

```
In [43]: # This mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive')

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'cs231n/assignments/assignment1/'
FOLDERNAME = 'Coursework/ENPM703/assignment2'
assert FOLDERNAME is not None, "[!] Enter the foldername."

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

# This downloads the CIFAR-10 dataset to your Drive
# if it doesn't already exist.
%cd /content/drive/My\ Drive/$FOLDERNAME/cs231n/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call `drive.mount("/content/drive", force_remount=True)`.  
/content/drive/My Drive/Coursework/ENPM703/assignment2/cs231n/datasets  
/content/drive/My Drive/Coursework/ENPM703/assignment2

## Multi-Layer Fully Connected Network

In this exercise, you will implement a fully connected network with an arbitrary number of hidden layers.

Read through the `FullyConnectedNet` class in the file `cs231n/classifiers/fc_net.py`.

Implement the network initialization, forward pass, and backward pass. Throughout this assignment, you will be implementing layers in `cs231n/layers.py`. You can re-use your implementations for `affine_forward`, `affine_backward`, `relu_forward`, `relu_backward`, and `softmax_loss` from Assignment 1. For right now, don't worry about implementing dropout or batch/layer normalization yet, as you will add those features later.

```
In [44]: # Setup cell.
import time
import numpy as np
import matplotlib.pyplot as plt
from cs231n.classifiers.fc_net import *
from cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient, eval_numerical_grad
```

```

from cs231n.solver import Solver

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

%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 [45]: # Load the (preprocessed) CIFAR-10 data.
data = get_CIFAR10_data()
for k, v in list(data.items()):
    print(f"{k}: {v.shape}")

```

```

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

```

## Initial Loss and Gradient Check

As a sanity check, run the following to check the initial loss and to gradient check the network both with and without regularization. This is a good way to see if the initial losses seem reasonable.

For gradient checking, you should expect to see errors around  $1e-7$  or less.

```

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

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
    )

    loss, grads = model.loss(X, y)

```

```

print("Initial loss: ", loss)

# Most of the errors should be on the order of e-7 or smaller.
# NOTE: It is fine however to see an error for W2 on the order of e-5
# for the check when reg = 0.0
for name in sorted(grads):
    f = lambda _: model.loss(X, y)[0]
    grad_num = eval_numerical_gradient(f, model.params[name], verbose=False)
    print(f"{name} relative error: {rel_error(grad_num, grads[name])}")

```

```

Running check with reg = 0
Initial loss: 2.3004790897684924
W1 relative error: 7.696803870986541e-08
W2 relative error: 1.7087519140575808e-05
W3 relative error: 2.9508423118300657e-07
b1 relative error: 4.660094650186831e-09
b2 relative error: 2.085654124402131e-09
b3 relative error: 6.598642296022133e-11
Running check with reg = 3.14
Initial loss: 7.052114776533016
W1 relative error: 3.904542008453064e-09
W2 relative error: 6.86942277940646e-08
W3 relative error: 2.1311298702113723e-08
b1 relative error: 1.1683196894962977e-08
b2 relative error: 1.7223751746766738e-09
b3 relative error: 1.3200479211447775e-10

```

As another sanity check, make sure your network can overfit on a small dataset of 50 images. First, we will try a three-layer network with 100 units in each hidden layer. In the following cell, tweak the **learning rate** and **weight initialization scale** to overfit and achieve 100% training accuracy within 20 epochs.

In [47]: *# TODO: Use a three-layer Net to overfit 50 training examples by  
# tweaking just the learning rate and initialization scale.*

```

num_train = 50
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 = 1e-2 # Experiment with this!
learning_rate = 1e-2 # Experiment with this!
model = FullyConnectedNet(
    [100, 100],
    weight_scale=weight_scale,
    dtype=np.float64
)
solver = Solver(
    model,
    small_data,
    print_every=10,

```

```

    num_epochs=20,
    batch_size=25,
    update_rule="sgd",
    optim_config={"learning_rate": learning_rate},
)
solver.train()

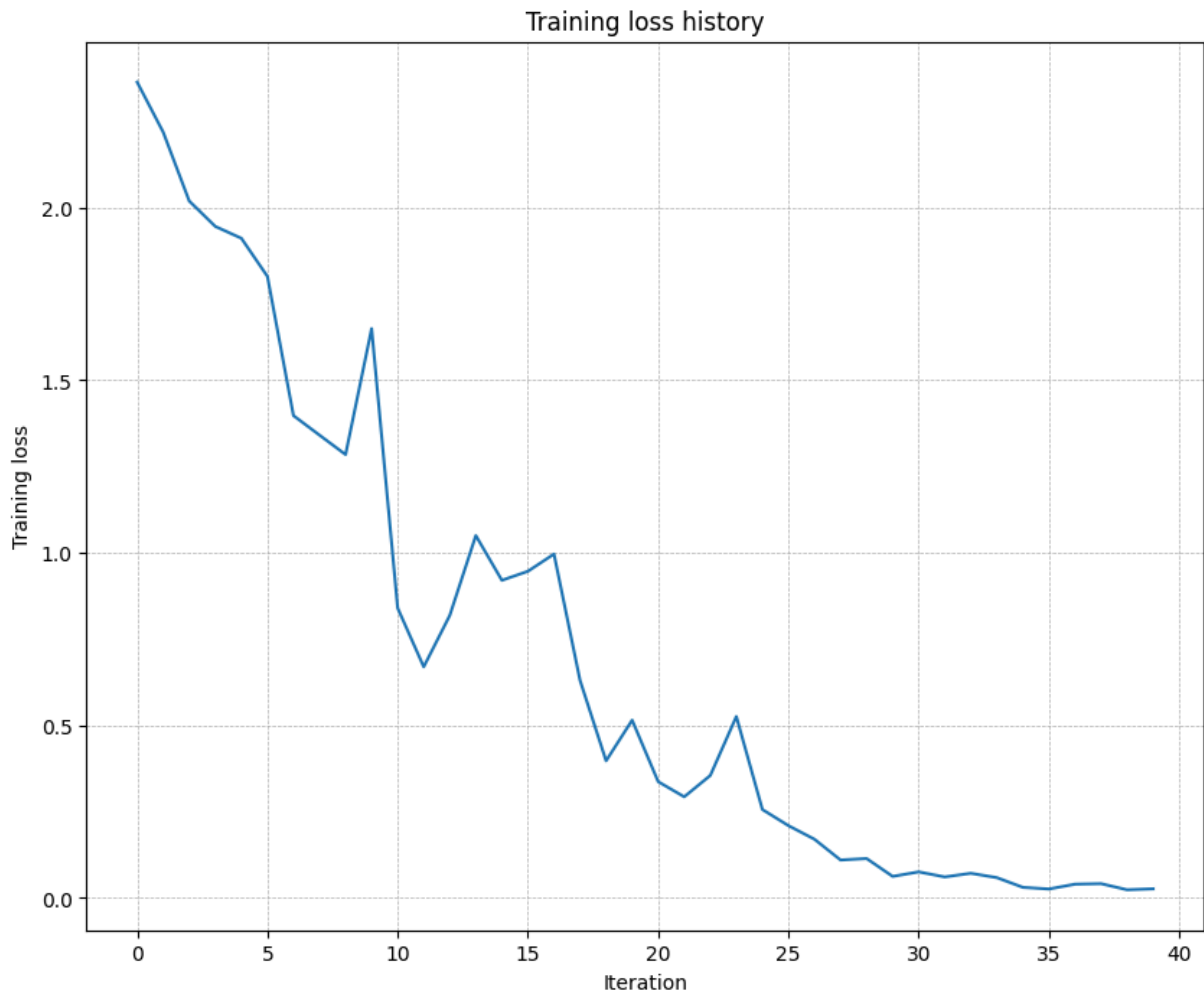
plt.plot(solver.loss_history)
plt.title("Training loss history")
plt.xlabel("Iteration")
plt.ylabel("Training loss")
plt.grid(linestyle='--', linewidth=0.5)
plt.show()

```

```

(Iteration 1 / 40) loss: 2.363364
(Epoch 0 / 20) train acc: 0.180000; val_acc: 0.108000
(Epoch 1 / 20) train acc: 0.320000; val_acc: 0.127000
(Epoch 2 / 20) train acc: 0.440000; val_acc: 0.172000
(Epoch 3 / 20) train acc: 0.500000; val_acc: 0.184000
(Epoch 4 / 20) train acc: 0.540000; val_acc: 0.181000
(Epoch 5 / 20) train acc: 0.740000; val_acc: 0.190000
(Iteration 11 / 40) loss: 0.839976
(Epoch 6 / 20) train acc: 0.740000; val_acc: 0.187000
(Epoch 7 / 20) train acc: 0.740000; val_acc: 0.183000
(Epoch 8 / 20) train acc: 0.820000; val_acc: 0.177000
(Epoch 9 / 20) train acc: 0.860000; val_acc: 0.200000
(Epoch 10 / 20) train acc: 0.920000; val_acc: 0.191000
(Iteration 21 / 40) loss: 0.337174
(Epoch 11 / 20) train acc: 0.960000; val_acc: 0.189000
(Epoch 12 / 20) train acc: 0.940000; val_acc: 0.180000
(Epoch 13 / 20) train acc: 1.000000; val_acc: 0.199000
(Epoch 14 / 20) train acc: 1.000000; val_acc: 0.199000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.195000
(Iteration 31 / 40) loss: 0.075911
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.182000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.201000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.207000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.185000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.192000

```



Now, try to use a five-layer network with 100 units on each layer to overfit on 50 training examples. Again, you will have to adjust the learning rate and weight initialization scale, but you should be able to achieve 100% training accuracy within 20 epochs.

In [55]: *# TODO: Use a five-layer Net to overfit 50 training examples by  
# tweaking just the learning rate and initialization scale.*

