: 0.322000
(Epoch 19 / 25) train acc: 0.664000; val_acc: 0.337000
(Epoch 20 / 25) train acc: 0.686000; val_acc: 0.331000
(Iteration 101 / 125) loss: 1.211559
(Epoch 21 / 25) train acc: 0.744000; val_acc: 0.340000
(Epoch 22 / 25) train acc: 0.758000; val_acc: 0.340000
(Epoch 23 / 25) train acc: 0.742000; val_acc: 0.316000
(Epoch 24 / 25) train acc: 0.760000; val_acc: 0.313000
(Epoch 25 / 25) train acc: 0.790000; val_acc: 0.316000
(Iteration 1 / 125) loss: 2.300607
(Epoch 0 / 25) train acc: 0.172000; val_acc: 0.167000
(Epoch 1 / 25) train acc: 0.202000; val_acc: 0.197000
(Epoch 2 / 25) train acc: 0.314000; val_acc: 0.233000
(Epoch 3 / 25) train acc: 0.370000; val_acc: 0.277000
(Epoch 4 / 25) train acc: 0.420000; val_acc: 0.285000
(Epoch 5 / 25) train acc: 0.454000; val_acc: 0.296000
(Epoch 6 / 25) train acc: 0.510000; val_acc: 0.304000
(Epoch 7 / 25) train acc: 0.608000; val_acc: 0.323000
(Epoch 8 / 25) train acc: 0.628000; val_acc: 0.285000
(Epoch 9 / 25) train acc: 0.694000; val_acc: 0.305000
(Epoch 10 / 25) train acc: 0.764000; val_acc: 0.314000
(Epoch 11 / 25) train acc: 0.794000; val_acc: 0.303000
(Epoch 12 / 25) train acc: 0.804000; val_acc: 0.307000
(Epoch 13 / 25) train acc: 0.840000; val_acc: 0.293000
(Epoch 14 / 25) train acc: 0.866000; val_acc: 0.310000
(Epoch 15 / 25) train acc: 0.904000; val_acc: 0.312000
(Epoch 16 / 25) train acc: 0.930000; val_acc: 0.299000
(Epoch 17 / 25) train acc: 0.944000; val_acc: 0.297000
(Epoch 18 / 25) train acc: 0.976000; val_acc: 0.255000
(Epoch 19 / 25) train acc: 0.940000; val_acc: 0.278000
(Epoch 20 / 25) train acc: 0.980000; val_acc: 0.303000
(Iteration 101 / 125) loss: 0.087768
(Epoch 21 / 25) train acc: 0.980000; val_acc: 0.291000
(Epoch 22 / 25) train acc: 0.984000; val_acc: 0.290000
(Epoch 23 / 25) train acc: 0.986000; val_acc: 0.284000
(Epoch 24 / 25) train acc: 0.992000; val_acc: 0.293000
(Epoch 25 / 25) train acc: 0.996000; val_acc: 0.288000
```

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

```
train_accs = []
val_accs = []
for dropout in dropout_choices:
```

```

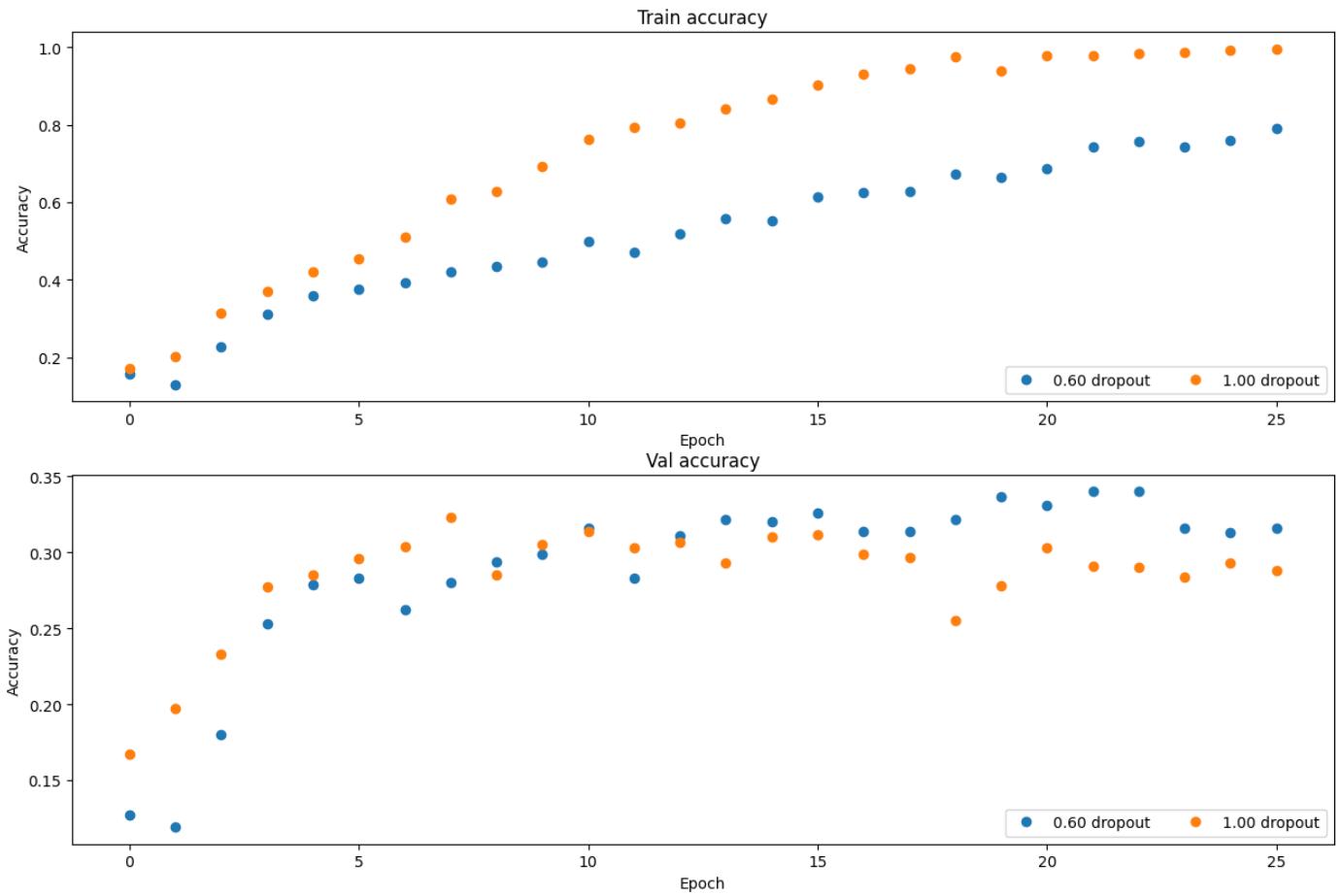
solver = solvers[dropout]
train_accs.append(solver.train_acc_history[-1])
val_accs.append(solver.val_acc_history[-1])

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

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

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

```



Question

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

Answer:

The orange when dropout = 1 represents not using dropout. So when dropout = 0.6, it is being used. So the blue represents dropout. Based on the experiment, yes dropout performs regularization because you can see both blue and orange have similar validation accuracy while the blue (dropout) doesn't have such a high rate of increase for train accuracy. This means that it is mitigating overfitting and performing regularization.

Final part of the assignment

Get over 55% validation accuracy on CIFAR-10 by using the layers you have implemented. You will be graded according to the following equation:

$\min(\text{floor}((X - 32\%)) / 23\%, 1)$ where if you get 55% or higher validation accuracy, you get full points.

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

optimizer = 'adam'
best_model = None

layer_dims = [500, 500, 500]
weight_scale = 0.01
learning_rate = 1e-3
lr_decay = 0.9

model = FullyConnectedNet(layer_dims, weight_scale=weight_scale,
                          use_batchnorm=True, dropout=0.6)

solver = Solver(model, data,
                num_epochs=10, batch_size=100,
                update_rule=optimizer,
                optim_config={
                    'learning_rate': learning_rate,
                },
                lr_decay=lr_decay,
                verbose=True, print_every=50)
solver.train()

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