```
num_train = 50
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'],
}

learning_rate = 2e-3 # Experiment with this!
weight_scale = 5e-2 # Experiment with this!
model = FullyConnectedNet(
    [100, 100, 100, 100],
    weight_scale=weight_scale,
    dtype=np.float64
)
```

```

solver = Solver(
    model,
    small_data,
    print_every=10,
    num_epochs=20,
    batch_size=25,
    update_rule='sgd',
    optim_config={'learning_rate': learning_rate},
)
solver.train()

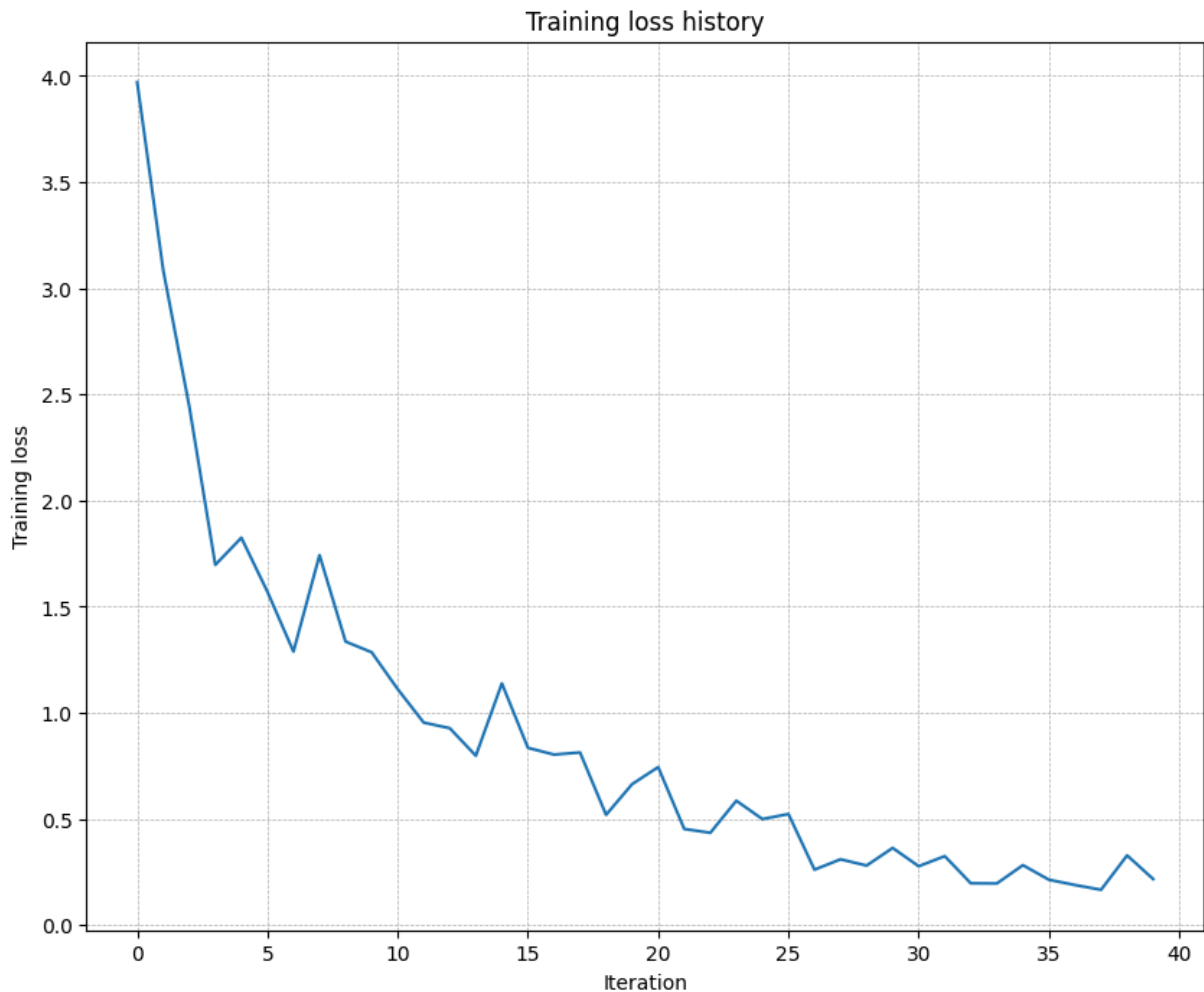
plt.plot(solver.loss_history)
plt.title('Training loss history')
plt.xlabel('Iteration')
plt.ylabel('Training loss')
plt.grid(linestyle='--', linewidth=0.5)
plt.show()

```

```

(Iteration 1 / 40) loss: 3.971181
(Epoch 0 / 20) train acc: 0.180000; val_acc: 0.139000
(Epoch 1 / 20) train acc: 0.260000; val_acc: 0.127000
(Epoch 2 / 20) train acc: 0.480000; val_acc: 0.120000
(Epoch 3 / 20) train acc: 0.560000; val_acc: 0.152000
(Epoch 4 / 20) train acc: 0.760000; val_acc: 0.158000
(Epoch 5 / 20) train acc: 0.780000; val_acc: 0.148000
(Iteration 11 / 40) loss: 1.111498
(Epoch 6 / 20) train acc: 0.800000; val_acc: 0.150000
(Epoch 7 / 20) train acc: 0.820000; val_acc: 0.144000
(Epoch 8 / 20) train acc: 0.840000; val_acc: 0.149000
(Epoch 9 / 20) train acc: 0.860000; val_acc: 0.151000
(Epoch 10 / 20) train acc: 0.880000; val_acc: 0.138000
(Iteration 21 / 40) loss: 0.743830
(Epoch 11 / 20) train acc: 0.940000; val_acc: 0.155000
(Epoch 12 / 20) train acc: 0.960000; val_acc: 0.159000
(Epoch 13 / 20) train acc: 0.960000; val_acc: 0.159000
(Epoch 14 / 20) train acc: 0.960000; val_acc: 0.147000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.159000
(Iteration 31 / 40) loss: 0.276849
(Epoch 16 / 20) train acc: 0.980000; val_acc: 0.166000
(Epoch 17 / 20) train acc: 0.980000; val_acc: 0.167000
(Epoch 18 / 20) train acc: 0.980000; val_acc: 0.168000
(Epoch 19 / 20) train acc: 0.980000; val_acc: 0.168000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.158000

```



## Inline Question 1:

Did you notice anything about the comparative difficulty of training the three-layer network vs. training the five-layer network? In particular, based on your experience, which network seemed more sensitive to the initialization scale? Why do you think that is the case?

## Answer:

When the three-layer and five-layer networks were trained side by side, it was found that the five-layer network was considerably more sensitive to the weight initialization scale. Its deeper focus increases the sensitivity since it increases the chance of seeing disappearing or bursting gradients.

## Training Notes

### Network in Three Layers:

- **Scale of Weight:  $1e-2$**

- **Training Accuracy:** By Epoch 15, **100%** was attained.
- **Loss:** Showing successful learning, it decreased from **2.31** to around **0.08**.

**Gradient Behavior:** Stabler gradients that make optimization simpler.

### Network of Five Layers:

#### Scale of Weight: 5e-2

- **Training Accuracy:** Starting at just **22%**, it reached **100%** by Epoch 19.
- **Loss:** Commenced at **3.08**, indicating initial training difficulties.

**Gradient Behavior:** Because of the possibility of disappearing or bursting gradients, increased sensitivity to weight scale resulted in more intricate tuning.

### Conclusion

Both networks attained 100% training accuracy; however, because of its depth, the five-layer network needed more cautious weight initialization. The three-layer network, on the other hand, provided a more seamless training process and a more noticeable ideal weight scale.

## Update rules

So far we have used vanilla stochastic gradient descent (SGD) as our update rule. More sophisticated update rules can make it easier to train deep networks. We will implement a few of the most commonly used update rules and compare them to vanilla SGD.

## SGD+Momentum

Stochastic gradient descent with momentum is a widely used update rule that tends to make deep networks converge faster than vanilla stochastic gradient descent. See the Momentum Update section at <http://cs231n.github.io/neural-networks-3/#sgd> for more information.

Open the file `cs231n/optim.py` and read the documentation at the top of the file to make sure you understand the API. Implement the SGD+momentum update rule in the function `sgd_momentum` and run the following to check your implementation. You should see errors less than e-8.

```
In [56]: from cs231n.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      ]])

# Should see relative errors around e-8 or less
print("next_w error: ", rel_error(next_w, expected_next_w))
print("velocity error: ", rel_error(expected_velocity, config["velocity"]))

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

```

Once you have done so, run the following to train a six-layer network with both SGD and SGD+momentum. You should see the SGD+momentum update rule converge faster.

```

In [57]: 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']:
    print('Running with ', 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': 5e-3},
        verbose=True,

```

```

    )
    solvers[update_rule] = solver
    solver.train()

fig, axes = plt.subplots(3, 1, figsize=(15, 15))

axes[0].set_title('Training loss')
axes[0].set_xlabel('Iteration')
axes[1].set_title('Training accuracy')
axes[1].set_xlabel('Epoch')
axes[2].set_title('Validation accuracy')
axes[2].set_xlabel('Epoch')

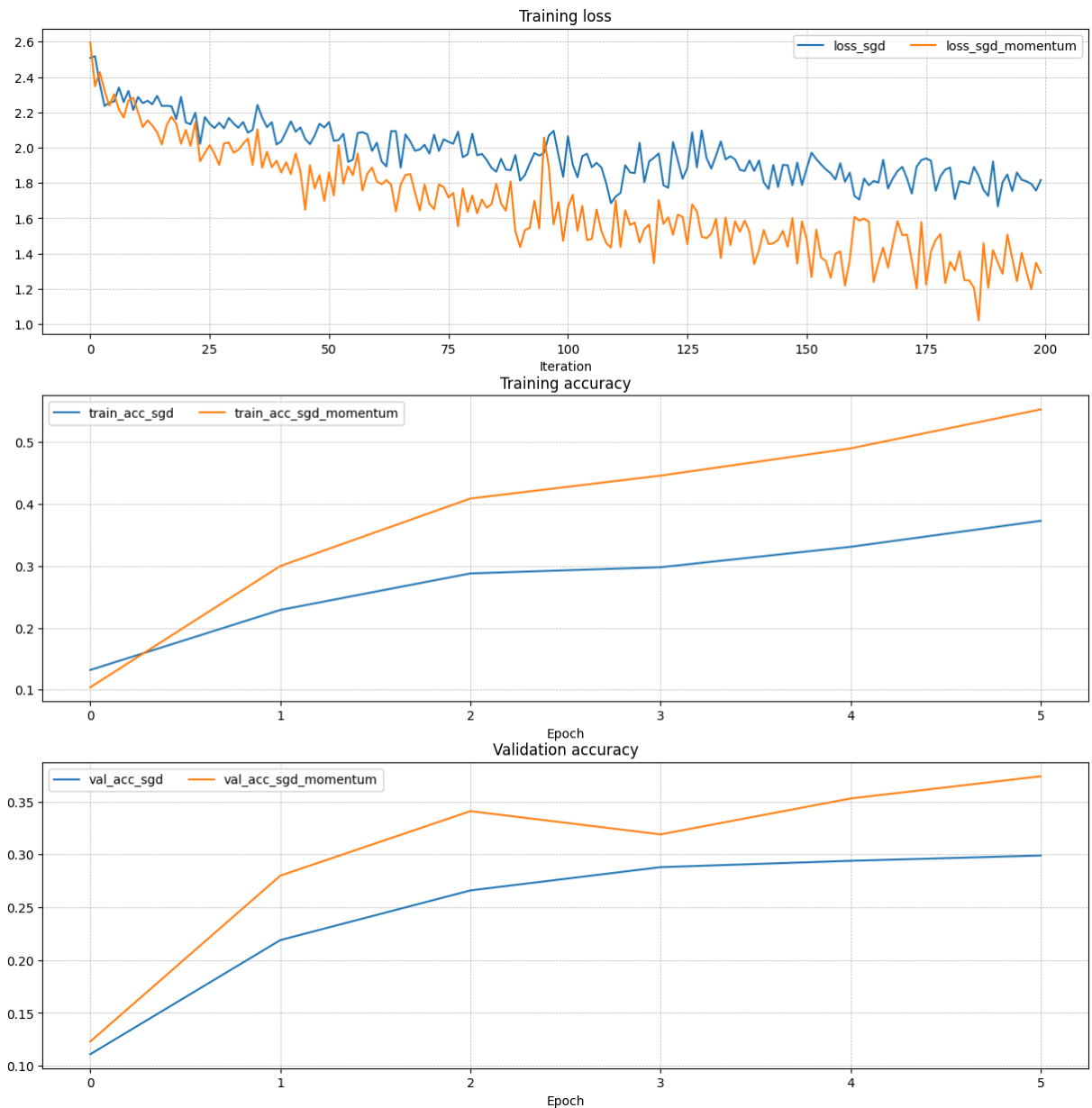
for update_rule, solver in solvers.items():
    axes[0].plot(solver.loss_history, label=f"loss_{update_rule}")
    axes[1].plot(solver.train_acc_history, label=f"train_acc_{update_rule}")
    axes[2].plot(solver.val_acc_history, label=f"val_acc_{update_rule}")

for ax in axes:
    ax.legend(loc="best", ncol=4)
    ax.grid(linestyle='--', linewidth=0.5)

plt.show()

```

```
Running with  sgd
(Iteration 1 / 200) loss: 2.508084
(Epoch 0 / 5) train acc: 0.132000; val_acc: 0.111000
(Iteration 11 / 200) loss: 2.287128
(Iteration 21 / 200) loss: 2.143093
(Iteration 31 / 200) loss: 2.137349
(Epoch 1 / 5) train acc: 0.229000; val_acc: 0.219000
(Iteration 41 / 200) loss: 2.035900
(Iteration 51 / 200) loss: 2.144458
(Iteration 61 / 200) loss: 2.027986
(Iteration 71 / 200) loss: 2.016312
(Epoch 2 / 5) train acc: 0.288000; val_acc: 0.266000
(Iteration 81 / 200) loss: 2.079021
(Iteration 91 / 200) loss: 1.813389
(Iteration 101 / 200) loss: 2.064546
(Iteration 111 / 200) loss: 1.723232
(Epoch 3 / 5) train acc: 0.298000; val_acc: 0.288000
(Iteration 121 / 200) loss: 1.787481
(Iteration 131 / 200) loss: 1.880817
(Iteration 141 / 200) loss: 1.927259
(Iteration 151 / 200) loss: 1.883699
(Epoch 4 / 5) train acc: 0.331000; val_acc: 0.294000
(Iteration 161 / 200) loss: 1.725596
(Iteration 171 / 200) loss: 1.891697
(Iteration 181 / 200) loss: 1.887813
(Iteration 191 / 200) loss: 1.667553
(Epoch 5 / 5) train acc: 0.373000; val_acc: 0.299000
Running with  sgd_momentum
(Iteration 1 / 200) loss: 2.595847
(Epoch 0 / 5) train acc: 0.104000; val_acc: 0.123000
(Iteration 11 / 200) loss: 2.202265
(Iteration 21 / 200) loss: 2.100635
(Iteration 31 / 200) loss: 1.971155
(Epoch 1 / 5) train acc: 0.300000; val_acc: 0.280000
(Iteration 41 / 200) loss: 1.858637
(Iteration 51 / 200) loss: 1.859528
(Iteration 61 / 200) loss: 1.809259
(Iteration 71 / 200) loss: 1.791472
(Epoch 2 / 5) train acc: 0.409000; val_acc: 0.341000
(Iteration 81 / 200) loss: 1.728877
(Iteration 91 / 200) loss: 1.436205
(Iteration 101 / 200) loss: 1.659388
(Iteration 111 / 200) loss: 1.699924
(Epoch 3 / 5) train acc: 0.446000; val_acc: 0.319000
(Iteration 121 / 200) loss: 1.568407
(Iteration 131 / 200) loss: 1.512991
(Iteration 141 / 200) loss: 1.420275
(Iteration 151 / 200) loss: 1.478698
(Epoch 4 / 5) train acc: 0.490000; val_acc: 0.353000
(Iteration 161 / 200) loss: 1.607737
(Iteration 171 / 200) loss: 1.503822
(Iteration 181 / 200) loss: 1.353460
(Iteration 191 / 200) loss: 1.349060
(Epoch 5 / 5) train acc: 0.553000; val_acc: 0.374000
```



## RMSProp and Adam

RMSProp [1] and Adam [2] are update rules that set per-parameter learning rates by using a running average of the second moments of gradients.

In the file `cs231n/optim.py`, implement the RMSProp update rule in the `rmsprop` function and implement the Adam update rule in the `adam` function, and check your implementations using the tests below.

**NOTE:** Please implement the *complete* Adam update rule (with the bias correction mechanism), not the first simplified version mentioned in the course notes.

[1] Tijmen Tieleman and Geoffrey Hinton. "Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude." COURSE: Neural

Networks for Machine Learning 4 (2012).

[2] Diederik Kingma and Jimmy Ba, "Adam: A Method for Stochastic Optimization", ICLR 2015.

```
In [58]: # Test RMSProp implementation
from cs231n.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)
cache = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)

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

expected_next_w = np.asarray([
    [-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 ]])

# You should see relative errors around e-7 or less
print('next_w error: ', rel_error(expected_next_w, next_w))
print('cache error: ', rel_error(expected_cache, config['cache']))

next_w error: 9.524687511038133e-08
cache error: 2.6477955807156126e-09
```

```
In [59]: # Test Adam implementation
from cs231n.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)
m = np.linspace(0.6, 0.9, num=N*D).reshape(N, D)
v = np.linspace(0.7, 0.5, num=N*D).reshape(N, D)

config = {'learning_rate': 1e-2, 'm': m, 'v': v, '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_v = 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],
    [ 0.54036881, 0.52995808, 0.51954735, 0.50913662, 0.49872589]])
```

```

[ 0.54159906, 0.53110598, 0.52061845, 0.51013645, 0.49966,   ]])
expected_m = 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       ]])

# You should see relative errors around e-7 or less
print('next_w error: ', rel_error(expected_next_w, next_w))
print('v error: ', rel_error(expected_v, config['v']))
print('m error: ', rel_error(expected_m, config['m']))

```

next\_w error: 1.1395691798535431e-07

v error: 4.208314038113071e-09

m error: 4.214963193114416e-09

Once you have debugged your RMSProp and Adam implementations, run the following to train a pair of deep networks using these new update rules:

```

In [60]: learning_rates = {'rmsprop': 1e-4, 'adam': 1e-3}
for update_rule in ['adam', 'rmsprop']:
    print('Running with ', 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=True
    )
    solvers[update_rule] = solver
    solver.train()
    print()

fig, axes = plt.subplots(3, 1, figsize=(15, 15))

axes[0].set_title('Training loss')
axes[0].set_xlabel('Iteration')
axes[1].set_title('Training accuracy')
axes[1].set_xlabel('Epoch')
axes[2].set_title('Validation accuracy')
axes[2].set_xlabel('Epoch')

for update_rule, solver in solvers.items():
    axes[0].plot(solver.loss_history, label=f"{update_rule}")
    axes[1].plot(solver.train_acc_history, label=f"{update_rule}")
    axes[2].plot(solver.val_acc_history, label=f"{update_rule}")

for ax in axes:
    ax.legend(loc='best', ncol=4)
    ax.grid(linestyle='--', linewidth=0.5)

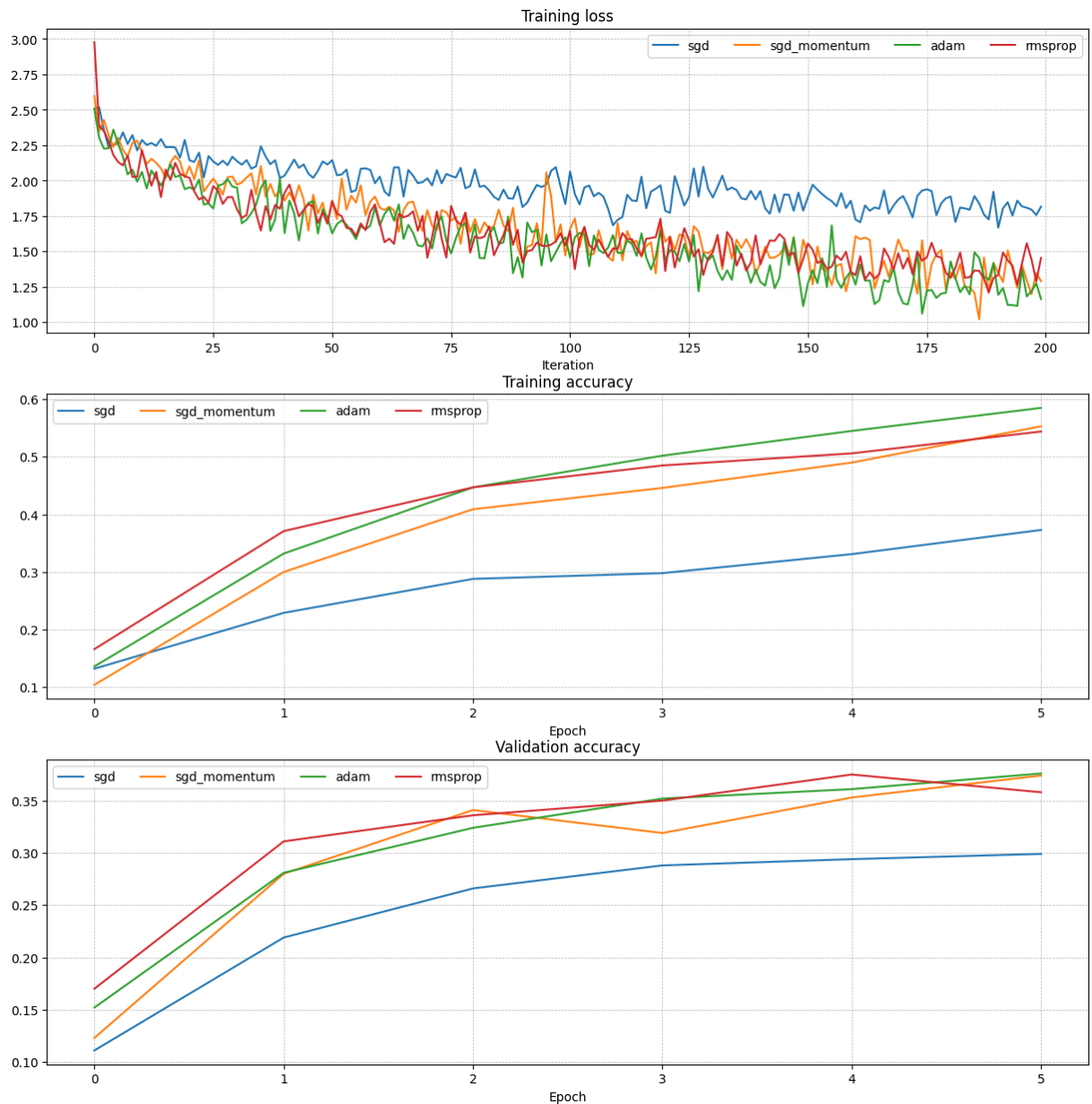
```

```
plt.show()
```

Running with adam  
(Iteration 1 / 200) loss: 2.503933  
(Epoch 0 / 5) train acc: 0.136000; val\_acc: 0.152000  
(Iteration 11 / 200) loss: 2.061754  
(Iteration 21 / 200) loss: 1.952930  
(Iteration 31 / 200) loss: 1.949460  
(Epoch 1 / 5) train acc: 0.332000; val\_acc: 0.281000  
(Iteration 41 / 200) loss: 1.629484  
(Iteration 51 / 200) loss: 1.627666  
(Iteration 61 / 200) loss: 1.683181  
(Iteration 71 / 200) loss: 1.593157  
(Epoch 2 / 5) train acc: 0.447000; val\_acc: 0.324000  
(Iteration 81 / 200) loss: 1.691935  
(Iteration 91 / 200) loss: 1.316360  
(Iteration 101 / 200) loss: 1.457656  
(Iteration 111 / 200) loss: 1.490421  
(Epoch 3 / 5) train acc: 0.502000; val\_acc: 0.352000  
(Iteration 121 / 200) loss: 1.437246  
(Iteration 131 / 200) loss: 1.472117  
(Iteration 141 / 200) loss: 1.284863  
(Iteration 151 / 200) loss: 1.276522  
(Epoch 4 / 5) train acc: 0.545000; val\_acc: 0.361000  
(Iteration 161 / 200) loss: 1.265843  
(Iteration 171 / 200) loss: 1.134311  
(Iteration 181 / 200) loss: 1.429556  
(Iteration 191 / 200) loss: 1.193171  
(Epoch 5 / 5) train acc: 0.585000; val\_acc: 0.376000

Running with rmsprop  
(Iteration 1 / 200) loss: 2.976445  
(Epoch 0 / 5) train acc: 0.166000; val\_acc: 0.170000  
(Iteration 11 / 200) loss: 2.215317  
(Iteration 21 / 200) loss: 2.019541  
(Iteration 31 / 200) loss: 1.796836  
(Epoch 1 / 5) train acc: 0.371000; val\_acc: 0.311000  
(Iteration 41 / 200) loss: 1.914885  
(Iteration 51 / 200) loss: 1.854153  
(Iteration 61 / 200) loss: 1.698718  
(Iteration 71 / 200) loss: 1.455506  
(Epoch 2 / 5) train acc: 0.447000; val\_acc: 0.336000  
(Iteration 81 / 200) loss: 1.607075  
(Iteration 91 / 200) loss: 1.421859  
(Iteration 101 / 200) loss: 1.644586  
(Iteration 111 / 200) loss: 1.517447  
(Epoch 3 / 5) train acc: 0.485000; val\_acc: 0.350000  
(Iteration 121 / 200) loss: 1.361570  
(Iteration 131 / 200) loss: 1.493687  
(Iteration 141 / 200) loss: 1.408904  
(Iteration 151 / 200) loss: 1.555976  
(Epoch 4 / 5) train acc: 0.506000; val\_acc: 0.375000  
(Iteration 161 / 200) loss: 1.332950  
(Iteration 171 / 200) loss: 1.378604  
(Iteration 181 / 200) loss: 1.314109  
(Iteration 191 / 200) loss: 1.334453  
(Epoch 5 / 5) train acc: 0.544000; val\_acc: 0.358000





## Inline Question 2:

AdaGrad, like Adam, is a per-parameter optimization method that uses the following update rule:

```
cache += dw**2
w += - learning_rate * dw / (np.sqrt(cache) + eps)
```

John notices that when he was training a network with AdaGrad that the updates became very small, and that his network was learning slowly. Using your knowledge of the AdaGrad update rule, why do you think the updates would become very small? Would Adam have the same issue?

**Answer:**

Smaller updates result from a decreasing step size over time due to the monotonic growth in the cache's squared gradients' cumulative value. Slower updates at the global minimum might be beneficial if the optimization function is convex. On the other hand, this may result in the algorithm becoming trapped in local minima for non-convex functions.

The Adam algorithm uses bias correction to lessen the effects of the growing time, therefore avoiding this problem.

$b_1 = m/(1-\beta_1^t)$  and  $b_2 = v/(1-\beta_2^t)$  are the update rules for Adam. The numerator of the updates will also rise when  $(t)$  increases, as both  $(b_1)$  and  $(b_2)$  increase as well. This keeps the updates from being too tiny.

## Train a Good Model!

Train the best fully connected model that you can on CIFAR-10, storing your best model in the `best_model` variable. We require you to get at least 50% accuracy on the validation set using a fully connected network.

If you are careful it should be possible to get accuracies above 55%, but we don't require it for this part and won't assign extra credit for doing so. Later in the assignment we will ask you to train the best convolutional network that you can on CIFAR-10, and we would prefer that you spend your effort working on convolutional networks rather than fully connected networks.

**Note:** You might find it useful to complete the `BatchNormalization.ipynb` and `Dropout.ipynb` notebooks before completing this part, since those techniques can help you train powerful models.

```
In [61]: best_model = None

#####
# TODO: Train the best FullyConnectedNet that you can on CIFAR-10. You might
# find batch/layer normalization and dropout useful. Store your best model i
# the best_model variable.
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

# pass
# Hyperparameters
learning_rates = 1.5e-3
weight_scales = 1e-2
reg = 1e-4

# Initialize model
model = FullyConnectedNet(
    [400, 300, 200, 100, 50],
```

```

    reg=reg,
    weight_scale=weight_scales,
    dtype=np.float64
)

# Set up solver
solver = Solver(
    model,
    data,
    num_epochs=10,
    batch_size=100,
    update_rule="adam",
    optim_config={"learning_rate": learning_rates},
    verbose=True,
    print_every=1000
)

# Train model
solver.train()

# Save best model
best_model = model

# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
#####
#                                     END OF YOUR CODE
#####

```

```

(Iteration 1 / 4900) loss: 2.309740
(Epoch 0 / 10) train acc: 0.103000; val_acc: 0.098000
(Epoch 1 / 10) train acc: 0.374000; val_acc: 0.355000
(Epoch 2 / 10) train acc: 0.418000; val_acc: 0.419000
(Iteration 1001 / 4900) loss: 1.801674
(Epoch 3 / 10) train acc: 0.484000; val_acc: 0.442000
(Epoch 4 / 10) train acc: 0.479000; val_acc: 0.443000
(Iteration 2001 / 4900) loss: 1.545944
(Epoch 5 / 10) train acc: 0.495000; val_acc: 0.472000
(Epoch 6 / 10) train acc: 0.532000; val_acc: 0.494000
(Iteration 3001 / 4900) loss: 1.353906
(Epoch 7 / 10) train acc: 0.547000; val_acc: 0.510000
(Epoch 8 / 10) train acc: 0.512000; val_acc: 0.504000
(Iteration 4001 / 4900) loss: 1.293499
(Epoch 9 / 10) train acc: 0.538000; val_acc: 0.516000
(Epoch 10 / 10) train acc: 0.544000; val_acc: 0.500000

```

## Test Your Model!

Run your best model on the validation and test sets. You should achieve at least 50% accuracy on the validation set.

```

In [62]: y_test_pred = np.argmax(best_model.loss(data['X_test']), axis=1)
         y_val_pred = np.argmax(best_model.loss(data['X_val']), axis=1)

```

```
print('Validation set accuracy: ', (y_val_pred == data['y_val']).mean())  
print('Test set accuracy: ', (y_test_pred == data['y_test']).mean())
```

Validation set accuracy: 0.516

Test set accuracy: 0.525

This notebook was converted to PDF with [convert.ploomber.io](https://convert.ploomber.io)

```

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

# TODO: Enter the foldername in your Drive where you have saved the
unzipped
# assignment folder, e.g. 'cs231n/assignments/assignment1/'
FOLDERNAME = 'Coursework/ENPM703/assignment2'
assert FOLDERNAME is not None, "[!] Enter the foldername."

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

# This downloads the CIFAR-10 dataset to your Drive
# if it doesn't already exist.
%cd /content/drive/My\ Drive/$FOLDERNAME/cs231n/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME

Drive already mounted at /content/drive; to attempt to forcibly
remount, call drive.mount("/content/drive", force_remount=True).
/content/drive/My Drive/Coursework/ENPM703/assignment2/cs231n/datasets
/content/drive/My Drive/Coursework/ENPM703/assignment2

```

## Batch Normalization

One way to make deep networks easier to train is to use more sophisticated optimization procedures such as SGD+momentum, RMSProp, or Adam. Another strategy is to change the architecture of the network to make it easier to train. One idea along these lines is batch normalization, proposed by [1] in 2015.

To understand the goal of batch normalization, it is important to first recognize that machine learning methods tend to perform better with input data consisting of uncorrelated features with zero mean and unit variance. When training a neural network, we can preprocess the data before feeding it to the network to explicitly decorrelate its features. This will ensure that the first layer of the network sees data that follows a nice distribution. However, even if we preprocess the input data, the activations at deeper layers of the network will likely no longer be decorrelated and will no longer have zero mean or unit variance, since they are output from earlier layers in the network. Even worse, during the training process the distribution of features at each layer of the network will shift as the weights of each layer are updated.

The authors of [1] hypothesize that the shifting distribution of features inside deep neural networks may make training deep networks more difficult. To overcome this problem, they propose to insert into the network layers that normalize batches. At training time, such a layer uses a minibatch of data to estimate the mean and standard deviation of each feature. These

estimated means and standard deviations are then used to center and normalize the features of the minibatch. A running average of these means and standard deviations is kept during training, and at test time these running averages are used to center and normalize features.

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

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

```
# Setup cell.
import time
import numpy as np
import matplotlib.pyplot as plt
from cs231n.classifiers.fc_net import *
from cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient,
eval_numerical_gradient_array
from cs231n.solver import Solver

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

%load_ext autoreload
%autoreload 2

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

def print_mean_std(x,axis=0):
    print(f"  means: {x.mean(axis=axis)}")
    print(f"  stds: {x.std(axis=axis)}\n")

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

# Load the (preprocessed) CIFAR-10 data.
data = get_CIFAR10_data()
for k, v in list(data.items()):
    print(f"{k}: {v.shape}")

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

```
y_val: (1000,)
X_test: (1000, 3, 32, 32)
y_test: (1000,)
```

## Batch Normalization: Forward Pass

In the file `cs231n/layers.py`, implement the batch normalization forward pass in the function `batchnorm_forward`. Once you have done so, run the following to test your implementation.

Referencing the paper linked to above in [1] may be helpful!

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

# Simulate the forward pass for a two-layer network.
np.random.seed(231)
N, D1, D2, D3 = 200, 50, 60, 3
X = np.random.randn(N, D1)
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)
a = np.maximum(0, X.dot(W1)).dot(W2)

print('Before batch normalization:')
print_mean_std(a,axis=0)

gamma = np.ones((D3,))
beta = np.zeros((D3,))

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

gamma = np.asarray([1.0, 2.0, 3.0])
beta = np.asarray([11.0, 12.0, 13.0])

# Now means should be close to beta and stds close to gamma.
print('After batch normalization (gamma=', gamma, ', beta=', beta,
')')
a_norm, _ = batchnorm_forward(a, gamma, beta, {'mode': 'train'})
print_mean_std(a_norm,axis=0)

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

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

```
means: [ 3.55271368e-17  1.71529457e-16 -2.76167977e-17]
stds:  [0.99999999  1.          1.          ]
```

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

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

```
# Check the test-time forward pass by running the training-time
# forward pass many times to warm up the running averages, and then
# checking the means and variances of activations after a test-time
# forward pass.
```

```
np.random.seed(231)
N, D1, D2, D3 = 200, 50, 60, 3
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)
```

```
bn_param = {'mode': 'train'}
gamma = np.ones(D3)
beta = np.zeros(D3)
```

```
for t in range(50):
    X = np.random.randn(N, D1)
    a = np.maximum(0, X.dot(W1)).dot(W2)
    batchnorm_forward(a, gamma, beta, bn_param)
```

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

```
# Means should be close to zero and stds close to one, but will be
# noisier than training-time forward passes.
```

```
print('After batch normalization (test-time):')
print_mean_std(a_norm,axis=0)
```

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

## Batch Normalization: Backward Pass

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

To derive the backward pass you should write out the computation graph for batch normalization and backprop through each of the intermediate nodes. Some intermediates may



have multiple outgoing branches; make sure to sum gradients across these branches in the backward pass.

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

```
# Gradient check batchnorm backward pass.
np.random.seed(231)
N, D = 4, 5
x = 5 * np.random.randn(N, D) + 12
gamma = np.random.randn(D)
beta = np.random.randn(D)
dout = np.random.randn(N, D)

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

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

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

# You should expect to see relative errors between 1e-13 and 1e-8.
print('dx error: ', rel_error(dx_num, dx))
print('dgamma error: ', rel_error(da_num, dgamma))
print('dbeta error: ', rel_error(db_num, dbeta))

dx error:  1.7029241291468676e-09
dgamma error:  7.420414216247087e-13
dbeta error:  2.8795057655839487e-12
```

## Batch Normalization: Alternative Backward Pass

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

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

In the forward pass, given a set of inputs  $X = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_N \end{bmatrix}$ ,

we first calculate the mean  $\mu$  and variance  $v$ . With  $\mu$  and  $v$  calculated, we can calculate the standard deviation  $\sigma$  and normalized data  $Y$ . The equations and graph illustration below describe the computation ( $y_i$  is the  $i$ -th element of the vector  $Y$ ).

$$\begin{aligned} & \mu = \frac{1}{N} \sum_{k=1}^N x_k & v = \frac{1}{N} \sum_{k=1}^N (x_k - \mu)^2 \\ & \sigma = \sqrt{v + \epsilon} & y_i = \frac{x_i - \mu}{\sigma} \end{aligned}$$

The meat of our problem during backpropagation is to compute  $\frac{\partial L}{\partial X}$ , given the upstream gradient we receive,  $\frac{\partial L}{\partial Y}$ . To do this, recall the chain rule in calculus gives us  $\frac{\partial L}{\partial X} = \frac{\partial L}{\partial Y} \cdot \frac{\partial Y}{\partial X}$ .

The unknown/hard part is  $\frac{\partial Y}{\partial X}$ . We can find this by first deriving step-by-step our local gradients at  $\frac{\partial v}{\partial X}$ ,  $\frac{\partial \mu}{\partial X}$ ,  $\frac{\partial \sigma}{\partial v}$ ,  $\frac{\partial Y}{\partial \sigma}$ , and  $\frac{\partial Y}{\partial \mu}$ , and then use the chain rule to compose these gradients (which appear in the form of vectors!) appropriately to compute  $\frac{\partial Y}{\partial X}$ .

If it's challenging to directly reason about the gradients over  $X$  and  $Y$  which require matrix multiplication, try reasoning about the gradients in terms of individual elements  $x_i$  and  $y_i$  first: in that case, you will need to come up with the derivations for  $\frac{\partial L}{\partial x_i}$ , by relying on the Chain Rule to first calculate the intermediate  $\frac{\partial \mu}{\partial x_i}$ ,  $\frac{\partial v}{\partial x_i}$ ,  $\frac{\partial \sigma}{\partial x_i}$ , then assemble these pieces to calculate  $\frac{\partial y_i}{\partial x_i}$ .

You should make sure each of the intermediary gradient derivations are all as simplified as possible, for ease of implementation.

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

```
np.random.seed(231)
N, D = 100, 500
x = 5 * np.random.randn(N, D) + 12
gamma = np.random.randn(D)
beta = np.random.randn(D)
dout = np.random.randn(N, D)
```

```

bn_param = {'mode': 'train'}
out, cache = batchnorm_forward(x, gamma, beta, bn_param)

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

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

dx difference:  9.20004371222927e-13
dgamma difference:  0.0
dbeta difference:  0.0
speedup: 1.10x

```

## Fully Connected Networks with Batch Normalization

Now that you have a working implementation for batch normalization, go back to your `FullyConnectedNet` in the file `cs231n/classifiers/fc_net.py`. Modify your implementation to add batch normalization.

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

**Hint:** You might find it useful to define an additional helper layer similar to those in the file `cs231n/layer_utils.py`.

```

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

# You should expect losses between 1e-4~1e-10 for W,
# losses between 1e-08~1e-10 for b,
# and losses between 1e-08~1e-09 for beta and gammas.
for reg in [0, 3.14]:
    print('Running check with reg = ', reg)
    model = FullyConnectedNet([H1, H2], input_dim=D, num_classes=C,
                              reg=reg, weight_scale=5e-2,
                              dtype=np.float64,

```

```

                                normalization='batchnorm')

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

for name in sorted(grads):
    f = lambda _: model.loss(X, y)[0]
    grad_num = eval_numerical_gradient(f, model.params[name],
    verbose=False, h=1e-5)
    print('%s relative error: %.2e' % (name, rel_error(grad_num,
    grads[name])))
    if reg == 0: print()

```

```

Running check with reg = 0
Initial loss: 2.2611955101340957
W1 relative error: 1.10e-04
W2 relative error: 3.17e-06
W3 relative error: 5.10e-10
b1 relative error: 2.22e-08
b2 relative error: 5.55e-09
b3 relative error: 1.01e-10
beta1 relative error: 7.33e-09
beta2 relative error: 1.89e-09
gamma1 relative error: 7.98e-09
gamma2 relative error: 2.27e-09

```

```

Running check with reg = 3.14
Initial loss: 6.996533220108303
W1 relative error: 1.98e-06
W2 relative error: 2.28e-06
W3 relative error: 1.11e-08
b1 relative error: 5.55e-09
b2 relative error: 2.22e-08
b3 relative error: 1.42e-10
beta1 relative error: 6.65e-09
beta2 relative error: 4.23e-09
gamma1 relative error: 5.94e-09
gamma2 relative error: 9.60e-09

```

## Batch Normalization for Deep Networks

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

```

np.random.seed(231)

# Try training a very deep net with batchnorm.
hidden_dims = [100, 100, 100, 100, 100]

```

```

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

weight_scale = 2e-2
bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale,
    normalization='batchnorm')
model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale,
    normalization=None)

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

print('\nSolver without batch norm:')
solver = Solver(model, small_data,
    num_epochs=10, batch_size=50,
    update_rule='adam',
    optim_config={
        'learning_rate': 1e-3,
    },
    verbose=True, print_every=20)
solver.train()

```

Solver with batch norm:

```

(Iteration 1 / 200) loss: 2.340974
(Epoch 0 / 10) train acc: 0.107000; val_acc: 0.115000
(Epoch 1 / 10) train acc: 0.314000; val_acc: 0.266000
(Iteration 21 / 200) loss: 2.039365
(Epoch 2 / 10) train acc: 0.389000; val_acc: 0.276000
(Iteration 41 / 200) loss: 2.036703
(Epoch 3 / 10) train acc: 0.500000; val_acc: 0.321000
(Iteration 61 / 200) loss: 1.775613
(Epoch 4 / 10) train acc: 0.532000; val_acc: 0.312000
(Iteration 81 / 200) loss: 1.274868
(Epoch 5 / 10) train acc: 0.581000; val_acc: 0.307000
(Iteration 101 / 200) loss: 1.261756
(Epoch 6 / 10) train acc: 0.662000; val_acc: 0.338000
(Iteration 121 / 200) loss: 1.049509

```

```
(Epoch 7 / 10) train acc: 0.677000; val_acc: 0.321000
(Iteration 141 / 200) loss: 1.138841
(Epoch 8 / 10) train acc: 0.702000; val_acc: 0.298000
(Iteration 161 / 200) loss: 0.802456
(Epoch 9 / 10) train acc: 0.790000; val_acc: 0.354000
(Iteration 181 / 200) loss: 1.008300
(Epoch 10 / 10) train acc: 0.773000; val_acc: 0.311000
```

Solver without batch norm:

```
(Iteration 1 / 200) loss: 2.302332
(Epoch 0 / 10) train acc: 0.129000; val_acc: 0.131000
(Epoch 1 / 10) train acc: 0.283000; val_acc: 0.250000
(Iteration 21 / 200) loss: 2.041970
(Epoch 2 / 10) train acc: 0.316000; val_acc: 0.277000
(Iteration 41 / 200) loss: 1.900473
(Epoch 3 / 10) train acc: 0.373000; val_acc: 0.282000
(Iteration 61 / 200) loss: 1.713156
(Epoch 4 / 10) train acc: 0.390000; val_acc: 0.310000
(Iteration 81 / 200) loss: 1.662209
(Epoch 5 / 10) train acc: 0.434000; val_acc: 0.300000
(Iteration 101 / 200) loss: 1.696058
(Epoch 6 / 10) train acc: 0.535000; val_acc: 0.345000
(Iteration 121 / 200) loss: 1.557986
(Epoch 7 / 10) train acc: 0.530000; val_acc: 0.304000
(Iteration 141 / 200) loss: 1.432189
(Epoch 8 / 10) train acc: 0.628000; val_acc: 0.339000
(Iteration 161 / 200) loss: 1.034116
(Epoch 9 / 10) train acc: 0.654000; val_acc: 0.342000
(Iteration 181 / 200) loss: 0.905796
(Epoch 10 / 10) train acc: 0.714000; val_acc: 0.331000
```

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

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

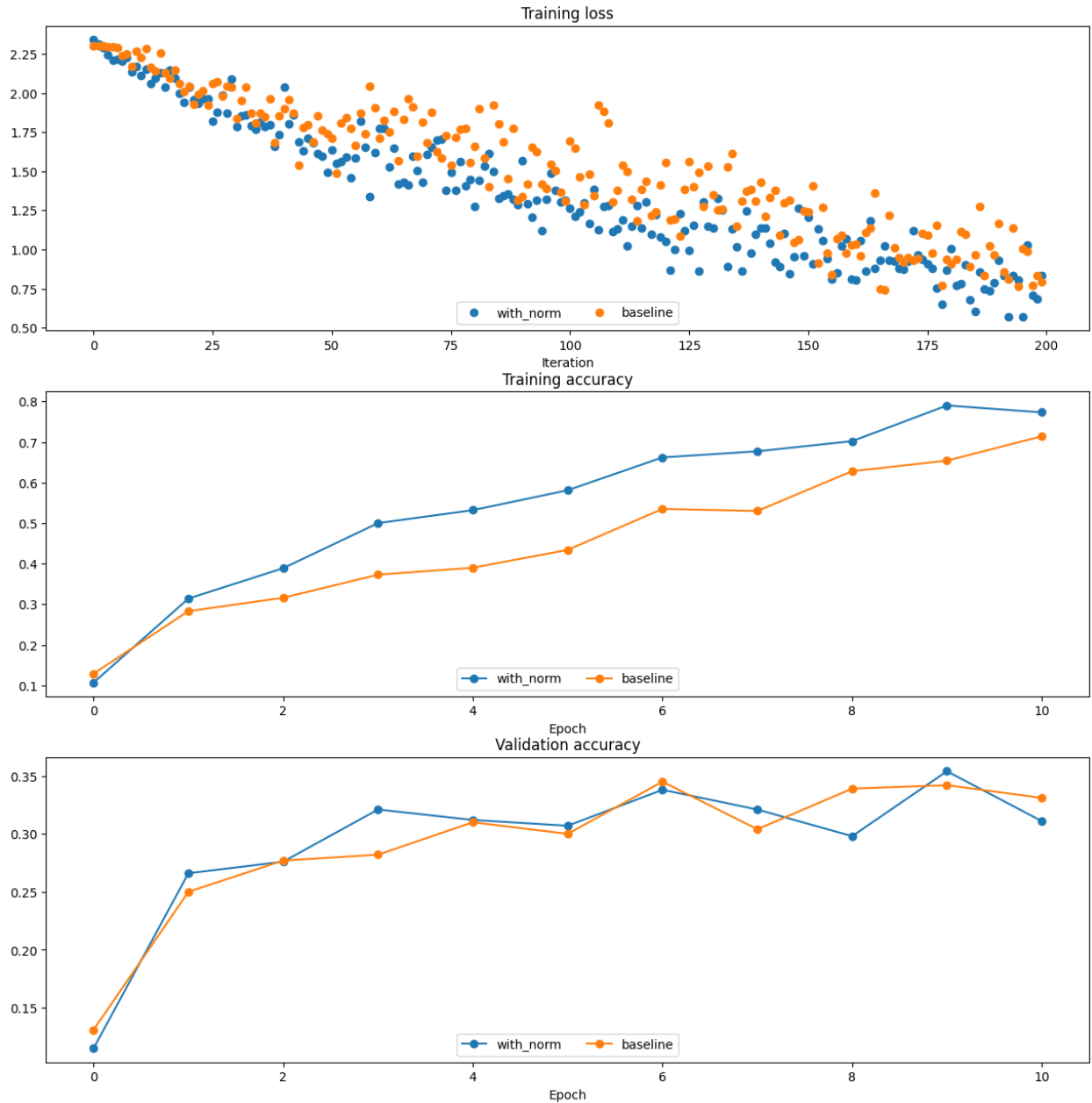
```

        label += str(labels[0])
    plt.plot(bl_plot, bl_marker, label=label)
    plt.legend(loc='lower center', ncol=num_bn+1)

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

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

```



## Batch Normalization and Initialization

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

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



```

np.random.seed(231)

# Try training a very deep net with batchnorm.
hidden_dims = [50, 50, 50, 50, 50, 50, 50]
num_train = 1000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

bn_solvers_ws = {}
solvers_ws = {}
weight_scales = np.logspace(-4, 0, num=20)
for i, weight_scale in enumerate(weight_scales):
    print('Running weight scale %d / %d' % (i + 1,
len(weight_scales)))
    bn_model = FullyConnectedNet(hidden_dims,
weight_scale=weight_scale, normalization='batchnorm')
    model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale,
normalization=None)

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

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

```

```

Running weight scale 1 / 20
Running weight scale 2 / 20
Running weight scale 3 / 20
Running weight scale 4 / 20
Running weight scale 5 / 20
Running weight scale 6 / 20
Running weight scale 7 / 20

```

```

Running weight scale 8 / 20
Running weight scale 9 / 20
Running weight scale 10 / 20
Running weight scale 11 / 20
Running weight scale 12 / 20
Running weight scale 13 / 20
Running weight scale 14 / 20
Running weight scale 15 / 20
Running weight scale 16 / 20
Running weight scale 17 / 20
Running weight scale 18 / 20
Running weight scale 19 / 20
Running weight scale 20 / 20

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

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

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

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

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

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

plt.subplot(3, 1, 3)
plt.title('Final training loss vs. weight initialization scale')
plt.xlabel('Weight initialization scale')

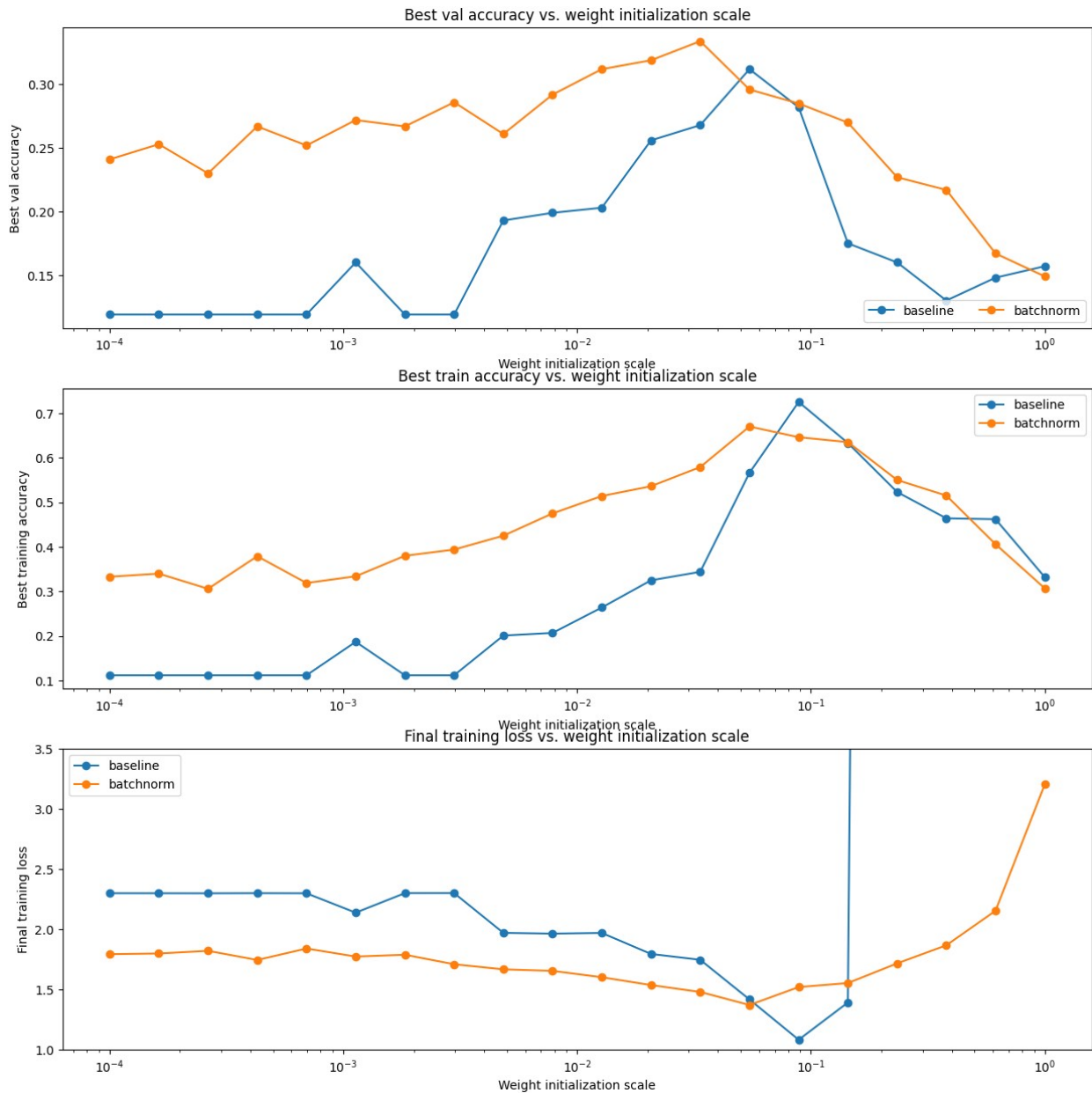
```

```

plt.ylabel('Final training loss')
plt.semilogx(weight_scales, final_train_loss, '-o', label='baseline')
plt.semilogx(weight_scales, bn_final_train_loss, '-o',
label='batchnorm')
plt.legend()
plt.gca().set_ylim(1.0, 3.5)

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

```



## Inline Question 1:

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

## Answer:

### Overall accuracy

- All things considered, batch normalization yields higher accuracy and loss outcomes. The spike may be viewed as a little overfit without batchnorm since, despite the fact that there is a spike (at weight initialization  $10^1$ ) indicating that training accuracy is greater without batchnorm, validation accuracy is the same as with batchnorm.
- Because the gradient descent stages are more focused on the loss minima, it seems sense that the batchnorm would yield superior results. Each input is linked to a network's parameters, and varying their scales might result in an unwieldy loss function architecture that emphasizes certain parameter gradients. Higher loss is therefore produced without the use of normalization.

### Robustness

- Notably, batch normalization yields results that are smoother across various initializations of the weight scale. Without it, a slight alteration in initialization might result in a notably different outcome.
- Batch normalization makes it easier to isolate centered data, which ensures better smoothness. To put it another way, compared to adjusting the weights for non-centered data, minor weight adjustments do not result in as significant output disparities.

### Stability

- Finally, even with extremely high or very low weight initialization, batch normalization yields accurate results. Large weights, such those on the scale of 100, or small weights, like  $10^4$ , make it impossible for the baseline to complete adequate gradient descent and result in skewed accuracy.
- This is due to the fact that big weight initialization produces large output values in the absence of normalization across numerous layers, whereas tiny initialization produces small values. When large error gradients build up, the weights are updated enormously, causing the gradient to burst; when low error gradients build up, the weights are updated small, causing the gradient to disappear. After each layer's raw output is normalized, the data is placed back on the scale where each data point's average is distributed from the

## Batch Normalization and Batch Size

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

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

```
def run_batchsize_experiments(normalization_mode):
    np.random.seed(231)

    # Try training a very deep net with batchnorm.
    hidden_dims = [100, 100, 100, 100, 100]
    num_train = 1000
    small_data = {
        'X_train': data['X_train'][:num_train],
        'y_train': data['y_train'][:num_train],
        'X_val': data['X_val'],
        'y_val': data['y_val'],
    }
    n_epochs=10
    weight_scale = 2e-2
    batch_sizes = [5,10,50]
    lr = 10**(-3.5)
    solver_bsize = batch_sizes[0]

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

    solver.train()

    bn_solvers = []
    for i in range(len(batch_sizes)):
        b_size=batch_sizes[i]
        print('Normalization: batch size = ',b_size)
        bn_model = FullyConnectedNet(hidden_dims,
weight_scale=weight_scale, normalization=normalization_mode)
        bn_solver = Solver(bn_model, small_data,
                        num_epochs=n_epochs, batch_size=b_size,
                        update_rule='adam',
                        optim_config={
                            'learning_rate': lr,
                        },
                        verbose=False)

        bn_solver.train()
        bn_solvers.append(bn_solver)
```

```

return bn_solvers, solver, batch_sizes

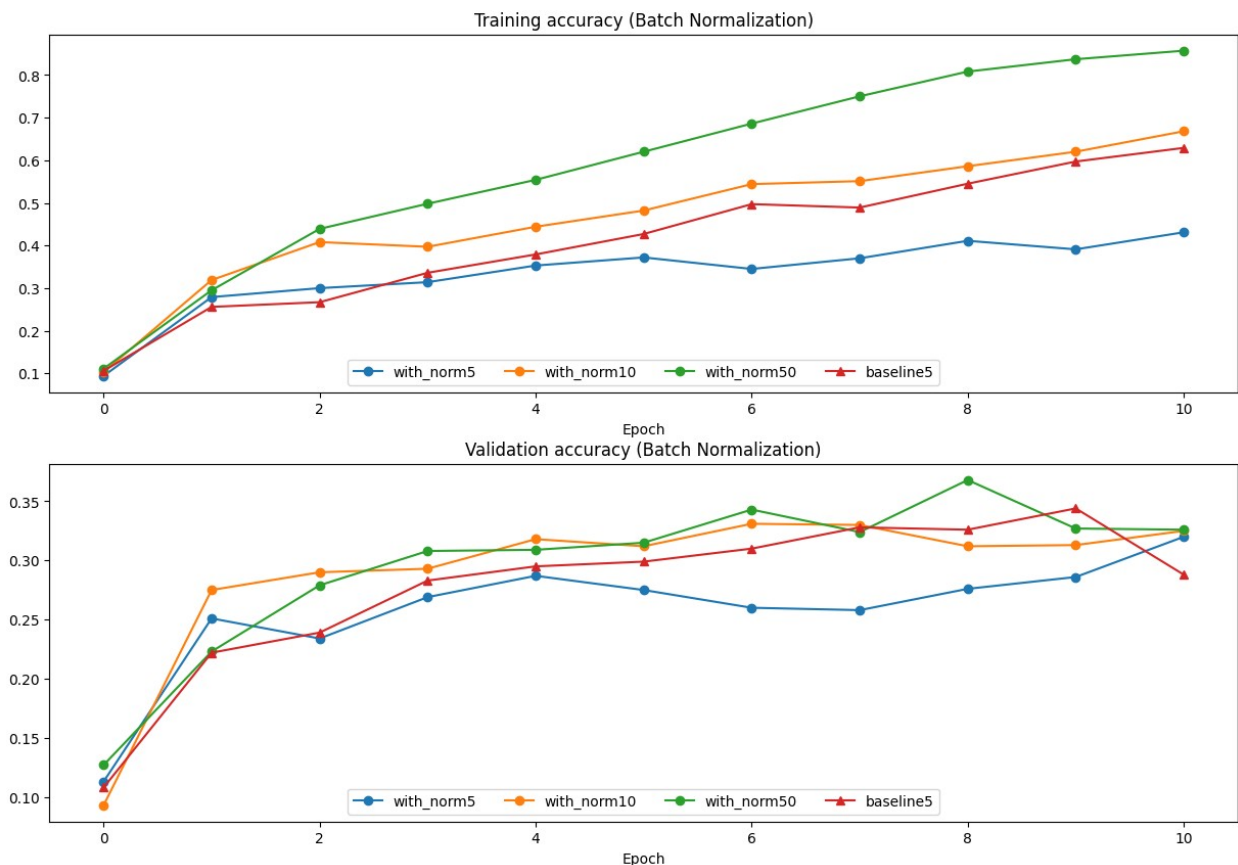
batch_sizes = [5,10,50]
bn_solvers_bsize, solver_bsize, batch_sizes =
run_batchsize_experiments('batchnorm')

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

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

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

```



## Inline Question 2:

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

### Answer:

#### Training accuracy

- It is evident that training accuracy is increased by utilizing a bigger batch size. The same batch size (i.e., 5) yields higher training accuracy without the use of batch norm as compared to baseline.
- It makes sense that training accuracy would increase as batch size increased. The model grows more adept at overfitting itself to its explicit characteristics as it observes more training samples, making classification simpler. Since the model trains on raw (precise) features rather than modified ones, it is expected to perform better without batch norm if the batch size is small (i.e., 5). Because mean and variance are estimated using a tiny batch size, the corrected features may be noisy.

#### Validation accuracy

- Though it is less reliant on batch size and yields results that are broadly comparable when there are enough epochs, the second plot displays a similar tendency (bigger batch - larger accuracy).
- Larger batch sizes lead the model to converge more quickly over a few epochs, improving accuracy; but, after a given number of epochs, the model finally converges with all batch sizes, thus performance levels out. Because the model cannot overfit the validation set as much as it can the training set, performance over time is less dependent on batch size than training accuracy.

## Layer Normalization

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

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

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

## Inline Question 3:

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

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

## Answer:

Batch normalization is the process of normalizing each unique feature based on the same characteristics (at the same indices) over several samples given various feature vectors.

Given a single feature vector, layer normalization involves normalizing each unique characteristic throughout the single sample depending on other features (at varying degrees).

Understanding the definitions makes it simple to see:

2 is analogous to layer normalization because, given a single vector of pixels, we scale each pixel with respect to all the other pixels in the same vector.

3 is analogous to batch normalization because, given multiple images, we shift each image with respect to all the other images (the mean of them).

## Layer Normalization: Implementation

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

Here's what you need to do:

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

Run the cell below to check your results.

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

Run the second cell below to check your results.



- Modify `cs231n/classifiers/fc_net.py` to add layer normalization to the `FullyConnectedNet`. When the `normalization` flag is set to `"layernorm"` in the constructor, you should insert a layer normalization layer before each ReLU nonlinearity.

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

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

# Simulate the forward pass for a two-layer network.
np.random.seed(231)
N, D1, D2, D3 = 4, 50, 60, 3
X = np.random.randn(N, D1)
W1 = np.random.randn(D1, D2)
W2 = np.random.randn(D2, D3)
a = np.maximum(0, X.dot(W1)).dot(W2)

print('Before layer normalization:')
print_mean_std(a,axis=1)

gamma = np.ones(D3)
beta = np.zeros(D3)

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

gamma = np.asarray([3.0,3.0,3.0])
beta = np.asarray([5.0,5.0,5.0])

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

Before layer normalization:
means: [-59.06673243 -47.60782686 -43.31137368 -26.40991744]
stds:  [10.07429373 28.39478981 35.28360729  4.01831507]

After layer normalization (gamma=1, beta=0)
means: [ 2.59052039e-16  0.00000000e+00  2.22044605e-16 -
5.55111512e-16]
stds:  [0.99999995 0.99999999 1.          0.99999969]

After layer normalization (gamma= [3. 3. 3.] , beta= [5. 5. 5.] )
means: [5. 5. 5. 5.]
stds:  [2.99999985 2.99999998 2.99999999 2.99999907]
```

```

# Gradient check batchnorm backward pass.
np.random.seed(231)
N, D = 4, 5
x = 5 * np.random.randn(N, D) + 12
gamma = np.random.randn(D)
beta = np.random.randn(D)
dout = np.random.randn(N, D)

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

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

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

# You should expect to see relative errors between 1e-12 and 1e-8.
print('dx error: ', rel_error(dx_num, dx))
print('dgamma error: ', rel_error(da_num, dgamma))
print('dbeta error: ', rel_error(db_num, dbeta))

dx error:  1.433615657860454e-09
dgamma error:  4.519489546032799e-12
dbeta error:  2.276445013433725e-12

```

## Layer Normalization and Batch Size

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

```

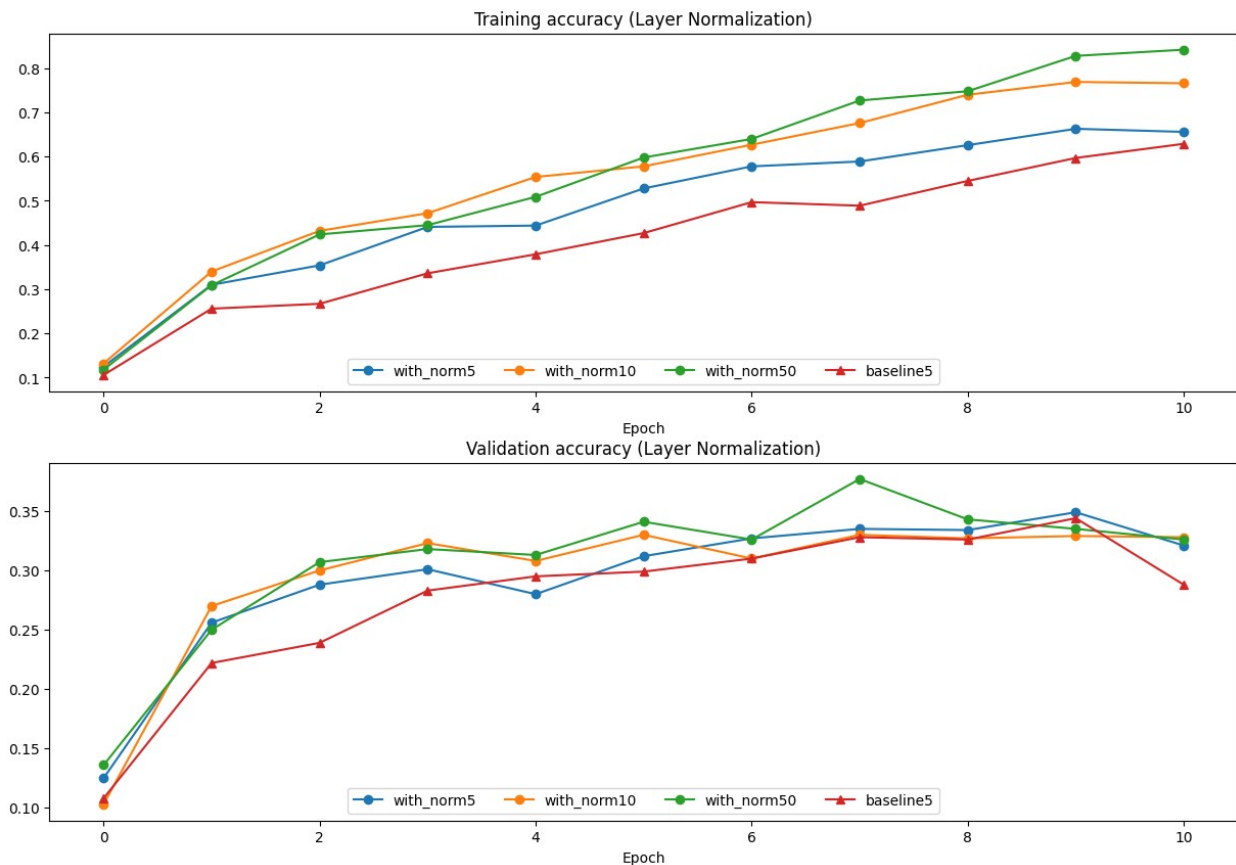
ln_solvers_bsize, solver_bsize, batch_sizes =
run_batchsize_experiments('layernorm')

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

```

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

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



## Inline Question 4:

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

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

## Answer:

1. FALSE: Having more layers won't affect the performance of layer normalization since data feature vectors are normalized per layer. By roughly centering each feature vector inside its space, normalization reduces the skewed topology of the loss function. Additionally, scaling each characteristic prevents the gradient from

disappearing or blowing up by averaging out extreme values (for example, if one is very high and another is very low).

2. TRUE: It is true that a feature vector with a small number of values will not accurately represent the mean and variance of those characteristics and will not place them on comparable scales. The performance would be loud as a result.
3. TRUE: The weights may be overly penalized by strong regularization, which would prevent them from appropriately weighting particular vector characteristics. Generally speaking, a strong regularization term tends to increase the loss and create simpler models.