```
(Iteration 1 / 4900) loss: 2.315670
(Epoch 0 / 10) train acc: 0.116000; val_acc: 0.116000
(Iteration 51 / 4900) loss: 1.942495
(Iteration 101 / 4900) loss: 1.732891
(Iteration 151 / 4900) loss: 1.610360
(Iteration 201 / 4900) loss: 1.794659
(Iteration 251 / 4900) loss: 1.614585
(Iteration 301 / 4900) loss: 1.614509
(Iteration 351 / 4900) loss: 1.440573
(Iteration 401 / 4900) loss: 1.639682
(Iteration 451 / 4900) loss: 1.613830
(Epoch 1 / 10) train acc: 0.452000; val_acc: 0.477000
(Iteration 501 / 4900) loss: 1.388291
(Iteration 551 / 4900) loss: 1.400066
(Iteration 601 / 4900) loss: 1.555558
(Iteration 651 / 4900) loss: 1.682169
(Iteration 701 / 4900) loss: 1.325679
(Iteration 751 / 4900) loss: 1.399323
(Iteration 801 / 4900) loss: 1.459228
(Iteration 851 / 4900) loss: 1.584442
(Iteration 901 / 4900) loss: 1.408745
(Iteration 951 / 4900) loss: 1.515232
(Epoch 2 / 10) train acc: 0.502000; val_acc: 0.488000
(Iteration 1001 / 4900) loss: 1.653320
(Iteration 1051 / 4900) loss: 1.454258
(Iteration 1101 / 4900) loss: 1.556873
(Iteration 1151 / 4900) loss: 1.539827
(Iteration 1201 / 4900) loss: 1.564294
(Iteration 1251 / 4900) loss: 1.226554
(Iteration 1301 / 4900) loss: 1.402670
(Iteration 1351 / 4900) loss: 1.310938
(Iteration 1401 / 4900) loss: 1.512348
(Iteration 1451 / 4900) loss: 1.405581
(Epoch 3 / 10) train acc: 0.553000; val_acc: 0.508000
(Iteration 1501 / 4900) loss: 1.431631
(Iteration 1551 / 4900) loss: 1.751771
(Iteration 1601 / 4900) loss: 1.314030
(Iteration 1651 / 4900) loss: 1.460319
(Iteration 1701 / 4900) loss: 1.384233
(Iteration 1751 / 4900) loss: 1.322205
(Iteration 1801 / 4900) loss: 1.355316
(Iteration 1851 / 4900) loss: 1.460610
(Iteration 1901 / 4900) loss: 1.445547
(Iteration 1951 / 4900) loss: 1.243986
(Epoch 4 / 10) train acc: 0.557000; val_acc: 0.537000
(Iteration 2001 / 4900) loss: 1.329682
(Iteration 2051 / 4900) loss: 1.259213
(Iteration 2101 / 4900) loss: 1.332589
(Iteration 2151 / 4900) loss: 1.332359
(Iteration 2201 / 4900) loss: 1.337638
(Iteration 2251 / 4900) loss: 1.124938
(Iteration 2301 / 4900) loss: 1.390980
(Iteration 2351 / 4900) loss: 1.345582
(Iteration 2401 / 4900) loss: 1.219983
(Epoch 5 / 10) train acc: 0.585000; val_acc: 0.546000
(Iteration 2451 / 4900) loss: 1.234222
(Iteration 2501 / 4900) loss: 1.251678
(Iteration 2551 / 4900) loss: 1.223066
(Iteration 2601 / 4900) loss: 1.093901
(Iteration 2651 / 4900) loss: 1.327221
(Iteration 2701 / 4900) loss: 1.168672
(Iteration 2751 / 4900) loss: 1.412426
```

```
(Iteration 2801 / 4900) loss: 1.127916
(Iteration 2851 / 4900) loss: 1.145958
(Iteration 2901 / 4900) loss: 1.297254
(Epoch 6 / 10) train acc: 0.597000; val_acc: 0.563000
(Iteration 2951 / 4900) loss: 1.348911
(Iteration 3001 / 4900) loss: 1.135094
(Iteration 3051 / 4900) loss: 1.212895
(Iteration 3101 / 4900) loss: 1.289152
(Iteration 3151 / 4900) loss: 1.243169
(Iteration 3201 / 4900) loss: 1.125017
(Iteration 3251 / 4900) loss: 1.283074
(Iteration 3301 / 4900) loss: 1.278779
(Iteration 3351 / 4900) loss: 1.190933
(Iteration 3401 / 4900) loss: 1.438955
(Epoch 7 / 10) train acc: 0.592000; val_acc: 0.555000
(Iteration 3451 / 4900) loss: 1.300583
(Iteration 3501 / 4900) loss: 1.444260
(Iteration 3551 / 4900) loss: 1.221397
(Iteration 3601 / 4900) loss: 1.090857
(Iteration 3651 / 4900) loss: 1.240972
(Iteration 3701 / 4900) loss: 1.373810
(Iteration 3751 / 4900) loss: 1.215900
(Iteration 3801 / 4900) loss: 1.163014
(Iteration 3851 / 4900) loss: 1.281768
(Iteration 3901 / 4900) loss: 1.124513
(Epoch 8 / 10) train acc: 0.629000; val_acc: 0.564000
(Iteration 3951 / 4900) loss: 1.246855
(Iteration 4001 / 4900) loss: 1.117939
(Iteration 4051 / 4900) loss: 1.240018
(Iteration 4101 / 4900) loss: 1.378874
(Iteration 4151 / 4900) loss: 1.297547
(Iteration 4201 / 4900) loss: 1.206859
(Iteration 4251 / 4900) loss: 1.295331
(Iteration 4301 / 4900) loss: 1.269939
(Iteration 4351 / 4900) loss: 1.419047
(Iteration 4401 / 4900) loss: 1.023463
(Epoch 9 / 10) train acc: 0.610000; val_acc: 0.564000
(Iteration 4451 / 4900) loss: 1.261891
(Iteration 4501 / 4900) loss: 1.186561
(Iteration 4551 / 4900) loss: 1.296652
(Iteration 4601 / 4900) loss: 1.202455
(Iteration 4651 / 4900) loss: 1.140189
(Iteration 4701 / 4900) loss: 1.165684
(Iteration 4751 / 4900) loss: 1.264138
(Iteration 4801 / 4900) loss: 1.149772
(Iteration 4851 / 4900) loss: 1.164538
(Epoch 10 / 10) train acc: 0.652000; val_acc: 0.564000
```

Validation accuracy reached 56.4%!!

In []: