

FullyConnectedNets

May 22, 2023

```
[ ]: # this mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive', force_remount=True)

# enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'CV7062610/assignments/assignment3/'
FOLDERNAME = 'CV7062610/assignments/assignment3/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

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

# this downloads the CIFAR-10 dataset to your Drive
# if it doesn't already exist.
%cd drive/My\ Drive/$FOLDERNAME/CV7062610/datasets/
!bash get_datasets.sh
%cd /content
```

Mounted at /content/drive
/content/drive/My Drive/CV7062610/assignments/assignment3/CV7062610/datasets
/content

```
[ ]: # As usual, a bit of setup
from __future__ import print_function
import time
import numpy as np
import matplotlib.pyplot as plt
from CV7062610.classifiers.fc_net import *
from CV7062610.data_utils import get_CIFAR10_data
from CV7062610.gradient_check import eval_numerical_gradient,
    eval_numerical_gradient_array
from CV7062610.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))))
```

===== You can safely ignore the message below if you are NOT working on ConvolutionalNetworks.ipynb =====

You will need to compile a Cython extension for a portion of this assignment.

The instructions to do this will be given in a section of the notebook below.

There will be an option for Colab users and another for Jupyter (local) users.

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

```
data = get_CIFAR10_data()
for k, v in list(data.items()):
    print('%s: ' % k, v.shape)
```

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

1 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.

2 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 CV7062610/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 $e-8$.

```
[ ]: from CV7062610.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.

```
[ ]: 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()
print()

plt.subplot(3, 1, 1)
plt.title('Training loss')
plt.xlabel('Iteration')

plt.subplot(3, 1, 2)
plt.title('Training accuracy')
plt.xlabel('Epoch')

plt.subplot(3, 1, 3)
plt.title('Validation accuracy')
plt.xlabel('Epoch')

for update_rule, solver in solvers.items():
    plt.subplot(3, 1, 1)
    plt.plot(solver.loss_history, 'o', label="loss_%s" % update_rule)

    plt.subplot(3, 1, 2)
    plt.plot(solver.train_acc_history, '-o', label="train_acc_%s" % update_rule)

    plt.subplot(3, 1, 3)
    plt.plot(solver.val_acc_history, '-o', label="val_acc_%s" % update_rule)

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

```

```

running with  sgd
(Iteration 1 / 200) loss: 2.631166
(Epoch 0 / 5) train acc: 0.082000; val_acc: 0.077000
(Iteration 11 / 200) loss: 2.286968
(Iteration 21 / 200) loss: 2.258384
(Iteration 31 / 200) loss: 2.161800

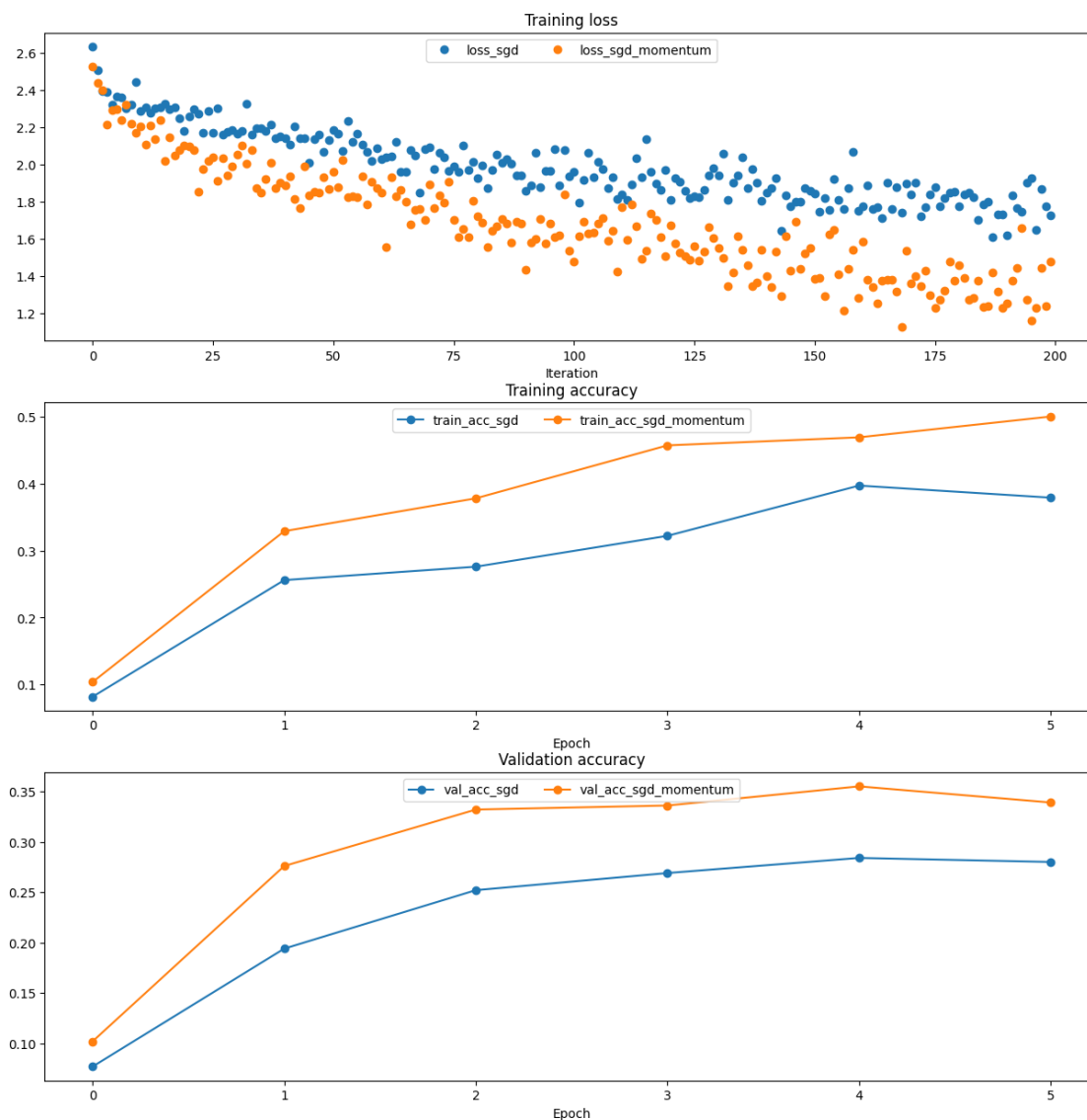
```

(Epoch 1 / 5) train acc: 0.256000; val_acc: 0.194000
(Iteration 41 / 200) loss: 2.138508
(Iteration 51 / 200) loss: 2.184243
(Iteration 61 / 200) loss: 2.025917
(Iteration 71 / 200) loss: 2.090448
(Epoch 2 / 5) train acc: 0.276000; val_acc: 0.252000
(Iteration 81 / 200) loss: 1.924731
(Iteration 91 / 200) loss: 1.855421
(Iteration 101 / 200) loss: 1.957003
(Iteration 111 / 200) loss: 1.838934
(Epoch 3 / 5) train acc: 0.322000; val_acc: 0.269000
(Iteration 121 / 200) loss: 1.805770
(Iteration 131 / 200) loss: 1.940704
(Iteration 141 / 200) loss: 1.846346
(Iteration 151 / 200) loss: 1.843771
(Epoch 4 / 5) train acc: 0.397000; val_acc: 0.284000
(Iteration 161 / 200) loss: 1.775256
(Iteration 171 / 200) loss: 1.838702
(Iteration 181 / 200) loss: 1.772504
(Iteration 191 / 200) loss: 1.618585
(Epoch 5 / 5) train acc: 0.379000; val_acc: 0.280000

running with sgd_momentum

(Iteration 1 / 200) loss: 2.525214
(Epoch 0 / 5) train acc: 0.104000; val_acc: 0.102000
(Iteration 11 / 200) loss: 2.203895
(Iteration 21 / 200) loss: 2.096467
(Iteration 31 / 200) loss: 2.052350
(Epoch 1 / 5) train acc: 0.329000; val_acc: 0.276000
(Iteration 41 / 200) loss: 1.884967
(Iteration 51 / 200) loss: 1.957779
(Iteration 61 / 200) loss: 1.846180
(Iteration 71 / 200) loss: 1.892333
(Epoch 2 / 5) train acc: 0.378000; val_acc: 0.332000
(Iteration 81 / 200) loss: 1.718259
(Iteration 91 / 200) loss: 1.430624
(Iteration 101 / 200) loss: 1.474311
(Iteration 111 / 200) loss: 1.767637
(Epoch 3 / 5) train acc: 0.457000; val_acc: 0.336000
(Iteration 121 / 200) loss: 1.670887
(Iteration 131 / 200) loss: 1.550273
(Iteration 141 / 200) loss: 1.396245
(Iteration 151 / 200) loss: 1.385227
(Epoch 4 / 5) train acc: 0.469000; val_acc: 0.355000
(Iteration 161 / 200) loss: 1.585418
(Iteration 171 / 200) loss: 1.360043
(Iteration 181 / 200) loss: 1.458696
(Iteration 191 / 200) loss: 1.250292

(Epoch 5 / 5) train acc: 0.500000; val_acc: 0.339000



3 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 `CV7062610/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.” COURSEARA: Neural Networks for Machine Learning 4 (2012).

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

```
[ ]: # Test RMSProp implementation
from CV7062610.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
```

```
[ ]: # Test Adam implementation
from CV7062610.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,],
[ 0.54159906, 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:

```

[ ]: 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()

plt.subplot(3, 1, 1)
plt.title('Training loss')
plt.xlabel('Iteration')

plt.subplot(3, 1, 2)

```



```

plt.title('Training accuracy')
plt.xlabel('Epoch')

plt.subplot(3, 1, 3)
plt.title('Validation accuracy')
plt.xlabel('Epoch')

for update_rule, solver in list(solvers.items()):
    plt.subplot(3, 1, 1)
    plt.plot(solver.loss_history, 'o', label=update_rule)

    plt.subplot(3, 1, 2)
    plt.plot(solver.train_acc_history, '-o', label=update_rule)

    plt.subplot(3, 1, 3)
    plt.plot(solver.val_acc_history, '-o', label=update_rule)

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

```

```

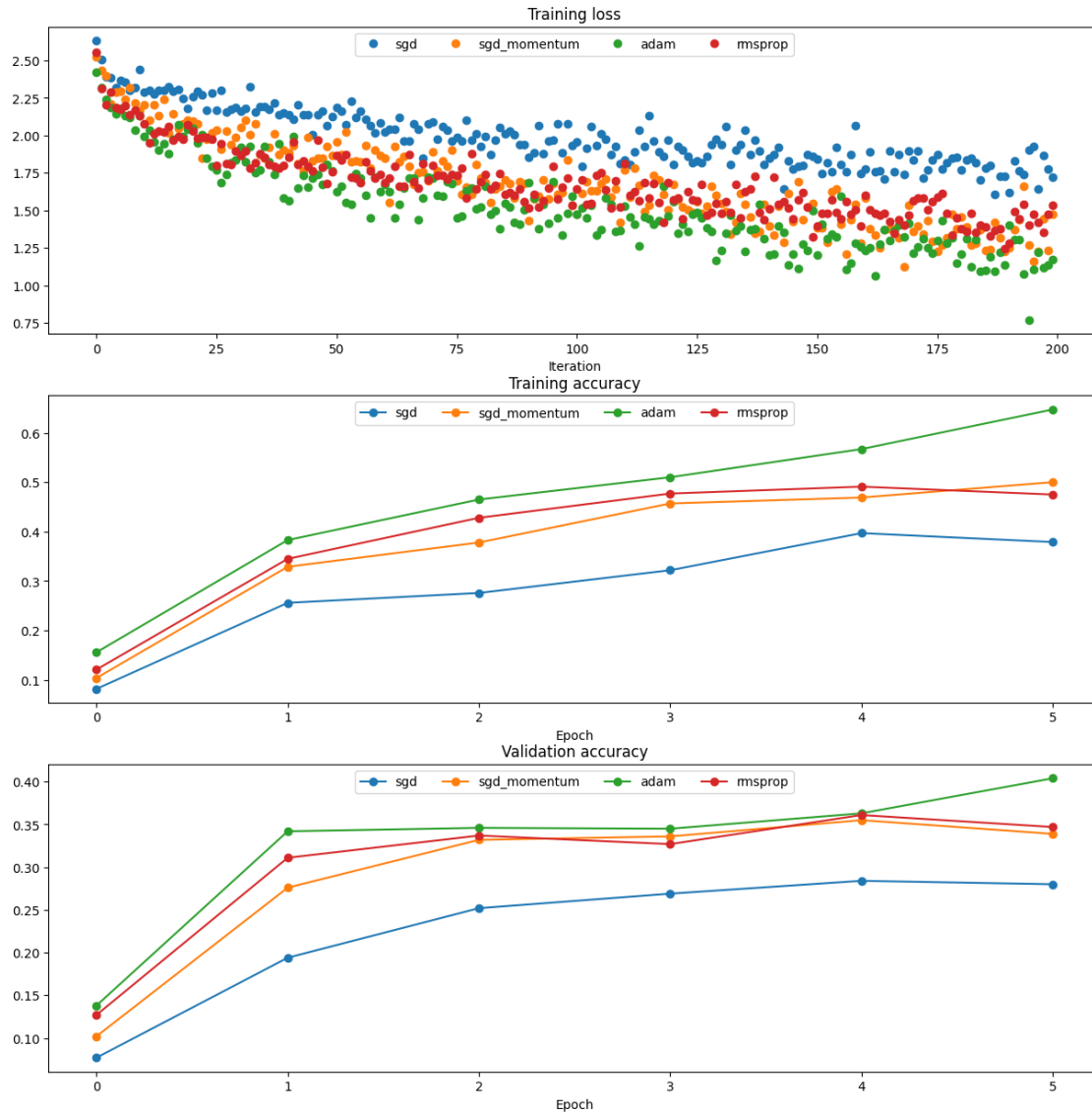
running with adam
(Iteration 1 / 200) loss: 2.423783
(Epoch 0 / 5) train acc: 0.156000; val_acc: 0.138000
(Iteration 11 / 200) loss: 1.991330
(Iteration 21 / 200) loss: 2.045923
(Iteration 31 / 200) loss: 1.921859
(Epoch 1 / 5) train acc: 0.383000; val_acc: 0.342000
(Iteration 41 / 200) loss: 1.565561
(Iteration 51 / 200) loss: 1.623416
(Iteration 61 / 200) loss: 1.674059
(Iteration 71 / 200) loss: 1.717843
(Epoch 2 / 5) train acc: 0.465000; val_acc: 0.346000
(Iteration 81 / 200) loss: 1.521725
(Iteration 91 / 200) loss: 1.686438
(Iteration 101 / 200) loss: 1.477843
(Iteration 111 / 200) loss: 1.531259
(Epoch 3 / 5) train acc: 0.510000; val_acc: 0.345000
(Iteration 121 / 200) loss: 1.559150
(Iteration 131 / 200) loss: 1.236174
(Iteration 141 / 200) loss: 1.205578
(Iteration 151 / 200) loss: 1.201568
(Epoch 4 / 5) train acc: 0.567000; val_acc: 0.363000
(Iteration 161 / 200) loss: 1.234008
(Iteration 171 / 200) loss: 1.217424

```

(Iteration 181 / 200) loss: 1.210178
(Iteration 191 / 200) loss: 1.404529
(Epoch 5 / 5) train acc: 0.647000; val_acc: 0.404000

running with rmsprop

(Iteration 1 / 200) loss: 2.553076
(Epoch 0 / 5) train acc: 0.121000; val_acc: 0.127000
(Iteration 11 / 200) loss: 2.077437
(Iteration 21 / 200) loss: 2.029830
(Iteration 31 / 200) loss: 1.867870
(Epoch 1 / 5) train acc: 0.345000; val_acc: 0.311000
(Iteration 41 / 200) loss: 1.855108
(Iteration 51 / 200) loss: 1.865094
(Iteration 61 / 200) loss: 1.687106
(Iteration 71 / 200) loss: 1.800815
(Epoch 2 / 5) train acc: 0.428000; val_acc: 0.337000
(Iteration 81 / 200) loss: 1.693280
(Iteration 91 / 200) loss: 1.517699
(Iteration 101 / 200) loss: 1.641903
(Iteration 111 / 200) loss: 1.814286
(Epoch 3 / 5) train acc: 0.477000; val_acc: 0.327000
(Iteration 121 / 200) loss: 1.571749
(Iteration 131 / 200) loss: 1.477997
(Iteration 141 / 200) loss: 1.509328
(Iteration 151 / 200) loss: 1.397429
(Epoch 4 / 5) train acc: 0.491000; val_acc: 0.361000
(Iteration 161 / 200) loss: 1.603371
(Iteration 171 / 200) loss: 1.565835
(Iteration 181 / 200) loss: 1.377710
(Iteration 191 / 200) loss: 1.280735
(Epoch 5 / 5) train acc: 0.475000; val_acc: 0.347000



3.1 Inline Question 1:

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?

3.2 Answer:

AdaGrad accumulates the squared gradients over time. As training progresses, the accumulated squared gradients in cache tend to increase. Consequently, the denominator np.sqrt(cache) becomes larger, leading to a smaller learning rate for the parameter update.

The effect of a smaller learning rate is that the updates become smaller as well. This effect can result in slower learning since the parameter w is updated with smaller steps, causing a longer time to converge to the optimal solution.

However, Adam is not suffering from the same issue as AdaGrad. This is because in Adam, the adaptive learning rate is determined by a combination of both the first moment estimate (m) and the second moment normalized estimate (v_hat). In addition, Adam updates the parameters based on a dynamic learning rate that is adjusted according to both the magnitude of the gradients and their direction. This allows Adam to adjust the parameters with appropriate step-sizes, which can help it overcome the problem of slow learning associated with AdaGrad.

4 Train a good model!

Train the best 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 net or a cnn net.

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.

```
[16]: 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 in #
# the best_model variable. #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

best_val = 0.0

# Define the hyperparameters to be searched over
hidden_sizes = [100, 100, 100]
num_epochs = 10
batch_size = 200
learning_rates = [1e-2, 5e-3, 1e-3, 5e-4, 1e-4]
weight_scales = [1e-2, 5e-3, 1e-3, 5e-4, 1e-4]
reg_strengths = [1e-3, 5e-3, 1e-4, 5e-4, 1e-5, 5e-5]
dropout_p = [0.625, 0.75, 0.85]
num_iter = 10

# Perform random search over hyperparameters
for i in range(num_iter):
    lr = learning_rates[np.random.randint(len(learning_rates))]
```

```

ws = weight_scales[np.random.randint(len(weight_scales))]
reg = reg_strengths[np.random.randint(len(reg_strengths))]
for drop in dropout_p:

    # Create a new FullyConnectedNet with the chosen hyperparameters
    model = FullyConnectedNet(hidden_sizes, weight_scale=ws, reg=reg,
    ↪normalization='batchnorm', dropout=drop)

    # Train the model
    solver = Solver(model, data,
                    num_epochs=num_epochs, batch_size=batch_size,
                    update_rule='adam',
                    optim_config={
                        'learning_rate': lr
                    },
                    verbose=False)

    solver.train()
    val_accuracy = solver.best_val_acc

    # Update the best validation accuracy and best model if necessary
    if best_val < val_accuracy:
        best_val = val_accuracy
        best_model = model

    # Print results
    print('combination - lr %f, ws %f, reg %f, dropout %f, val accuracy: %f'
    ↪% (
        lr, ws, reg, drop, val_accuracy))

print('best validation accuracy achieved: %f' % best_val)

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

```

```

combination - lr 0.010000, ws 0.005000, reg 0.001000, dropout 0.625000, val
accuracy: 0.382000
combination - lr 0.010000, ws 0.005000, reg 0.001000, dropout 0.750000, val
accuracy: 0.397000
combination - lr 0.010000, ws 0.005000, reg 0.001000, dropout 0.850000, val
accuracy: 0.408000
combination - lr 0.001000, ws 0.010000, reg 0.000010, dropout 0.625000, val
accuracy: 0.520000
combination - lr 0.001000, ws 0.010000, reg 0.000010, dropout 0.750000, val
accuracy: 0.534000
combination - lr 0.001000, ws 0.010000, reg 0.000010, dropout 0.850000, val

```

accuracy: 0.557000
 combination - lr 0.010000, ws 0.001000, reg 0.000050, dropout 0.625000, val
 accuracy: 0.478000
 combination - lr 0.010000, ws 0.001000, reg 0.000050, dropout 0.750000, val
 accuracy: 0.497000
 combination - lr 0.010000, ws 0.001000, reg 0.000050, dropout 0.850000, val
 accuracy: 0.514000
 combination - lr 0.000500, ws 0.000100, reg 0.000100, dropout 0.625000, val
 accuracy: 0.505000
 combination - lr 0.000500, ws 0.000100, reg 0.000100, dropout 0.750000, val
 accuracy: 0.531000
 combination - lr 0.000500, ws 0.000100, reg 0.000100, dropout 0.850000, val
 accuracy: 0.555000
 combination - lr 0.000100, ws 0.001000, reg 0.000050, dropout 0.625000, val
 accuracy: 0.491000
 combination - lr 0.000100, ws 0.001000, reg 0.000050, dropout 0.750000, val
 accuracy: 0.529000
 combination - lr 0.000100, ws 0.001000, reg 0.000050, dropout 0.850000, val
 accuracy: 0.532000
 combination - lr 0.001000, ws 0.001000, reg 0.000500, dropout 0.625000, val
 accuracy: 0.505000
 combination - lr 0.001000, ws 0.001000, reg 0.000500, dropout 0.750000, val
 accuracy: 0.529000
 combination - lr 0.001000, ws 0.001000, reg 0.000500, dropout 0.850000, val
 accuracy: 0.550000
 combination - lr 0.010000, ws 0.010000, reg 0.000010, dropout 0.625000, val
 accuracy: 0.492000
 combination - lr 0.010000, ws 0.010000, reg 0.000010, dropout 0.750000, val
 accuracy: 0.517000
 combination - lr 0.010000, ws 0.010000, reg 0.000010, dropout 0.850000, val
 accuracy: 0.529000
 combination - lr 0.010000, ws 0.001000, reg 0.000100, dropout 0.625000, val
 accuracy: 0.454000
 combination - lr 0.010000, ws 0.001000, reg 0.000100, dropout 0.750000, val
 accuracy: 0.472000
 combination - lr 0.010000, ws 0.001000, reg 0.000100, dropout 0.850000, val
 accuracy: 0.492000
 combination - lr 0.000500, ws 0.000500, reg 0.000010, dropout 0.625000, val
 accuracy: 0.506000
 combination - lr 0.000500, ws 0.000500, reg 0.000010, dropout 0.750000, val
 accuracy: 0.528000
 combination - lr 0.000500, ws 0.000500, reg 0.000010, dropout 0.850000, val
 accuracy: 0.562000
 combination - lr 0.005000, ws 0.000500, reg 0.005000, dropout 0.625000, val
 accuracy: 0.387000
 combination - lr 0.005000, ws 0.000500, reg 0.005000, dropout 0.750000, val
 accuracy: 0.382000
 combination - lr 0.005000, ws 0.000500, reg 0.005000, dropout 0.850000, val

```
accuracy: 0.409000  
best validation accuracy achieved: 0.562000
```

5 Test your model!

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

```
[17]: 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: 0.562  
Test set accuracy: 0.566
```

BatchNormalization

May 22, 2023

```
[1]: # this mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive', force_remount=True)

# enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'CV7062610/assignments/assignment3/'
FOLDERNAME = 'CV7062610/assignments/assignment3/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

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

# this downloads the CIFAR-10 dataset to your Drive
# if it doesn't already exist.
%cd drive/My\ Drive/$FOLDERNAME/CV7062610/datasets/
!bash get_datasets.sh
%cd /content
```

```
Mounted at /content/drive
/content/drive/My Drive/CV7062610/assignments/assignment3/CV7062610/datasets
/content
```

1 Batch Normalization

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

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

since they are output from earlier layers in the network. Even worse, during the training process the distribution of features at each layer of the network will shift as the weights of each layer are updated.

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

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

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

```
[2]: # As usual, a bit of setup
import time
import numpy as np
import matplotlib.pyplot as plt
from CV7062610.classifiers.fc_net import *
from CV7062610.data_utils import get_CIFAR10_data
from CV7062610.gradient_check import eval_numerical_gradient, \
    eval_numerical_gradient_array
from CV7062610.solver import Solver

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

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

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

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

===== You can safely ignore the message below if you are NOT working on ConvolutionalNetworks.ipynb =====

You will need to compile a Cython extension for a portion of this assignment.

The instructions to do this will be given in a section of the notebook below.

There will be an option for Colab users and another for Jupyter (local) users.

```
[3]: # Load the (preprocessed) CIFAR10 data.
data = get_CIFAR10_data()
for k, v in data.items():
    print('%s: ' % k, v.shape)
```

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

1.1 Batch normalization: forward

In the file CV7062610/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!

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

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

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

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

```

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

```

Before batch normalization:

```

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

```

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

```

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

```

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

```

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

```

[5]: *# 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]
```

1.2 Batch normalization: backward

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

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

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

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

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

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

_, cache = batchnorm_forward(x, gamma, beta, bn_param)
dx, dgamma, dbeta = batchnorm_backward(dout, cache)
#You should expect to see relative errors between 1e-13 and 1e-8
print('dx error: ', rel_error(dx_num, dx))
print('dgamma error: ', rel_error(da_num, dgamma))
print('dbeta error: ', rel_error(db_num, dbeta))
```

```
dx error:  1.7029261167605239e-09
dgamma error:  7.420414216247087e-13
dbeta error:  2.8795057655839487e-12
```

1.3 Batch normalization: alternative backward

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

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

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

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

$$\sigma = \sqrt{v + \epsilon} \qquad y_i = \frac{x_i - \mu}{\sigma} \qquad (2)$$

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.

```
[7]: 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:  5.964155941709756e-13
dgamma difference:  0.0
dbeta difference:  0.0
speedup: 1.72x

```

1.4 Fully Connected Nets with Batch Normalization

Now that you have a working implementation for batch normalization, go back to your `FullyConnectedNet` in the file `CV7062610/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 `CV7062610/layer_utils.py`. If you decide to do so, do it in the file `CV7062610/classifiers/fc_net.py`.

```

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

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

```

```

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

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

```

```

Running check with reg = 0
Initial loss: 2.2611955101340957
W1 relative error: 1.10e-04
W2 relative error: 2.85e-06
W3 relative error: 4.05e-10
b1 relative error: 3.11e-07
b2 relative error: 1.94e-08
b3 relative error: 1.01e-10
beta1 relative error: 7.33e-09
beta2 relative error: 1.89e-09
gamma1 relative error: 6.96e-09
gamma2 relative error: 1.96e-09

```

```

Running check with reg = 3.14
Initial loss: 6.996533220108303
W1 relative error: 1.98e-06
W2 relative error: 2.28e-06
W3 relative error: 1.11e-08
b1 relative error: 1.38e-08
b2 relative error: 7.99e-07
b3 relative error: 1.73e-10
beta1 relative error: 6.65e-09
beta2 relative error: 3.48e-09
gamma1 relative error: 8.80e-09
gamma2 relative error: 5.28e-09

```

2 Batchnorm for deep networks

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

```

[9]: 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: 0.313000; val_acc: 0.266000
(Iteration 21 / 200) loss: 2.039345
(Epoch 2 / 10) train acc: 0.396000; val_acc: 0.281000
(Iteration 41 / 200) loss: 2.047471
(Epoch 3 / 10) train acc: 0.484000; val_acc: 0.315000
(Iteration 61 / 200) loss: 1.739554
(Epoch 4 / 10) train acc: 0.525000; val_acc: 0.318000
(Iteration 81 / 200) loss: 1.246973
(Epoch 5 / 10) train acc: 0.595000; val_acc: 0.335000
(Iteration 101 / 200) loss: 1.354765
(Epoch 6 / 10) train acc: 0.636000; val_acc: 0.332000

```



```

(Iteration 121 / 200) loss: 1.015961
(Epoch 7 / 10) train acc: 0.672000; val_acc: 0.330000
(Iteration 141 / 200) loss: 1.148975
(Epoch 8 / 10) train acc: 0.687000; val_acc: 0.307000
(Iteration 161 / 200) loss: 0.620747
(Epoch 9 / 10) train acc: 0.785000; val_acc: 0.324000
(Iteration 181 / 200) loss: 0.743336
(Epoch 10 / 10) train acc: 0.741000; val_acc: 0.313000

```

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: 1.696059
(Epoch 6 / 10) train acc: 0.535000; val_acc: 0.345000
(Iteration 121 / 200) loss: 1.557987
(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: 0.905795
(Epoch 10 / 10) train acc: 0.712000; val_acc: 0.328000

```

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.

```

[10]: 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)

```

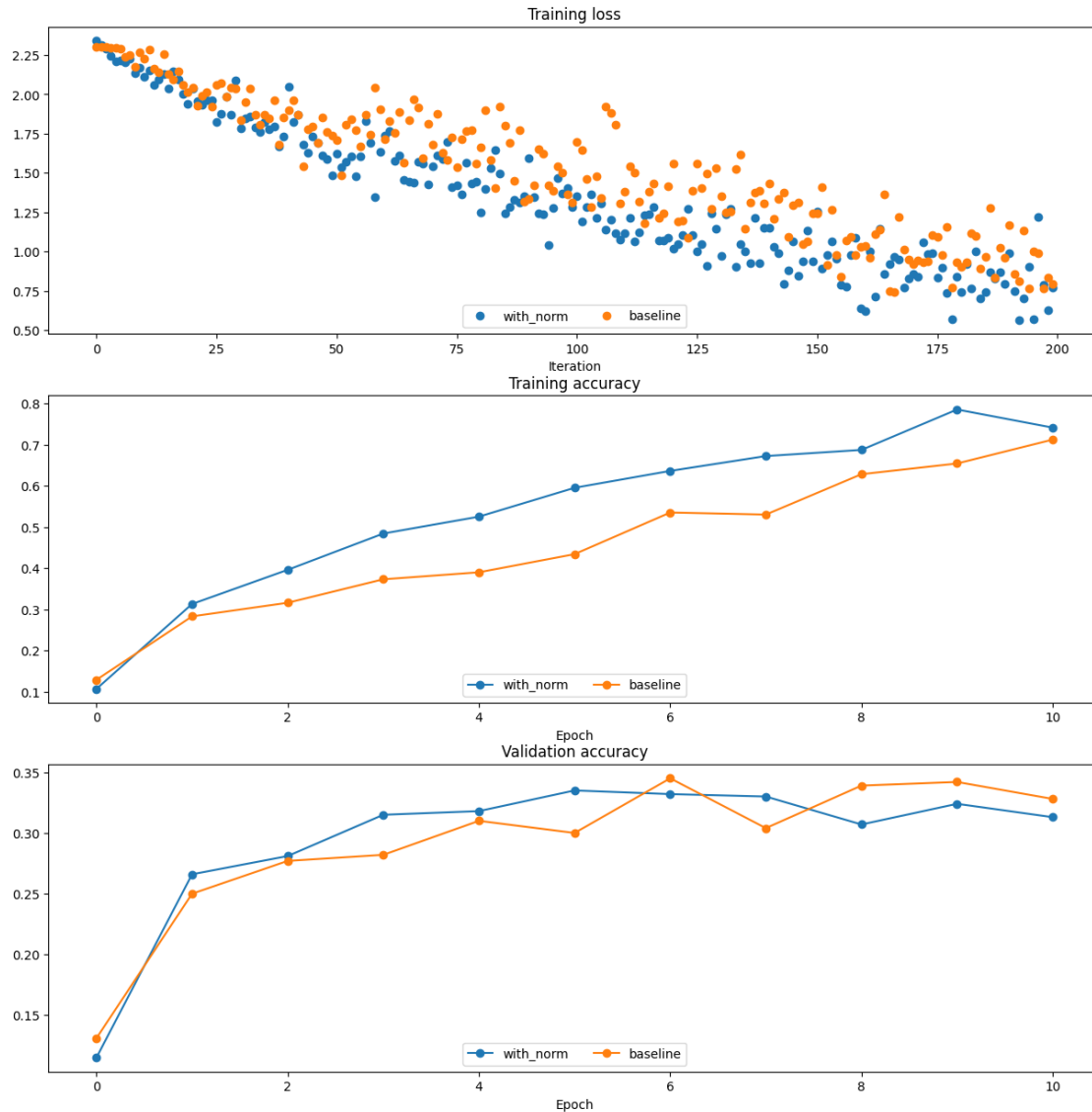
```

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()

```



3 Batch normalization and initialization

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

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

```
[11]: 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

```

```

[12]: # 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')

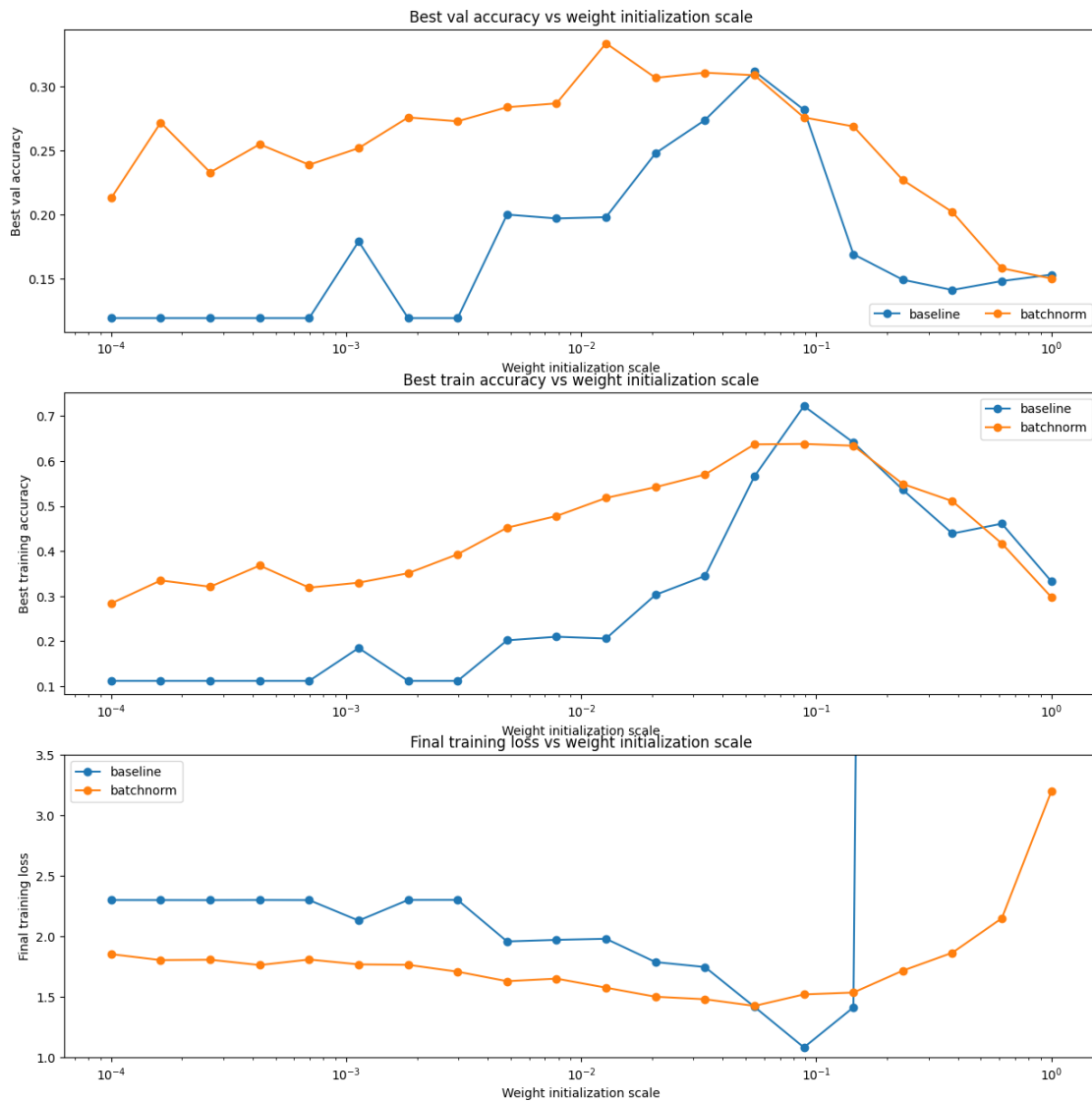
plt.subplot(3, 1, 2)
plt.title('Best train accuracy vs weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Best training accuracy')
plt.semilogx(weight_scales, best_train_accs, '-o', label='baseline')
plt.semilogx(weight_scales, bn_best_train_accs, '-o', label='batchnorm')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



3.1 Inline Question 1:

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

3.2 Answer:

For the first plot, we can see that the model experiences overfitting. In addition, the batch norm model outperforms the baseline model due to the regularizing properties of batch normalization.

The second plot, highlights the issue of vanishing gradients due to small initial weights. The baseline model exhibits greater sensitivity to this problem, resulting in lower accuracy. Finding an appropriate weight scale for the baseline model is challenging. In contrast, the batchnorm model exhibits reduced sensitivity to weight initialization, allowing for improved performance across different weight scales.

The third plot shows the issue of exploding gradients, which is particularly evident in the baseline model when using weight scale values larger than 1e-1. In contrast, the batchnorm model does not suffer from that problem and remains stable at various weight scales.

In general, batch normalization can help to avoid the problem of vanishing and exploding gradients. The reason for that is because, its normalizing the inputs to each layer during training, it helps to maintain a more stable distribution of values throughout the network. The normalization prevents very large and small values. Additionally, batch normalization has regularization effects that can help reduce overfitting. It introduces some noise to the training process by normalizing each mini-batch independently. This noise acts as a form of regularization, similar to dropout or weight decay, and can improve the generalization capability of the model.

4 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.

```
[13]: 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

```

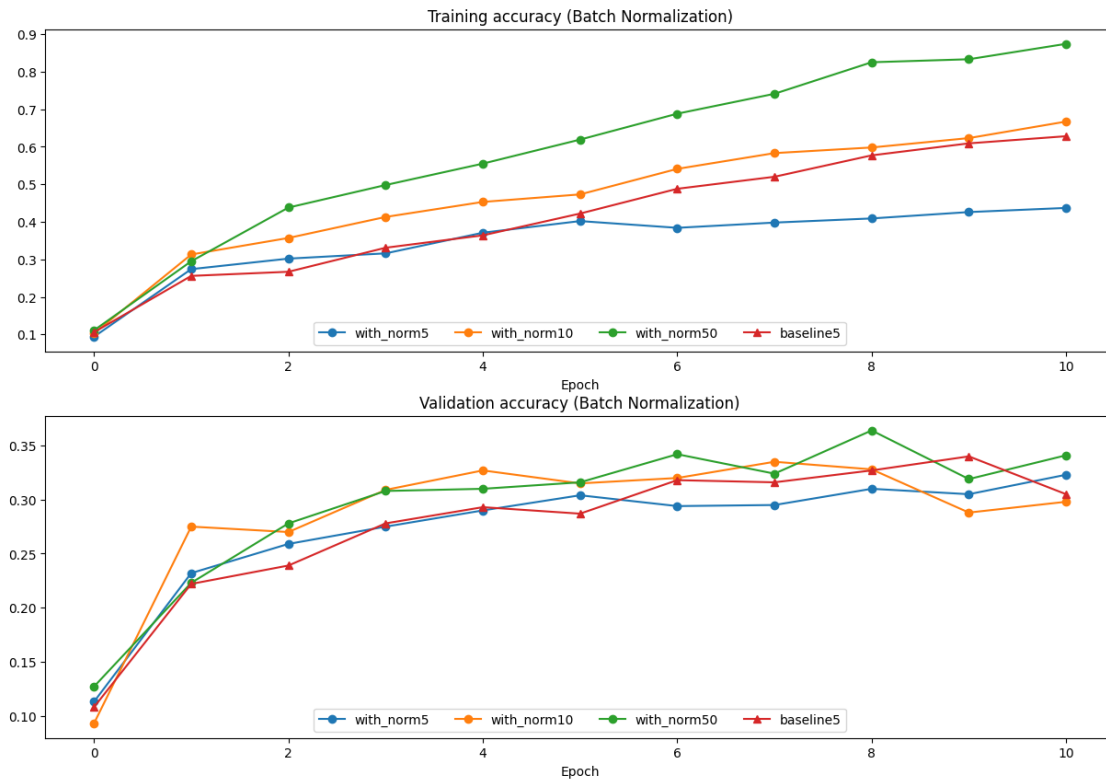
```

[14]: 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()
```



4.1 Inline Question 2:

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

4.2 Answer:

As we can see, the batch size directly affects the performance of the batch normalization. As the batch size increases the results improve. It doesn't mean that we need to put the biggest value we can. What I mean is, the direct link between the batch size and the performance of the network is clear. When the batch size is very small like "with_norm5" in our results, even the baseline model outperforms the batchnorm model. The cause of this is that we are trying to calculate the statistics of a batch, i.e., mean and variance, in order to find an approximation of the statistics of the entire dataset. Hence, when using small batch sizes, these statistics can be very noisy. However, when using a large batch size we can obtain better approximation.

5 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.

5.1 Inline Question 3:

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

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

5.2 Answer:

option 1 is **not** directly analogous to either batch normalization or layer normalization. It represents a different type of normalization specific to the rows of pixels within an image.

Option 2 is analogous to layer normalization. Layer normalization normalizes the inputs across the features/channels of a layer, and scaling the RGB channels of each image to sum up to 1 is a similar type of normalization applied at the layer level.

Option 3 is analogous to batch normalization. Batch normalization subtracts the mean of each feature/channel within a mini-batch, in order to normalize the data distribution.

option 4 is **not** analogous to either batch normalization or layer normalization. It represents a binary thresholding operation applied to the RGB values, which is not related to normalization techniques.

6 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 `CV7062610/layers.py`, implement the forward pass for layer normalization in the function `layernorm_forward`.

Run the cell below to check your results. * In CV7062610/layers.py, implement the backward pass for layer normalization in the function layernorm_backward.

Run the second cell below to check your results. * Modify CV7062610/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.

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

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

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

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

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

Before layer normalization:

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

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

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

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

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

```
[16]: # 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
```

7 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!

```
[17]: 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()

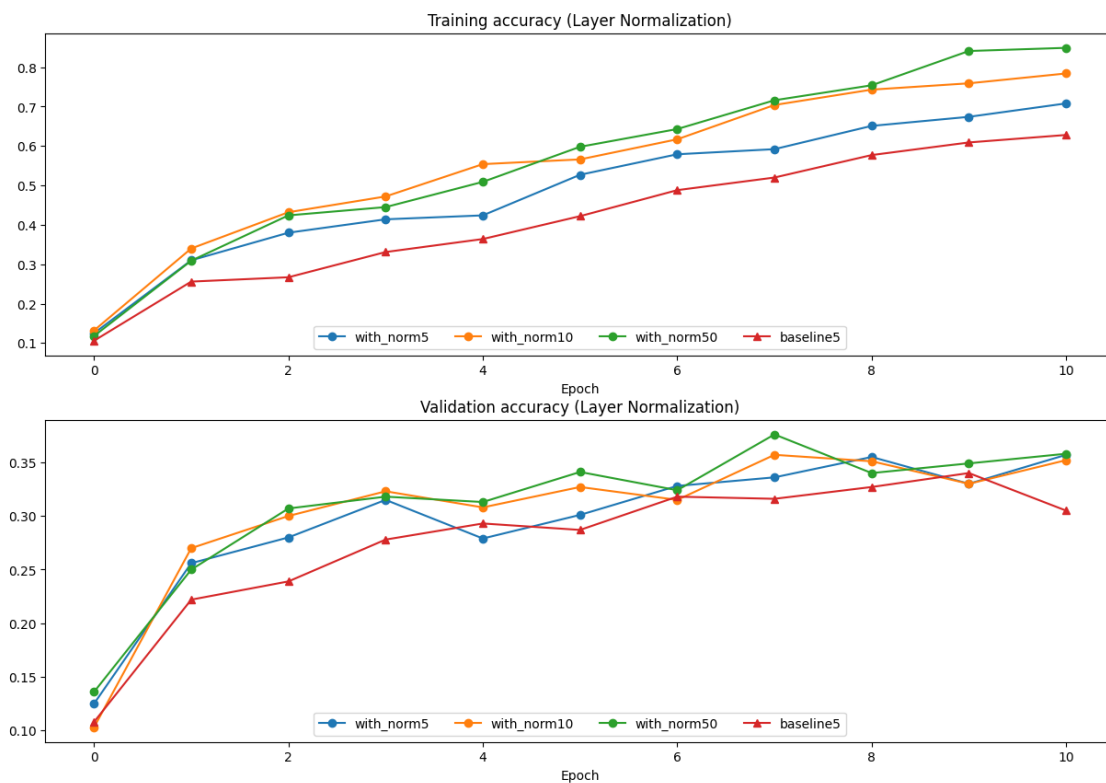
```

No normalization: batch size = 5

Normalization: batch size = 5

Normalization: batch size = 10

Normalization: batch size = 50



7.1 Inline Question 4:

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

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

7.2 Answer:

1. False, Layer normalization assumes that the statistics (mean and variance) of the input data are relatively consistent across the features. In deep networks, as the input propagates through multiple layers, the statistics of the data can change significantly. Layer normalization might struggle to normalize the data effectively in such cases, leading to suboptimal performance. But still preforms well! In our previous example, we built a 5 layers network, which can be considered as a deep network, and the resualts were well.
2. True, Having a small dimension of features affects the performance of layer normalization. If the number of features is very small, there may not be enough diversity in the data to accurately estimate the statistics. This can result in inadequate normalization and hinder the effectiveness of layer normalization.
3. True, Layer normalization introduces additional learnable parameters (gamma and beta) that allow scaling and shifting the normalized data. If we will use a high regularization term, it can restrict the freedom of these parameters, making it harder for the model to adapt and normalize the data effectively. This can lead to suboptimal performance and hinder the benefits of layer normalization.

I want to note that, the effectiveness of layer normalization can depend on the specific characteristics of the dataset and the nature of the problem being solved.

8 Spatial batch normalization: forward

In the file CV7062610/layers.py, implement the forward pass for spatial batch normalization in the function `spatial_batchnorm_forward`. Check your implementation by running the following:

```
[18]: np.random.seed(231)
# Check the training-time forward pass by checking means and variances
# of features both before and after spatial batch normalization

N, C, H, W = 2, 3, 4, 5
x = 4 * np.random.randn(N, C, H, W) + 10

print('Before spatial batch normalization:')
print('  Shape: ', x.shape)
print('  Means: ', x.mean(axis=(0, 2, 3)))
print('  Stds: ', x.std(axis=(0, 2, 3)))

# Means should be close to zero and stds close to one
gamma, beta = np.ones(C), np.zeros(C)
bn_param = {'mode': 'train'}
out, _ = spatial_batchnorm_forward(x, gamma, beta, bn_param)
print('After spatial batch normalization:')
print('  Shape: ', out.shape)
print('  Means: ', out.mean(axis=(0, 2, 3)))
print('  Stds: ', out.std(axis=(0, 2, 3)))
```

```

# Means should be close to beta and stds close to gamma
gamma, beta = np.asarray([3, 4, 5]), np.asarray([6, 7, 8])
out, _ = spatial_batchnorm_forward(x, gamma, beta, bn_param)
print('After spatial batch normalization (nontrivial gamma, beta):')
print('  Shape: ', out.shape)
print('  Means: ', out.mean(axis=(0, 2, 3)))
print('  Stds: ', out.std(axis=(0, 2, 3)))

```

Before spatial batch normalization:

```

Shape: (2, 3, 4, 5)
Means: [9.33463814 8.90909116 9.11056338]
Stds: [3.61447857 3.19347686 3.5168142 ]

```

After spatial batch normalization:

```

Shape: (2, 3, 4, 5)
Means: [ 6.18949336e-16  5.99520433e-16 -1.22124533e-16]
Stds: [0.99999962 0.99999951 0.9999996 ]

```

After spatial batch normalization (nontrivial gamma, beta):

```

Shape: (2, 3, 4, 5)
Means: [6. 7. 8.]
Stds: [2.99999885 3.99999804 4.99999798]

```

```

[19]: np.random.seed(231)
# 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, C, H, W = 10, 4, 11, 12

bn_param = {'mode': 'train'}
gamma = np.ones(C)
beta = np.zeros(C)
for t in range(50):
    x = 2.3 * np.random.randn(N, C, H, W) + 13
    spatial_batchnorm_forward(x, gamma, beta, bn_param)
bn_param['mode'] = 'test'
x = 2.3 * np.random.randn(N, C, H, W) + 13
a_norm, _ = spatial_batchnorm_forward(x, 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 spatial batch normalization (test-time):')
print('  means: ', a_norm.mean(axis=(0, 2, 3)))
print('  stds: ', a_norm.std(axis=(0, 2, 3)))

```

After spatial batch normalization (test-time):

```

means: [-0.08034406 0.07562881 0.05716371 0.04378383]
stds: [0.96718744 1.0299714 1.02887624 1.00585577]

```

9 Spatial batch normalization: backward

In the file CV7062610/layers.py, implement the backward pass for spatial batch normalization in the function `spatial_batchnorm_backward`. Run the following to check your implementation using a numeric gradient check:

```
[20]: np.random.seed(231)
N, C, H, W = 2, 3, 4, 5
x = 5 * np.random.randn(N, C, H, W) + 12
gamma = np.random.randn(C)
beta = np.random.randn(C)
dout = np.random.randn(N, C, H, W)

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

#You should expect errors of magnitudes between 1e-12~1e-06
_, cache = spatial_batchnorm_forward(x, gamma, beta, bn_param)
dx, dgamma, dbeta = spatial_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))
```

```
dx error:  2.786648197756335e-07
dgamma error:  7.0974817113608705e-12
dbeta error:  3.275608725278405e-12
```


Dropout

May 22, 2023

```
[2]: # this mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive', force_remount=True)

# enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'CV7062610/assignments/assignment3/'
FOLDERNAME = 'CV7062610/assignments/assignment3/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

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

# this downloads the CIFAR-10 dataset to your Drive
# if it doesn't already exist.
%cd drive/My\ Drive/$FOLDERNAME/CV7062610/datasets/
!bash get_datasets.sh
%cd /content
```

```
Mounted at /content/drive
/content/drive/My Drive/CV7062610/assignments/assignment3/CV7062610/datasets
/content
```

1 Dropout

Dropout [1] is a technique for regularizing neural networks by randomly setting some output activations to zero during the forward pass. In this exercise you will implement a dropout layer and modify your fully-connected network to optionally use dropout.

[1] Geoffrey E. Hinton et al, “Improving neural networks by preventing co-adaptation of feature detectors”, arXiv 2012

```
[3]: # As usual, a bit of setup
from __future__ import print_function
import time
import numpy as np
```

```

import matplotlib.pyplot as plt
from CV7062610.classifiers.fc_net import *
from CV7062610.data_utils import get_CIFAR10_data
from CV7062610.gradient_check import eval_numerical_gradient, \
    eval_numerical_gradient_array
from CV7062610.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))))

```

===== You can safely ignore the message below if you are NOT working on ConvolutionalNetworks.ipynb =====

You will need to compile a Cython extension for a portion of this assignment.

The instructions to do this will be given in a section of the notebook below.

There will be an option for Colab users and another for Jupyter (local) users.

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

```

data = get_CIFAR10_data()
for k, v in data.items():
    print('%s: ' % k, v.shape)

```

```

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

```

2 Dropout forward pass

In the file CV7062610/layers.py, implement the forward pass for dropout. Since dropout behaves differently during training and testing, make sure to implement the operation for both modes.

Once you have done so, run the cell below to test your implementation.

```
[5]: np.random.seed(231)
x = np.random.randn(500, 500) + 10

for p in [0.25, 0.4, 0.7]:
    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())
    print()
```

```
Running tests with p = 0.25
Mean of input: 10.000207878477502
Mean of train-time output: 10.014059116977283
Mean of test-time output: 10.000207878477502
Fraction of train-time output set to zero: 0.749784
Fraction of test-time output set to zero: 0.0
```

```
Running tests with p = 0.4
Mean of input: 10.000207878477502
Mean of train-time output: 9.977917658761159
Mean of test-time output: 10.000207878477502
Fraction of train-time output set to zero: 0.600796
Fraction of test-time output set to zero: 0.0
```

```
Running tests with p = 0.7
Mean of input: 10.000207878477502
Mean of train-time output: 9.987811912159426
Mean of test-time output: 10.000207878477502
Fraction of train-time output set to zero: 0.30074
Fraction of test-time output set to zero: 0.0
```

3 Dropout backward pass

In the file CV7062610/layers.py, implement the backward pass for dropout. After doing so, run the following cell to numerically gradient-check your implementation.

```
[6]: np.random.seed(231)
x = np.random.randn(10, 10) + 10
dout = np.random.randn(*x.shape)

dropout_param = {'mode': 'train', 'p': 0.2, '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], x, dout)

# Error should be around e-10 or less
print('dx relative error: ', rel_error(dx, dx_num))
```

dx relative error: 5.44560814873387e-11

3.1 Inline Question 1:

What happens if we do not divide the values being passed through inverse dropout by p in the dropout layer? Why does that happen?

3.2 Answer:

If we do not divide the values through the inverted dropout by p , the values are increased during training. This occurs because each output neuron during training is scaled by p to maintain the expected output value. However, if we do not divide by p in the backward pass, the gradients will be larger than expected.

Dividing by p during backward pass causes the gradients to retain the same scale as in testing. This keeps the gradients consistent between training and testing, which is important for proper gradient descent optimization.

4 Fully-connected nets with Dropout

In the file `CV7062610/classifiers/fc_net.py`, modify your implementation to use dropout. Specifically, if the constructor of the network receives a value that is not 1 for the `dropout` parameter, then the net should add a dropout layer immediately after every ReLU nonlinearity. After doing so, run the following to numerically gradient-check your implementation.

```
[7]: 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 dropout in [1, 0.75, 0.5]:
    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)

# Relative errors should be around e-6 or less; Note that it's fine
# if for dropout=1 you have W2 error be on the order of e-5.
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])))
print()

```

```

Running check with dropout = 1
Initial loss: 2.3004790897684924
W1 relative error: 1.48e-07
W2 relative error: 2.21e-05
W3 relative error: 3.53e-07
b1 relative error: 5.38e-09
b2 relative error: 2.09e-09
b3 relative error: 5.80e-11

```

```

Running check with dropout = 0.75
Initial loss: 2.302371489704412
W1 relative error: 1.90e-07
W2 relative error: 4.76e-06
W3 relative error: 2.60e-08
b1 relative error: 4.73e-09
b2 relative error: 1.82e-09
b3 relative error: 1.70e-10

```

```

Running check with dropout = 0.5
Initial loss: 2.3042759220785896
W1 relative error: 3.11e-07
W2 relative error: 1.84e-08
W3 relative error: 5.35e-08
b1 relative error: 5.37e-09
b2 relative error: 2.99e-09
b3 relative error: 1.13e-10

```

5 Regularization experiment

As an experiment, we will train a pair of two-layer networks on 500 training examples: one will use no dropout, and one will use a keep probability of 0.25. We will then visualize the training and validation accuracies of the two networks over time.

```

[8]: # Train two identical nets, one with dropout and one without
np.random.seed(231)
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 = [1, 0.25]
for dropout in dropout_choices:
    model = FullyConnectedNet([500], dropout=dropout)
    print(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
    print()

```

1

```

(Iteration 1 / 125) loss: 7.856643
(Epoch 0 / 25) train acc: 0.260000; val_acc: 0.184000
(Epoch 1 / 25) train acc: 0.416000; val_acc: 0.258000
(Epoch 2 / 25) train acc: 0.482000; val_acc: 0.276000
(Epoch 3 / 25) train acc: 0.532000; val_acc: 0.277000
(Epoch 4 / 25) train acc: 0.600000; val_acc: 0.271000
(Epoch 5 / 25) train acc: 0.708000; val_acc: 0.299000
(Epoch 6 / 25) train acc: 0.722000; val_acc: 0.282000
(Epoch 7 / 25) train acc: 0.832000; val_acc: 0.255000
(Epoch 8 / 25) train acc: 0.880000; val_acc: 0.268000
(Epoch 9 / 25) train acc: 0.902000; val_acc: 0.277000
(Epoch 10 / 25) train acc: 0.898000; val_acc: 0.261000
(Epoch 11 / 25) train acc: 0.924000; val_acc: 0.263000
(Epoch 12 / 25) train acc: 0.960000; val_acc: 0.300000
(Epoch 13 / 25) train acc: 0.972000; val_acc: 0.314000
(Epoch 14 / 25) train acc: 0.972000; val_acc: 0.310000
(Epoch 15 / 25) train acc: 0.974000; val_acc: 0.314000
(Epoch 16 / 25) train acc: 0.994000; val_acc: 0.303000
(Epoch 17 / 25) train acc: 0.970000; val_acc: 0.304000

```

```
(Epoch 18 / 25) train acc: 0.992000; val_acc: 0.312000
(Epoch 19 / 25) train acc: 0.992000; val_acc: 0.309000
(Epoch 20 / 25) train acc: 0.992000; val_acc: 0.289000
(Iteration 101 / 125) loss: 0.001969
(Epoch 21 / 25) train acc: 0.996000; val_acc: 0.291000
(Epoch 22 / 25) train acc: 1.000000; val_acc: 0.306000
(Epoch 23 / 25) train acc: 0.996000; val_acc: 0.309000
(Epoch 24 / 25) train acc: 0.998000; val_acc: 0.314000
(Epoch 25 / 25) train acc: 0.998000; val_acc: 0.305000
```

0.25

```
(Iteration 1 / 125) loss: 17.318478
(Epoch 0 / 25) train acc: 0.230000; val_acc: 0.177000
(Epoch 1 / 25) train acc: 0.378000; val_acc: 0.243000
(Epoch 2 / 25) train acc: 0.402000; val_acc: 0.254000
(Epoch 3 / 25) train acc: 0.502000; val_acc: 0.276000
(Epoch 4 / 25) train acc: 0.528000; val_acc: 0.298000
(Epoch 5 / 25) train acc: 0.562000; val_acc: 0.296000
(Epoch 6 / 25) train acc: 0.626000; val_acc: 0.291000
(Epoch 7 / 25) train acc: 0.622000; val_acc: 0.297000
(Epoch 8 / 25) train acc: 0.688000; val_acc: 0.313000
(Epoch 9 / 25) train acc: 0.712000; val_acc: 0.297000
(Epoch 10 / 25) train acc: 0.724000; val_acc: 0.306000
(Epoch 11 / 25) train acc: 0.768000; val_acc: 0.307000
(Epoch 12 / 25) train acc: 0.774000; val_acc: 0.284000
(Epoch 13 / 25) train acc: 0.828000; val_acc: 0.308000
(Epoch 14 / 25) train acc: 0.812000; val_acc: 0.346000
(Epoch 15 / 25) train acc: 0.850000; val_acc: 0.338000
(Epoch 16 / 25) train acc: 0.844000; val_acc: 0.307000
(Epoch 17 / 25) train acc: 0.858000; val_acc: 0.302000
(Epoch 18 / 25) train acc: 0.860000; val_acc: 0.318000
(Epoch 19 / 25) train acc: 0.884000; val_acc: 0.316000
(Epoch 20 / 25) train acc: 0.862000; val_acc: 0.315000
(Iteration 101 / 125) loss: 4.293572
(Epoch 21 / 25) train acc: 0.886000; val_acc: 0.330000
(Epoch 22 / 25) train acc: 0.898000; val_acc: 0.314000
(Epoch 23 / 25) train acc: 0.934000; val_acc: 0.323000
(Epoch 24 / 25) train acc: 0.918000; val_acc: 0.322000
(Epoch 25 / 25) train acc: 0.922000; val_acc: 0.324000
```

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

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

```

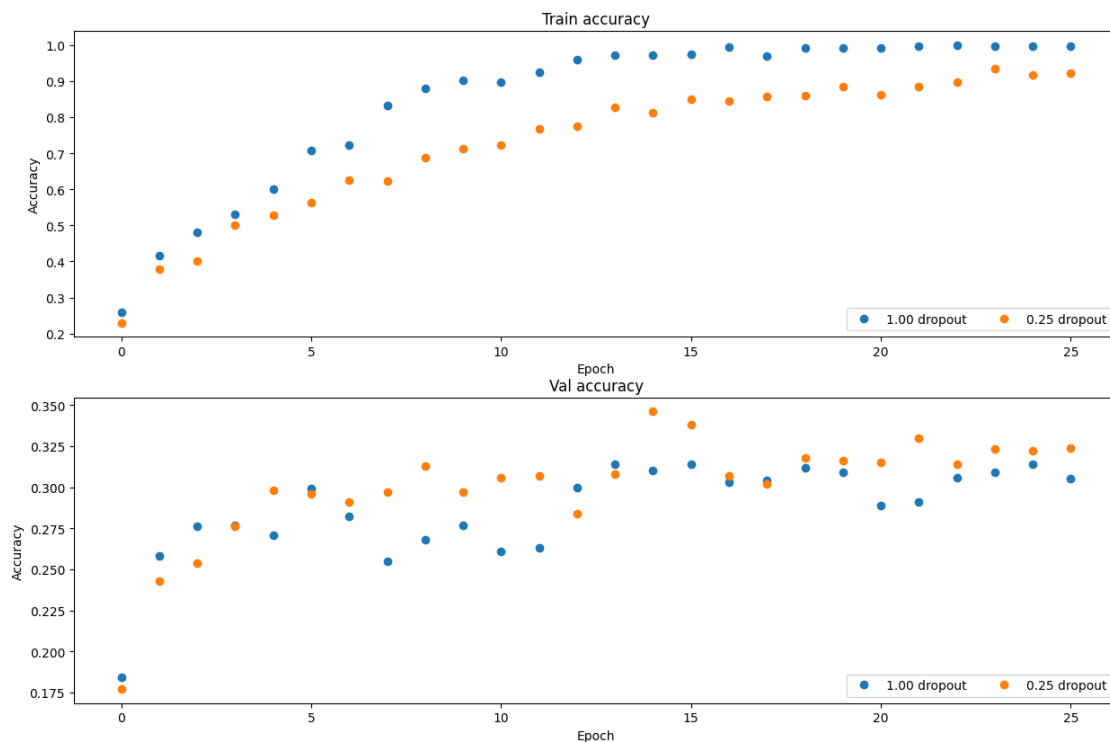
train_accs.append(solver.train_acc_history[-1])
val_accs.append(solver.val_acc_history[-1])

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()

```



5.1 Inline Question 2:

Compare the validation and training accuracies with and without dropout – what do your results suggest about dropout as a regularizer?

5.2 Answer:

The findings indicate that the model is experiencing overfitting to the training data.

Without Dropout, there is a significant difference between the accuracy achieved on the training data and the accuracy on the validation data. This disparity suggests that the model is memorizing the training examples and failing to generalize well to new, unseen data. In other words the model is overfitting.

Dropout address this issue by introducing regularization to the model. It aims to reduce overfitting by randomly deactivating a fraction of the neurons during training, by setting zeroing them, and forcing the model to rely on different combinations of neurons and preventing excessive specialization. In simple words, we want to force the model to use different set of neurons each time so, it will encourage the model to learn more robust and generalized features, thereby reducing overfitting.

By using Dropout, as we can see above, the training accuracy decreases slightly, but the validation accuracy improves. This indicates that Dropout has helped the model to perform better on the validation set. In other words, its ability to generalize to unseen data.

5.3 Inline Question 3:

Suppose we are training a deep fully-connected network for image classification, with dropout after hidden layers (parameterized by keep probability p). If we are concerned about overfitting, how should we modify p (if at all) when we decide to decrease the size of the hidden layers (that is, the number of nodes in each layer)?

5.4 Answer:

Reducing the size of hidden layers decreases the capacity and complexity of the model, which can help to prevent overfitting. However, this reduction may also reduce the ability of the model to learn and generalize well to the data. By increasing the keep probability (p) of the dropout, fewer neurons are dropped during training, giving a stronger regularization effect. This can help to obtain a more robust representation of the model and prevent overfitting.

PyTorch

May 22, 2023

```
[1]: # this mounts your Google Drive to the Colab VM.
from google.colab import drive
drive.mount('/content/drive', force_remount=True)

# enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'CV7062610/assignments/assignment3/'
FOLDERNAME = 'CV7062610/assignments/assignment3/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

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

# this downloads the CIFAR-10 dataset to your Drive
# if it doesn't already exist.
%cd drive/My\ Drive/$FOLDERNAME/CV7062610/datasets/
!bash get_datasets.sh
%cd /content
```

```
Mounted at /content/drive
/content/drive/My Drive/CV7062610/assignments/assignment3/CV7062610/datasets
/content
```

1 What's this PyTorch business?

You've written a lot of code in this assignment to provide a whole host of neural network functionality. Dropout, Batch Norm, and 2D convolutions are some of the workhorses of deep learning in computer vision. You've also worked hard to make your code efficient and vectorized.

For the last part of this assignment, though, we're going to leave behind your beautiful codebase and instead migrate to one of two popular deep learning frameworks: in this instance, PyTorch (or TensorFlow, if you choose to use that notebook).

1.0.1 What is PyTorch?

PyTorch is a system for executing dynamic computational graphs over Tensor objects that behave similarly as numpy ndarray. It comes with a powerful automatic differentiation engine that removes the need for manual back-propagation.

1.0.2 Why?

- Our code will now run on GPUs! Much faster training. When using a framework like PyTorch or TensorFlow you can harness the power of the GPU for your own custom neural network architectures without having to write CUDA code directly (which is beyond the scope of this class).
- We want you to be ready to use one of these frameworks for your project so you can experiment more efficiently than if you were writing every feature you want to use by hand.
- We want you to stand on the shoulders of giants! TensorFlow and PyTorch are both excellent frameworks that will make your lives a lot easier, and now that you understand their guts, you are free to use them :)
- We want you to be exposed to the sort of deep learning code you might run into in academia or industry.

1.0.3 PyTorch versions

This notebook assumes that you are using **PyTorch version 1.4**. In some of the previous versions (e.g. before 0.4), Tensors had to be wrapped in Variable objects to be used in autograd; however Variables have now been deprecated. In addition 1.0+ versions separate a Tensor's datatype from its device, and use numpy-style factories for constructing Tensors rather than directly invoking Tensor constructors.

1.1 How will I learn PyTorch?

Justin Johnson has made an excellent [tutorial](#) for PyTorch.

You can also find the detailed [API doc](#) here. If you have other questions that are not addressed by the API docs, the [PyTorch forum](#) is a much better place to ask than StackOverflow.

1.2 Install PyTorch 1.4 (ONLY IF YOU ARE WORKING LOCALLY)

1. Have the latest version of Anaconda installed on your machine.
2. Create a new conda environment starting from Python 3.7. In this setup example, we'll call it `torch_env`.
3. Run the command: `conda activate torch_env`
4. Run the command: `pip install torch==1.4 torchvision==0.5.0`

2 Table of Contents

This assignment has 5 parts. You will learn PyTorch on **three different levels of abstraction**, which will help you understand it better and prepare you for the final project.

1. Part I, Preparation: we will use CIFAR-10 dataset.

2. Part II, Barebones PyTorch: **Abstraction level 1**, we will work directly with the lowest-level PyTorch Tensors.
3. Part III, PyTorch Module API: **Abstraction level 2**, we will use `nn.Module` to define arbitrary neural network architecture.
4. Part IV, PyTorch Sequential API: **Abstraction level 3**, we will use `nn.Sequential` to define a linear feed-forward network very conveniently.
5. Part V, CIFAR-10 open-ended challenge: please implement your own network to get as high accuracy as possible on CIFAR-10. You can experiment with any layer, optimizer, hyperparameters or other advanced features.

Here is a table of comparison:

API	Flexibility	Convenience
Barebone	High	Low
<code>nn.Module</code>	High	Medium
<code>nn.Sequential</code>	Low	High

3 Part I. Preparation

First, we load the CIFAR-10 dataset. This might take a couple minutes the first time you do it, but the files should stay cached after that.

In previous parts of the assignment we had to write our own code to download the CIFAR-10 dataset, preprocess it, and iterate through it in minibatches; PyTorch provides convenient tools to automate this process for us.

```
[2]: import torch
# print(torch.__version__.split('.')[2])
assert '.'.join(torch.__version__.split('.')[2]) == '2.0'
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader
from torch.utils.data import sampler

import torchvision.datasets as dset
import torchvision.transforms as T

import numpy as np
```

```
[3]: NUM_TRAIN = 49000

# The torchvision.transforms package provides tools for preprocessing data
# and for performing data augmentation; here we set up a transform to
# preprocess the data by subtracting the mean RGB value and dividing by the
# standard deviation of each RGB value; we've hardcoded the mean and std.
transform = T.Compose([
    T.ToTensor(),
```

```

        T.Normalize((0.4914, 0.4822, 0.4465), (0.2023, 0.1994, 0.2010))
    ])

# We set up a Dataset object for each split (train / val / test); Datasets load
# training examples one at a time, so we wrap each Dataset in a DataLoader which
# iterates through the Dataset and forms minibatches. We divide the CIFAR-10
# training set into train and val sets by passing a Sampler object to the
# DataLoader telling how it should sample from the underlying Dataset.
cifar10_train = dset.CIFAR10('./CV7062610/datasets', train=True, download=True,
                             transform=transform)
loader_train = DataLoader(cifar10_train, batch_size=64,
                          sampler=sampler.SubsetRandomSampler(range(NUM_TRAIN)))

cifar10_val = dset.CIFAR10('./CV7062610/datasets', train=True, download=True,
                           transform=transform)
loader_val = DataLoader(cifar10_val, batch_size=64,
                       sampler=sampler.SubsetRandomSampler(range(NUM_TRAIN,
↪50000))))

cifar10_test = dset.CIFAR10('./CV7062610/datasets', train=False, download=True,
                            transform=transform)
loader_test = DataLoader(cifar10_test, batch_size=64)

```

Downloading <https://www.cs.toronto.edu/~kriz/cifar-10-python.tar.gz> to
./CV7062610/datasets/cifar-10-python.tar.gz

100%| | 170498071/170498071 [00:01<00:00, 92934244.62it/s]

Extracting ./CV7062610/datasets/cifar-10-python.tar.gz to ./CV7062610/datasets

Files already downloaded and verified

Files already downloaded and verified

You have an option to **use GPU by setting the flag to True below**. It is not necessary to use GPU for this assignment. Note that if your computer does not have CUDA enabled, `torch.cuda.is_available()` will return False and this notebook will fallback to CPU mode.

The global variables `dtype` and `device` will control the data types throughout this assignment.

3.1 Colab Users

If you are using Colab, you need to manually switch to a GPU device. You can do this by clicking Runtime -> Change runtime type and selecting GPU under Hardware Accelerator. Note that you have to rerun the cells from the top since the kernel gets restarted upon switching runtimes.

```

[4]: USE_GPU = True

dtype = torch.float32 # we will be using float throughout this tutorial

if USE_GPU and torch.cuda.is_available():
    device = torch.device('cuda')

```

```

else:
    device = torch.device('cpu')

# Constant to control how frequently we print train loss
print_every = 100

print('using device:', device)

```

using device: cuda

4 Part II. Barebones PyTorch

PyTorch ships with high-level APIs to help us define model architectures conveniently, which we will cover in Part II of this tutorial. In this section, we will start with the barebone PyTorch elements to understand the autograd engine better. After this exercise, you will come to appreciate the high-level model API more.

We will start with a simple fully-connected ReLU network with two hidden layers and no biases for CIFAR classification. This implementation computes the forward pass using operations on PyTorch Tensors, and uses PyTorch autograd to compute gradients. It is important that you understand every line, because you will write a harder version after the example.

When we create a PyTorch Tensor with `requires_grad=True`, then operations involving that Tensor will not just compute values; they will also build up a computational graph in the background, allowing us to easily backpropagate through the graph to compute gradients of some Tensors with respect to a downstream loss. Concretely if `x` is a Tensor with `x.requires_grad == True` then after backpropagation `x.grad` will be another Tensor holding the gradient of `x` with respect to the scalar loss at the end.

4.0.1 PyTorch Tensors: Flatten Function

A PyTorch Tensor is conceptionally similar to a numpy array: it is an n -dimensional grid of numbers, and like numpy PyTorch provides many functions to efficiently operate on Tensors. As a simple example, we provide a `flatten` function below which reshapes image data for use in a fully-connected neural network.

Recall that image data is typically stored in a Tensor of shape $N \times C \times H \times W$, where:

- N is the number of datapoints
- C is the number of channels
- H is the height of the intermediate feature map in pixels
- W is the width of the intermediate feature map in pixels

This is the right way to represent the data when we are doing something like a 2D convolution, that needs spatial understanding of where the intermediate features are relative to each other. When we use fully connected affine layers to process the image, however, we want each datapoint to be represented by a single vector – it’s no longer useful to segregate the different channels, rows, and columns of the data. So, we use a “flatten” operation to collapse the $C \times H \times W$ values per representation into a single long vector. The `flatten` function below first reads in the N , C , H , and W values from a given batch of data, and then returns a “view” of that data. “View” is analogous

to numpy's "reshape" method: it reshapes x's dimensions to be N x ??, where ?? is allowed to be anything (in this case, it will be C x H x W, but we don't need to specify that explicitly).

```
[5]: def flatten(x):
      N = x.shape[0] # read in N, C, H, W
      return x.view(N, -1) # "flatten" the C * H * W values into a single vector
      ↪ per image

def test_flatten():
    x = torch.arange(12).view(2, 1, 3, 2)
    print('Before flattening: ', x)
    print('After flattening: ', flatten(x))

test_flatten()
```

```
Before flattening: tensor([[[[ 0,  1],
                             [ 2,  3],
                             [ 4,  5]]],

                        [[[ 6,  7],
                             [ 8,  9],
                             [10, 11]]]])

After flattening: tensor([[ 0,  1,  2,  3,  4,  5],
                          [ 6,  7,  8,  9, 10, 11]])
```

4.0.2 Barebones PyTorch: Two-Layer Network

Here we define a function `two_layer_fc` which performs the forward pass of a two-layer fully-connected ReLU network on a batch of image data. After defining the forward pass we check that it doesn't crash and that it produces outputs of the right shape by running zeros through the network.

You don't have to write any code here, but it's important that you read and understand the implementation.

```
[6]: import torch.nn.functional as F # useful stateless functions

def two_layer_fc(x, params):
    """
    A fully-connected neural networks; the architecture is:
    NN is fully connected -> ReLU -> fully connected layer.
    Note that this function only defines the forward pass;
    PyTorch will take care of the backward pass for us.

    The input to the network will be a minibatch of data, of shape
    (N, d1, ..., dM) where d1 * ... * dM = D. The hidden layer will have H
    ↪ units,
```

and the output layer will produce scores for C classes.

Inputs:

- x : A PyTorch Tensor of shape (N, d_1, \dots, d_M) giving a minibatch of input data.
- $params$: A list $[w_1, w_2]$ of PyTorch Tensors giving weights for the network; w_1 has shape (D, H) and w_2 has shape (H, C) .

Returns:

- $scores$: A PyTorch Tensor of shape (N, C) giving classification scores for the input data x .

"""

first we flatten the image

$x = \text{flatten}(x)$ *# shape: [batch_size, C x H x W]*

$w_1, w_2 = params$

Forward pass: compute predicted y using operations on Tensors. Since w_1 and

*# w_2 have `requires_grad=True`, operations involving these Tensors will cause
PyTorch to build a computational graph, allowing automatic computation of
gradients. Since we are no longer implementing the backward pass by hand, we*

don't need to keep references to intermediate values.

you can also use `.clamp(min=0)`, equivalent to `F.relu()`

$x = F.relu(x.mm(w_1))$

$x = x.mm(w_2)$

return x

def `two_layer_fc_test()`:

`hidden_layer_size = 42`

$x = \text{torch.zeros}((64, 50), \text{dtype}=\text{dtype})$ *# minibatch size 64, feature dimension 50*

$w_1 = \text{torch.zeros}(50, \text{hidden_layer_size}, \text{dtype}=\text{dtype})$

$w_2 = \text{torch.zeros}(\text{hidden_layer_size}, 10, \text{dtype}=\text{dtype})$

$scores = \text{two_layer_fc}(x, [w_1, w_2])$

print(`scores.size()`) *# you should see [64, 10]*

`two_layer_fc_test()`

`torch.Size([64, 10])`

4.0.3 Barebones PyTorch: Three-Layer ConvNet

Here you will complete the implementation of the function `three_layer_convnet`, which will perform the forward pass of a three-layer convolutional network. Like above, we can immediately test our implementation by passing zeros through the network. The network should have the following

architecture:

1. A convolutional layer (with bias) with `channel_1` filters, each with shape `KW1 x KH1`, and zero-padding of two
2. ReLU nonlinearity
3. A convolutional layer (with bias) with `channel_2` filters, each with shape `KW2 x KH2`, and zero-padding of one
4. ReLU nonlinearity
5. Fully-connected layer with bias, producing scores for `C` classes.

Note that we have **no softmax activation** here after our fully-connected layer: this is because PyTorch's cross entropy loss performs a softmax activation for you, and by bundling that step in makes computation more efficient.

HINT: For convolutions: <http://pytorch.org/docs/stable/nn.html#torch.nn.functional.conv2d>; pay attention to the shapes of convolutional filters!

```
[7]: def three_layer_convnet(x, params):  
    """  
    Performs the forward pass of a three-layer convolutional network with the  
    architecture defined above.  
  
    Inputs:  
    - x: A PyTorch Tensor of shape (N, 3, H, W) giving a minibatch of images  
    - params: A list of PyTorch Tensors giving the weights and biases for the  
      network; should contain the following:  
      - conv_w1: PyTorch Tensor of shape (channel_1, 3, KH1, KW1) giving weights  
        for the first convolutional layer  
      - conv_b1: PyTorch Tensor of shape (channel_1,) giving biases for the_  
↪first  
        convolutional layer  
      - conv_w2: PyTorch Tensor of shape (channel_2, channel_1, KH2, KW2) giving  
        weights for the second convolutional layer  
      - conv_b2: PyTorch Tensor of shape (channel_2,) giving biases for the_  
↪second  
        convolutional layer  
      - fc_w: PyTorch Tensor giving weights for the fully-connected layer. Can_  
↪you  
        figure out what the shape should be?  
      - fc_b: PyTorch Tensor giving biases for the fully-connected layer. Can_  
↪you  
        figure out what the shape should be?  
  
    Returns:  
    - scores: PyTorch Tensor of shape (N, C) giving classification scores for x  
    """  
    conv_w1, conv_b1, conv_w2, conv_b2, fc_w, fc_b = params  
    scores = None
```

```

    #####
    # TODO: Implement the forward pass for the three-layer ConvNet.
    #
    #####
    # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

    conv_layer1 = F.conv2d(x, weight=conv_w1, bias=conv_b1, padding=2)
    layer1_relu = F.relu(conv_layer1)
    conv_layer2 = F.conv2d(layer1_relu, weight=conv_w2, bias=conv_b2, padding=1)
    layer2_relu = F.relu(conv_layer2)
    flat = flatten(layer2_relu)
    scores = flat.mm(fc_w) + fc_b

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

```

After defining the forward pass of the ConvNet above, run the following cell to test your implementation.

When you run this function, scores should have shape (64, 10).

```

[8]: def three_layer_convnet_test():
    x = torch.zeros((64, 3, 32, 32), dtype=dtype) # minibatch size 64, image
    size [3, 32, 32]

    conv_w1 = torch.zeros((6, 3, 5, 5), dtype=dtype) # [out_channel,
    in_channel, kernel_H, kernel_W]
    conv_b1 = torch.zeros((6,)) # out_channel
    conv_w2 = torch.zeros((9, 6, 3, 3), dtype=dtype) # [out_channel,
    in_channel, kernel_H, kernel_W]
    conv_b2 = torch.zeros((9,)) # out_channel

    # you must calculate the shape of the tensor after two conv layers, before
    the fully-connected layer
    fc_w = torch.zeros((9 * 32 * 32, 10))
    fc_b = torch.zeros(10)

    scores = three_layer_convnet(x, [conv_w1, conv_b1, conv_w2, conv_b2, fc_w,
    fc_b])

```

```
print(scores.size()) # you should see [64, 10]
three_layer_convnet_test()
```

```
torch.Size([64, 10])
```

4.0.4 Barebones PyTorch: Initialization

Let's write a couple utility methods to initialize the weight matrices for our models.

- `random_weight(shape)` initializes a weight tensor with the Kaiming normalization method.
- `zero_weight(shape)` initializes a weight tensor with all zeros. Useful for instantiating bias parameters.

The `random_weight` function uses the Kaiming normal initialization method, described in:

He et al, *Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification*, ICCV 2015, <https://arxiv.org/abs/1502.01852>

```
[9]: def random_weight(shape):
    """
    Create random Tensors for weights; setting requires_grad=True means that we
    want to compute gradients for these Tensors during the backward pass.
    We use Kaiming normalization: sqrt(2 / fan_in)
    """
    if len(shape) == 2: # FC weight
        fan_in = shape[0]
    else:
        fan_in = np.prod(shape[1:]) # conv weight [out_channel, in_channel, kH,
        ↪ kW]
    # randn is standard normal distribution generator.
    w = torch.randn(shape, device=device, dtype=dtype) * np.sqrt(2. / fan_in)
    w.requires_grad = True
    return w

def zero_weight(shape):
    return torch.zeros(shape, device=device, dtype=dtype, requires_grad=True)

# create a weight of shape [3 x 5]
# you should see the type `torch.cuda.FloatTensor` if you use GPU.
# Otherwise it should be `torch.FloatTensor`
random_weight((3, 5))
```

```
[9]: tensor([[[-0.8248, -0.1313, -0.4047,  0.1306, -0.2706],
              [-0.2564,  0.2223, -1.2444,  0.8059,  0.2418],
              [-0.5468,  0.9315,  1.7524, -1.0418, -0.4623]], device='cuda:0',
              requires_grad=True)
```

4.0.5 Barebones PyTorch: Check Accuracy

When training the model we will use the following function to check the accuracy of our model on the training or validation sets.

When checking accuracy we don't need to compute any gradients; as a result we don't need PyTorch to build a computational graph for us when we compute scores. To prevent a graph from being built we scope our computation under a `torch.no_grad()` context manager.

```
[10]: def check_accuracy_part2(loader, model_fn, params):
    """
    Check the accuracy of a classification model.

    Inputs:
    - loader: A DataLoader for the data split we want to check
    - model_fn: A function that performs the forward pass of the model,
      with the signature scores = model_fn(x, params)
    - params: List of PyTorch Tensors giving parameters of the model

    Returns: Nothing, but prints the accuracy of the model
    """
    split = 'val' if loader.dataset.train else 'test'
    print('Checking accuracy on the %s set' % split)
    num_correct, num_samples = 0, 0
    with torch.no_grad():
        for x, y in loader:
            x = x.to(device=device, dtype=dtype) # move to device, e.g. GPU
            y = y.to(device=device, dtype=torch.int64)
            scores = model_fn(x, params)
            _, preds = scores.max(1)
            num_correct += (preds == y).sum()
            num_samples += preds.size(0)
        acc = float(num_correct) / num_samples
        print('Got %d / %d correct (%.2f%%)' % (num_correct, num_samples, 100 *
↪acc))
```

4.0.6 BareBones PyTorch: Training Loop

We can now set up a basic training loop to train our network. We will train the model using stochastic gradient descent without momentum. We will use `torch.functional.cross_entropy` to compute the loss; you can [read about it here](#).

The training loop takes as input the neural network function, a list of initialized parameters (`[w1, w2]` in our example), and learning rate.

```
[11]: def train_part2(model_fn, params, learning_rate):
    """
    Train a model on CIFAR-10.
```

Inputs:

- *model_fn*: A Python function that performs the forward pass of the model. It should have the signature `scores = model_fn(x, params)` where `x` is a PyTorch Tensor of image data, `params` is a list of PyTorch Tensors giving model weights, and `scores` is a PyTorch Tensor of shape (N, C) giving scores for the elements in `x`.
- *params*: List of PyTorch Tensors giving weights for the model
- *learning_rate*: Python scalar giving the learning rate to use for SGD

Returns: Nothing

```
"""
for t, (x, y) in enumerate(loader_train):
    # Move the data to the proper device (GPU or CPU)
    x = x.to(device=device, dtype=dtype)
    y = y.to(device=device, dtype=torch.long)

    # Forward pass: compute scores and loss
    scores = model_fn(x, params)
    loss = F.cross_entropy(scores, y)

    # Backward pass: PyTorch figures out which Tensors in the computational
    # graph has requires_grad=True and uses backpropagation to compute the
    # gradient of the loss with respect to these Tensors, and stores the
    # gradients in the .grad attribute of each Tensor.
    loss.backward()

    # Update parameters. We don't want to backpropagate through the
    # parameter updates, so we scope the updates under a torch.no_grad()
    # context manager to prevent a computational graph from being built.
    with torch.no_grad():
        for w in params:
            w -= learning_rate * w.grad

            # Manually zero the gradients after running the backward pass
            w.grad.zero_()

    if t % print_every == 0:
        print('Iteration %d, loss = %.4f' % (t, loss.item()))
        check_accuracy_part2(loader_val, model_fn, params)
        print()
```

4.0.7 BareBones PyTorch: Train a Two-Layer Network

Now we are ready to run the training loop. We need to explicitly allocate tensors for the fully connected weights, `w1` and `w2`.

Each minibatch of CIFAR has 64 examples, so the tensor shape is `[64, 3, 32, 32]`.

After flattening, `x` shape should be `[64, 3 * 32 * 32]`. This will be the size of the first dimension of `w1`. The second dimension of `w1` is the hidden layer size, which will also be the first dimension of `w2`.

Finally, the output of the network is a 10-dimensional vector that represents the probability distribution over 10 classes.

You don't need to tune any hyperparameters but you should see accuracies above 40% after training for one epoch.

```
[12]: hidden_layer_size = 4000
      learning_rate = 1e-2

      w1 = random_weight((3 * 32 * 32, hidden_layer_size))
      w2 = random_weight((hidden_layer_size, 10))

      train_part2(two_layer_fc, [w1, w2], learning_rate)
```

```
Iteration 0, loss = 3.5647
Checking accuracy on the val set
Got 183 / 1000 correct (18.30%)
```

```
Iteration 100, loss = 2.3661
Checking accuracy on the val set
Got 275 / 1000 correct (27.50%)
```

```
Iteration 200, loss = 1.6487
Checking accuracy on the val set
Got 400 / 1000 correct (40.00%)
```

```
Iteration 300, loss = 1.7627
Checking accuracy on the val set
Got 416 / 1000 correct (41.60%)
```

```
Iteration 400, loss = 1.6810
Checking accuracy on the val set
Got 420 / 1000 correct (42.00%)
```

```
Iteration 500, loss = 1.9586
Checking accuracy on the val set
Got 389 / 1000 correct (38.90%)
```

```
Iteration 600, loss = 2.1049
Checking accuracy on the val set
Got 395 / 1000 correct (39.50%)
```

```
Iteration 700, loss = 1.6632
Checking accuracy on the val set
Got 442 / 1000 correct (44.20%)
```

4.0.8 BareBones PyTorch: Training a ConvNet

In the below you should use the functions defined above to train a three-layer convolutional network on CIFAR. The network should have the following architecture:

1. Convolutional layer (with bias) with 32 5x5 filters, with zero-padding of 2
2. ReLU
3. Convolutional layer (with bias) with 16 3x3 filters, with zero-padding of 1
4. ReLU
5. Fully-connected layer (with bias) to compute scores for 10 classes

You should initialize your weight matrices using the `random_weight` function defined above, and you should initialize your bias vectors using the `zero_weight` function above.

You don't need to tune any hyperparameters, but if everything works correctly you should achieve an accuracy above 42% after one epoch.

```
[13]: learning_rate = 3e-3

channel_1 = 32
channel_2 = 16

conv_w1 = None
conv_b1 = None
conv_w2 = None
conv_b2 = None
fc_w = None
fc_b = None

#####
# TODO: Initialize the parameters of a three-layer ConvNet.                                #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

conv_w1 = random_weight((channel_1, 3, 5, 5))
conv_b1 = zero_weight((channel_1,))
conv_w2 = random_weight((channel_2, 32, 3, 3))
conv_b2 = zero_weight((channel_2,))
fc_w = random_weight((channel_2*32*32, 10))
fc_b = zero_weight((10,))

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

params = [conv_w1, conv_b1, conv_w2, conv_b2, fc_w, fc_b]
```

```
train_part2(three_layer_convnet, params, learning_rate)
```

```
Iteration 0, loss = 2.9809
Checking accuracy on the val set
Got 123 / 1000 correct (12.30%)
```

```
Iteration 100, loss = 1.7261
Checking accuracy on the val set
Got 352 / 1000 correct (35.20%)
```

```
Iteration 200, loss = 1.8249
Checking accuracy on the val set
Got 407 / 1000 correct (40.70%)
```

```
Iteration 300, loss = 1.5912
Checking accuracy on the val set
Got 435 / 1000 correct (43.50%)
```

```
Iteration 400, loss = 1.6609
Checking accuracy on the val set
Got 448 / 1000 correct (44.80%)
```

```
Iteration 500, loss = 1.3858
Checking accuracy on the val set
Got 452 / 1000 correct (45.20%)
```

```
Iteration 600, loss = 1.5697
Checking accuracy on the val set
Got 478 / 1000 correct (47.80%)
```

```
Iteration 700, loss = 1.5425
Checking accuracy on the val set
Got 480 / 1000 correct (48.00%)
```

5 Part III. PyTorch Module API

Barebone PyTorch requires that we track all the parameter tensors by hand. This is fine for small networks with a few tensors, but it would be extremely inconvenient and error-prone to track tens or hundreds of tensors in larger networks.

PyTorch provides the `nn.Module` API for you to define arbitrary network architectures, while tracking every learnable parameters for you. In Part II, we implemented SGD ourselves. PyTorch also provides the `torch.optim` package that implements all the common optimizers, such as RMSProp, Adagrad, and Adam. It even supports approximate second-order methods like L-BFGS! You can refer to the [doc](#) for the exact specifications of each optimizer.

To use the Module API, follow the steps below:

1. Subclass `nn.Module`. Give your network class an intuitive name like `TwoLayerFC`.
2. In the constructor `__init__()`, define all the layers you need as class attributes. Layer objects like `nn.Linear` and `nn.Conv2d` are themselves `nn.Module` subclasses and contain learnable parameters, so that you don't have to instantiate the raw tensors yourself. `nn.Module` will track these internal parameters for you. Refer to the [doc](#) to learn more about the dozens of builtin layers. **Warning:** don't forget to call the `super().__init__()` first!
3. In the `forward()` method, define the *connectivity* of your network. You should use the attributes defined in `__init__` as function calls that take tensor as input and output the "transformed" tensor. Do *not* create any new layers with learnable parameters in `forward()`! All of them must be declared upfront in `__init__`.

After you define your Module subclass, you can instantiate it as an object and call it just like the NN forward function in part II.

5.0.1 Module API: Two-Layer Network

Here is a concrete example of a 2-layer fully connected network:

```
[14]: class TwoLayerFC(nn.Module):
    def __init__(self, input_size, hidden_size, num_classes):
        super().__init__()
        # assign layer objects to class attributes
        self.fc1 = nn.Linear(input_size, hidden_size)
        # nn.init package contains convenient initialization methods
        # http://pytorch.org/docs/master/nn.html#torch-nn-init
        nn.init.kaiming_normal_(self.fc1.weight)
        self.fc2 = nn.Linear(hidden_size, num_classes)
        nn.init.kaiming_normal_(self.fc2.weight)

    def forward(self, x):
        # forward always defines connectivity
        x = flatten(x)
        scores = self.fc2(F.relu(self.fc1(x)))
        return scores

def test_TwoLayerFC():
    input_size = 50
    x = torch.zeros((64, input_size), dtype=dtype) # minibatch size 64,
    ↪ feature dimension 50
    model = TwoLayerFC(input_size, 42, 10)
    scores = model(x)
    print(scores.size()) # you should see [64, 10]
test_TwoLayerFC()
```

```
torch.Size([64, 10])
```

5.0.2 Module API: Three-Layer ConvNet

It's your turn to implement a 3-layer ConvNet followed by a fully connected layer. The network architecture should be the same as in Part II:

1. Convolutional layer with `channel_1` 5x5 filters with zero-padding of 2
2. ReLU
3. Convolutional layer with `channel_2` 3x3 filters with zero-padding of 1
4. ReLU
5. Fully-connected layer to `num_classes` classes

You should initialize the weight matrices of the model using the Kaiming normal initialization method.

HINT: <http://pytorch.org/docs/stable/nn.html#conv2d>

After you implement the three-layer ConvNet, the `test_ThreeLayerConvNet` function will run your implementation; it should print (64, 10) for the shape of the output scores.

```
[15]: class ThreeLayerConvNet(nn.Module):
    def __init__(self, in_channel, channel_1, channel_2, num_classes):
        super().__init__()
        #####
        # TODO: Set up the layers you need for a three-layer ConvNet with the #
        # architecture defined above.                                         #
        #####
        # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

        self.conv1 = nn.Conv2d(in_channel, channel_1, kernel_size=5, padding=2)
        self.conv2 = nn.Conv2d(channel_1, channel_2, kernel_size=3, padding=1)
        self.fc = nn.Linear(channel_2 * 32 * 32, num_classes)

        # initialize weights
        nn.init.kaiming_normal_(self.conv1.weight)
        nn.init.zeros_(self.conv1.bias)
        nn.init.kaiming_normal_(self.conv2.weight)
        nn.init.zeros_(self.conv2.bias)
        nn.init.kaiming_normal_(self.fc.weight)
        nn.init.zeros_(self.fc.bias)

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

    def forward(self, x):
        scores = None
        #####
        # TODO: Implement the forward function for a 3-layer ConvNet. you      #
```

```

        # should use the layers you defined in __init__ and specify the      #
        # connectivity of those layers in forward()                          #
        #####
        # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

        x = F.relu(self.conv1(x))
        x = F.relu(self.conv2(x))
        scores = self.fc(flatten(x))

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

def test_ThreeLayerConvNet():
    x = torch.zeros((64, 3, 32, 32), dtype=dtype) # minibatch size 64, image
    ↪size [3, 32, 32]
    model = ThreeLayerConvNet(in_channel=3, channel_1=12, channel_2=8,
    ↪num_classes=10)
    scores = model(x)
    print(scores.size()) # you should see [64, 10]
test_ThreeLayerConvNet()

```

```
torch.Size([64, 10])
```

5.0.3 Module API: Check Accuracy

Given the validation or test set, we can check the classification accuracy of a neural network.

This version is slightly different from the one in part II. You don't manually pass in the parameters anymore.

```

[16]: def check_accuracy_part34(loader, model):
        if loader.dataset.train:
            print('Checking accuracy on validation set')
        else:
            print('Checking accuracy on test set')
        num_correct = 0
        num_samples = 0
        model.eval() # set model to evaluation mode
        with torch.no_grad():
            for x, y in loader:
                x = x.to(device=device, dtype=dtype) # move to device, e.g. GPU
                y = y.to(device=device, dtype=torch.long)
                scores = model(x)
                _, preds = scores.max(1)

```

```

        num_correct += (preds == y).sum()
        num_samples += preds.size(0)
    acc = float(num_correct) / num_samples
    print('Got %d / %d correct (%.2f)' % (num_correct, num_samples, 100 *
↪acc))

```

5.0.4 Module API: Training Loop

We also use a slightly different training loop. Rather than updating the values of the weights ourselves, we use an Optimizer object from the `torch.optim` package, which abstract the notion of an optimization algorithm and provides implementations of most of the algorithms commonly used to optimize neural networks.

```

[17]: def train_part34(model, optimizer, epochs=1):
    """
    Train a model on CIFAR-10 using the PyTorch Module API.

    Inputs:
    - model: A PyTorch Module giving the model to train.
    - optimizer: An Optimizer object we will use to train the model
    - epochs: (Optional) A Python integer giving the number of epochs to train
↪for

    Returns: Nothing, but prints model accuracies during training.
    """
    model = model.to(device=device) # move the model parameters to CPU/GPU
    for e in range(epochs):
        for t, (x, y) in enumerate(loader_train):
            model.train() # put model to training mode
            x = x.to(device=device, dtype=dtype) # move to device, e.g. GPU
            y = y.to(device=device, dtype=torch.long)

            scores = model(x)
            loss = F.cross_entropy(scores, y)

            # Zero out all of the gradients for the variables which the
↪optimizer
            # will update.
            optimizer.zero_grad()

            # This is the backwards pass: compute the gradient of the loss with
            # respect to each parameter of the model.
            loss.backward()

            # Actually update the parameters of the model using the gradients
            # computed by the backwards pass.
            optimizer.step()

```

```

if t % print_every == 0:
    print('Iteration %d, loss = %.4f' % (t, loss.item()))
    check_accuracy_part34(loader_val, model)
    print()

```

5.0.5 Module API: Train a Two-Layer Network

Now we are ready to run the training loop. In contrast to part II, we don't explicitly allocate parameter tensors anymore.

Simply pass the input size, hidden layer size, and number of classes (i.e. output size) to the constructor of `TwoLayerFC`.

You also need to define an optimizer that tracks all the learnable parameters inside `TwoLayerFC`.

You don't need to tune any hyperparameters, but you should see model accuracies above 40% after training for one epoch.

```

[18]: hidden_layer_size = 4000
      learning_rate = 1e-2
      model = TwoLayerFC(3 * 32 * 32, hidden_layer_size, 10)
      optimizer = optim.SGD(model.parameters(), lr=learning_rate)

      train_part34(model, optimizer)

```

```

Iteration 0, loss = 3.1123
Checking accuracy on validation set
Got 138 / 1000 correct (13.80)

```

```

Iteration 100, loss = 2.7970
Checking accuracy on validation set
Got 275 / 1000 correct (27.50)

```

```

Iteration 200, loss = 2.2086
Checking accuracy on validation set
Got 352 / 1000 correct (35.20)

```

```

Iteration 300, loss = 1.8768
Checking accuracy on validation set
Got 388 / 1000 correct (38.80)

```

```

Iteration 400, loss = 1.9888
Checking accuracy on validation set
Got 417 / 1000 correct (41.70)

```

```

Iteration 500, loss = 1.8703
Checking accuracy on validation set
Got 388 / 1000 correct (38.80)

```

```
Iteration 600, loss = 1.5631
Checking accuracy on validation set
Got 435 / 1000 correct (43.50)
```

```
Iteration 700, loss = 1.7312
Checking accuracy on validation set
Got 431 / 1000 correct (43.10)
```

5.0.6 Module API: Train a Three-Layer ConvNet

You should now use the Module API to train a three-layer ConvNet on CIFAR. This should look very similar to training the two-layer network! You don't need to tune any hyperparameters, but you should achieve above 45% after training for one epoch.

You should train the model using stochastic gradient descent without momentum.

```
[19]: learning_rate = 3e-3
channel_1 = 32
channel_2 = 16

model = None
optimizer = None
#####
# TODO: Instantiate your ThreeLayerConvNet model and a corresponding optimizer #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

model = ThreeLayerConvNet(in_channel=3, channel_1=channel_1,
    ↪channel_2=channel_2, num_classes=10)
optimizer = torch.optim.SGD(model.parameters(), lr=learning_rate)

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

train_part34(model, optimizer)
```

```
Iteration 0, loss = 2.8772
Checking accuracy on validation set
Got 117 / 1000 correct (11.70)
```

```
Iteration 100, loss = 1.9511
Checking accuracy on validation set
Got 373 / 1000 correct (37.30)
```

```
Iteration 200, loss = 1.9240
```

```
Checking accuracy on validation set
Got 422 / 1000 correct (42.20)
```

```
Iteration 300, loss = 1.4985
Checking accuracy on validation set
Got 437 / 1000 correct (43.70)
```

```
Iteration 400, loss = 1.3987
Checking accuracy on validation set
Got 449 / 1000 correct (44.90)
```

```
Iteration 500, loss = 1.6882
Checking accuracy on validation set
Got 479 / 1000 correct (47.90)
```

```
Iteration 600, loss = 1.6165
Checking accuracy on validation set
Got 468 / 1000 correct (46.80)
```

```
Iteration 700, loss = 1.6246
Checking accuracy on validation set
Got 482 / 1000 correct (48.20)
```

6 Part IV. PyTorch Sequential API

Part III introduced the PyTorch Module API, which allows you to define arbitrary learnable layers and their connectivity.

For simple models like a stack of feed forward layers, you still need to go through 3 steps: subclass `nn.Module`, assign layers to class attributes in `__init__`, and call each layer one by one in `forward()`. Is there a more convenient way?

Fortunately, PyTorch provides a container Module called `nn.Sequential`, which merges the above steps into one. It is not as flexible as `nn.Module`, because you cannot specify more complex topology than a feed-forward stack, but it's good enough for many use cases.

6.0.1 Sequential API: Two-Layer Network

Let's see how to rewrite our two-layer fully connected network example with `nn.Sequential`, and train it using the training loop defined above.

Again, you don't need to tune any hyperparameters here, but you should achieve above 40% accuracy after one epoch of training.

```
[20]: # We need to wrap `flatten` function in a module in order to stack it
      # in nn.Sequential
      class Flatten(nn.Module):
          def forward(self, x):
```

```

        return flatten(x)

hidden_layer_size = 4000
learning_rate = 1e-2

model = nn.Sequential(
    Flatten(),
    nn.Linear(3 * 32 * 32, hidden_layer_size),
    nn.ReLU(),
    nn.Linear(hidden_layer_size, 10),
)

# you can use Nesterov momentum in optim.SGD
optimizer = optim.SGD(model.parameters(), lr=learning_rate,
                       momentum=0.9, nesterov=True)

train_part34(model, optimizer)

```

Iteration 0, loss = 2.4199
 Checking accuracy on validation set
 Got 136 / 1000 correct (13.60)

Iteration 100, loss = 1.6227
 Checking accuracy on validation set
 Got 413 / 1000 correct (41.30)

Iteration 200, loss = 1.6452
 Checking accuracy on validation set
 Got 411 / 1000 correct (41.10)

Iteration 300, loss = 1.8250
 Checking accuracy on validation set
 Got 398 / 1000 correct (39.80)

Iteration 400, loss = 1.9432
 Checking accuracy on validation set
 Got 421 / 1000 correct (42.10)

Iteration 500, loss = 1.7115
 Checking accuracy on validation set
 Got 437 / 1000 correct (43.70)

Iteration 600, loss = 1.8687
 Checking accuracy on validation set
 Got 417 / 1000 correct (41.70)

Iteration 700, loss = 1.6250

Checking accuracy on validation set
Got 458 / 1000 correct (45.80)

6.0.2 Sequential API: Three-Layer ConvNet

Here you should use `nn.Sequential` to define and train a three-layer ConvNet with the same architecture we used in Part III:

1. Convolutional layer (with bias) with 32 5x5 filters, with zero-padding of 2
2. ReLU
3. Convolutional layer (with bias) with 16 3x3 filters, with zero-padding of 1
4. ReLU
5. Fully-connected layer (with bias) to compute scores for 10 classes

You should initialize your weight matrices using the `random_weight` function defined above, and you should initialize your bias vectors using the `zero_weight` function above.

You should optimize your model using stochastic gradient descent with Nesterov momentum 0.9.

Again, you don't need to tune any hyperparameters but you should see accuracy above 55% after one epoch of training.

```
[35]: channel_1 = 32
channel_2 = 16
learning_rate = 1e-2

model = None
optimizer = None

#####
# TODO: Rewrite the 2-layer ConvNet with bias from Part III with the Sequential API.
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

model = nn.Sequential(
    nn.Conv2d(3, channel_1, kernel_size=5, padding=2),
    nn.ReLU(),
    nn.Conv2d(channel_1, channel_2, kernel_size=3, padding=1),
    nn.ReLU(),
    nn.Flatten(),
    nn.Linear(channel_2*32*32, 10)
)

# Initialize weights using random_weight and biases using zero_weight
model.weight = random_weight((channel_2*32*32, 10))
model.bias = zero_weight((10,))
```

```
optimizer = optim.SGD(model.parameters(), lr = learning_rate, momentum = 0.9,
↳nesterov = True)

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

train_part34(model, optimizer)
```

Iteration 0, loss = 2.3173
Checking accuracy on validation set
Got 114 / 1000 correct (11.40)

Iteration 100, loss = 1.7750
Checking accuracy on validation set
Got 432 / 1000 correct (43.20)

Iteration 200, loss = 1.4821
Checking accuracy on validation set
Got 437 / 1000 correct (43.70)

Iteration 300, loss = 1.2806
Checking accuracy on validation set
Got 492 / 1000 correct (49.20)

Iteration 400, loss = 1.3001
Checking accuracy on validation set
Got 541 / 1000 correct (54.10)

Iteration 500, loss = 1.4178
Checking accuracy on validation set
Got 533 / 1000 correct (53.30)

Iteration 600, loss = 1.3421
Checking accuracy on validation set
Got 560 / 1000 correct (56.00)

Iteration 700, loss = 1.3538
Checking accuracy on validation set
Got 566 / 1000 correct (56.60)

7 Part V. CIFAR-10 open-ended challenge

In this section, you can experiment with whatever ConvNet architecture you'd like on CIFAR-10.

Now it's your job to experiment with architectures, hyperparameters, loss functions, and optimizers to train a model that achieves **at least 70%** accuracy on the CIFAR-10 **validation** set within 10 epochs. You can use the `check_accuracy` and `train` functions from above. You can use either `nn.Module` or `nn.Sequential` API.

Describe what you did at the end of this notebook.

Here are the official API documentation for each component. One note: what we call in the class “spatial batch norm” is called “BatchNorm2D” in PyTorch.

- Layers in torch.nn package: <http://pytorch.org/docs/stable/nn.html>
- Activations: <http://pytorch.org/docs/stable/nn.html#non-linear-activations>
- Loss functions: <http://pytorch.org/docs/stable/nn.html#loss-functions>
- Optimizers: <http://pytorch.org/docs/stable/optim.html>

7.0.1 Things you might try:

- **Filter size:** Above we used 5x5; would smaller filters be more efficient?
- **Number of filters:** Above we used 32 filters. Do more or fewer do better?
- **Pooling vs Strided Convolution:** Do you use max pooling or just stride convolutions?
- **Batch normalization:** Try adding spatial batch normalization after convolution layers and vanilla batch normalization after affine layers. Do your networks train faster?
- **Network architecture:** The network above has two layers of trainable parameters. Can you do better with a deep network? Good architectures to try include:
 - [conv-relu-pool]xN -> [affine]xM -> [softmax or SVM]
 - [conv-relu-conv-relu-pool]xN -> [affine]xM -> [softmax or SVM]
 - [batchnorm-relu-conv]xN -> [affine]xM -> [softmax or SVM]
- **Global Average Pooling:** Instead of flattening and then having multiple affine layers, perform convolutions until your image gets small (7x7 or so) and then perform an average pooling operation to get to a 1x1 image picture (1, 1, Filter#), which is then reshaped into a (Filter#) vector. This is used in [Google's Inception Network](#) (See Table 1 for their architecture).
- **Regularization:** Add l2 weight regularization, or perhaps use Dropout.

7.0.2 Tips for training

For each network architecture that you try, you should tune the learning rate and other hyperparameters. When doing this there are a couple important things to keep in mind:

- If the parameters are working well, you should see improvement within a few hundred iterations
- Remember the coarse-to-fine approach for hyperparameter tuning: start by testing a large range of hyperparameters for just a few training iterations to find the combinations of parameters that are working at all.
- Once you have found some sets of parameters that seem to work, search more finely around these parameters. You may need to train for more epochs.
- You should use the validation set for hyperparameter search, and save your test set for evaluating your architecture on the best parameters as selected by the validation set.

7.0.3 Going above and beyond

If you are feeling adventurous there are many other features you can implement to try and improve your performance. You are **not required** to implement any of these, but don't miss the fun if you have time!

- Alternative optimizers: you can try Adam, Adagrad, RMSprop, etc.
- Alternative activation functions such as leaky ReLU, parametric ReLU, ELU, or MaxOut.
- Model ensembles
- Data augmentation
- New Architectures
 - [ResNets](#) where the input from the previous layer is added to the output.
 - [DenseNets](#) where inputs into previous layers are concatenated together.
 - [This blog has an in-depth overview](#)

7.0.4 Have fun and happy training!

```
[39]: #####
# TODO:
# ↪#
# Experiment with any architectures, optimizers, and hyperparameters.
# Achieve AT LEAST 70% accuracy on the *validation set* within 10 epochs.
#
# Note that you can use the check_accuracy function to evaluate on either
# the test set or the validation set, by passing either loader_test or
# loader_val as the second argument to check_accuracy. You should not touch
# the test set until you have finished your architecture and hyperparameter
# tuning, and only run the test set once at the end to report a final value.
#####
model = None
optimizer = None

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

out_channels_1, out_channels_2, out_channels_3, num_classes = 16, 32, 64, 10
in_channels_1, in_channels_2, in_channels_3 = 3, 16, 32
filter_size_1, filter_size_2, filter_size_3 = 5, 3, 3

# Define the model architecture
model = nn.Sequential(
    nn.Conv2d(in_channels_1, out_channels_1, kernel_size=filter_size_1,
    ↪padding=2),
    nn.BatchNorm2d(out_channels_1),
    nn.ReLU(),
    nn.MaxPool2d(2),

    nn.Conv2d(in_channels_2, out_channels_2, kernel_size=filter_size_2,
    ↪padding=1),
```

```

        nn.BatchNorm2d(out_channels_2),
        nn.ReLU(),
        nn.MaxPool2d(2),

        nn.Conv2d(in_channels_3, out_channels_3, kernel_size=filter_size_3,
padding=1),
        nn.BatchNorm2d(out_channels_3),
        nn.ReLU(),
        nn.MaxPool2d(2),

        Flatten(),
        nn.Dropout(0.2, inplace=True),
        nn.Linear(out_channels_3 * 4 * 4, num_classes)
    )

optimizer = optim.Adam(model.parameters(), lr=1e-3)

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

# You should get at least 70% accuracy
train_part34(model, optimizer, epochs=10)

```

Iteration 0, loss = 2.4397
Checking accuracy on validation set
Got 101 / 1000 correct (10.10)

Iteration 100, loss = 1.7240
Checking accuracy on validation set
Got 455 / 1000 correct (45.50)

Iteration 200, loss = 1.2997
Checking accuracy on validation set
Got 520 / 1000 correct (52.00)

Iteration 300, loss = 1.3064
Checking accuracy on validation set
Got 514 / 1000 correct (51.40)

Iteration 400, loss = 1.3229
Checking accuracy on validation set
Got 562 / 1000 correct (56.20)

Iteration 500, loss = 1.2340
Checking accuracy on validation set
Got 584 / 1000 correct (58.40)

Iteration 600, loss = 1.3612
Checking accuracy on validation set
Got 566 / 1000 correct (56.60)

Iteration 700, loss = 1.1360
Checking accuracy on validation set
Got 636 / 1000 correct (63.60)

Iteration 0, loss = 1.0861
Checking accuracy on validation set
Got 628 / 1000 correct (62.80)

Iteration 100, loss = 1.0962
Checking accuracy on validation set
Got 616 / 1000 correct (61.60)

Iteration 200, loss = 0.9910
Checking accuracy on validation set
Got 612 / 1000 correct (61.20)

Iteration 300, loss = 0.9068
Checking accuracy on validation set
Got 660 / 1000 correct (66.00)

Iteration 400, loss = 1.3202
Checking accuracy on validation set
Got 658 / 1000 correct (65.80)

Iteration 500, loss = 0.8837
Checking accuracy on validation set
Got 650 / 1000 correct (65.00)

Iteration 600, loss = 0.8091
Checking accuracy on validation set
Got 663 / 1000 correct (66.30)

Iteration 700, loss = 0.9584
Checking accuracy on validation set
Got 689 / 1000 correct (68.90)

Iteration 0, loss = 0.8416
Checking accuracy on validation set
Got 664 / 1000 correct (66.40)

Iteration 100, loss = 0.6066
Checking accuracy on validation set
Got 692 / 1000 correct (69.20)

Iteration 200, loss = 0.7117
Checking accuracy on validation set
Got 692 / 1000 correct (69.20)

Iteration 300, loss = 0.9997
Checking accuracy on validation set
Got 706 / 1000 correct (70.60)

Iteration 400, loss = 0.8614
Checking accuracy on validation set
Got 699 / 1000 correct (69.90)

Iteration 500, loss = 0.9000
Checking accuracy on validation set
Got 709 / 1000 correct (70.90)

Iteration 600, loss = 1.0176
Checking accuracy on validation set
Got 696 / 1000 correct (69.60)

Iteration 700, loss = 0.8563
Checking accuracy on validation set
Got 707 / 1000 correct (70.70)

Iteration 0, loss = 0.8915
Checking accuracy on validation set
Got 705 / 1000 correct (70.50)

Iteration 100, loss = 0.7452
Checking accuracy on validation set
Got 684 / 1000 correct (68.40)

Iteration 200, loss = 0.7224
Checking accuracy on validation set
Got 723 / 1000 correct (72.30)

Iteration 300, loss = 0.5925
Checking accuracy on validation set
Got 712 / 1000 correct (71.20)

Iteration 400, loss = 0.9086
Checking accuracy on validation set
Got 705 / 1000 correct (70.50)

Iteration 500, loss = 0.9365
Checking accuracy on validation set
Got 712 / 1000 correct (71.20)

Iteration 600, loss = 0.5817
Checking accuracy on validation set
Got 697 / 1000 correct (69.70)

Iteration 700, loss = 0.6105
Checking accuracy on validation set
Got 718 / 1000 correct (71.80)

Iteration 0, loss = 0.8264
Checking accuracy on validation set
Got 727 / 1000 correct (72.70)

Iteration 100, loss = 0.8810
Checking accuracy on validation set
Got 737 / 1000 correct (73.70)

Iteration 200, loss = 0.7194
Checking accuracy on validation set
Got 733 / 1000 correct (73.30)

Iteration 300, loss = 0.6493
Checking accuracy on validation set
Got 723 / 1000 correct (72.30)

Iteration 400, loss = 0.6984
Checking accuracy on validation set
Got 712 / 1000 correct (71.20)

Iteration 500, loss = 0.6289
Checking accuracy on validation set
Got 736 / 1000 correct (73.60)

Iteration 600, loss = 0.8370
Checking accuracy on validation set
Got 751 / 1000 correct (75.10)

Iteration 700, loss = 0.6240
Checking accuracy on validation set
Got 732 / 1000 correct (73.20)

Iteration 0, loss = 0.9003
Checking accuracy on validation set
Got 727 / 1000 correct (72.70)

Iteration 100, loss = 0.6298
Checking accuracy on validation set
Got 730 / 1000 correct (73.00)

Iteration 200, loss = 0.7169
Checking accuracy on validation set
Got 739 / 1000 correct (73.90)

Iteration 300, loss = 0.7580
Checking accuracy on validation set
Got 725 / 1000 correct (72.50)

Iteration 400, loss = 0.6466
Checking accuracy on validation set
Got 715 / 1000 correct (71.50)

Iteration 500, loss = 0.7211
Checking accuracy on validation set
Got 739 / 1000 correct (73.90)

Iteration 600, loss = 0.6684
Checking accuracy on validation set
Got 744 / 1000 correct (74.40)

Iteration 700, loss = 0.7333
Checking accuracy on validation set
Got 727 / 1000 correct (72.70)

Iteration 0, loss = 0.7814
Checking accuracy on validation set
Got 726 / 1000 correct (72.60)

Iteration 100, loss = 0.9284
Checking accuracy on validation set
Got 746 / 1000 correct (74.60)

Iteration 200, loss = 0.7801
Checking accuracy on validation set
Got 749 / 1000 correct (74.90)

Iteration 300, loss = 0.6281
Checking accuracy on validation set
Got 749 / 1000 correct (74.90)

Iteration 400, loss = 0.7621
Checking accuracy on validation set
Got 747 / 1000 correct (74.70)

Iteration 500, loss = 0.5694
Checking accuracy on validation set
Got 738 / 1000 correct (73.80)

Iteration 600, loss = 0.7617
Checking accuracy on validation set
Got 742 / 1000 correct (74.20)

Iteration 700, loss = 0.7810
Checking accuracy on validation set
Got 746 / 1000 correct (74.60)

Iteration 0, loss = 0.6197
Checking accuracy on validation set
Got 755 / 1000 correct (75.50)

Iteration 100, loss = 0.6291
Checking accuracy on validation set
Got 749 / 1000 correct (74.90)

Iteration 200, loss = 0.8020
Checking accuracy on validation set
Got 757 / 1000 correct (75.70)

Iteration 300, loss = 0.5861
Checking accuracy on validation set
Got 738 / 1000 correct (73.80)

Iteration 400, loss = 0.7545
Checking accuracy on validation set
Got 742 / 1000 correct (74.20)

Iteration 500, loss = 0.6857
Checking accuracy on validation set
Got 744 / 1000 correct (74.40)

Iteration 600, loss = 0.5311
Checking accuracy on validation set
Got 736 / 1000 correct (73.60)

Iteration 700, loss = 1.0659
Checking accuracy on validation set
Got 748 / 1000 correct (74.80)

Iteration 0, loss = 0.8196
Checking accuracy on validation set
Got 737 / 1000 correct (73.70)

Iteration 100, loss = 0.7805
Checking accuracy on validation set
Got 750 / 1000 correct (75.00)

Iteration 200, loss = 0.6436
Checking accuracy on validation set
Got 761 / 1000 correct (76.10)

Iteration 300, loss = 0.6205
Checking accuracy on validation set
Got 733 / 1000 correct (73.30)

Iteration 400, loss = 0.5842
Checking accuracy on validation set
Got 731 / 1000 correct (73.10)

Iteration 500, loss = 0.6374
Checking accuracy on validation set
Got 750 / 1000 correct (75.00)

Iteration 600, loss = 0.6735
Checking accuracy on validation set
Got 752 / 1000 correct (75.20)

Iteration 700, loss = 0.6306
Checking accuracy on validation set
Got 745 / 1000 correct (74.50)

Iteration 0, loss = 0.6283
Checking accuracy on validation set
Got 769 / 1000 correct (76.90)

Iteration 100, loss = 0.9043
Checking accuracy on validation set
Got 735 / 1000 correct (73.50)

Iteration 200, loss = 0.8386
Checking accuracy on validation set
Got 757 / 1000 correct (75.70)

Iteration 300, loss = 0.4231
Checking accuracy on validation set
Got 747 / 1000 correct (74.70)

Iteration 400, loss = 0.6017
Checking accuracy on validation set
Got 755 / 1000 correct (75.50)

```
Iteration 500, loss = 0.5795
Checking accuracy on validation set
Got 732 / 1000 correct (73.20)
```

```
Iteration 600, loss = 0.7307
Checking accuracy on validation set
Got 742 / 1000 correct (74.20)
```

```
Iteration 700, loss = 0.7427
Checking accuracy on validation set
Got 750 / 1000 correct (75.00)
```

7.1 Describe what you did

In the cell below you should write an explanation of what you did, any additional features that you implemented, and/or any graphs that you made in the process of training and evaluating your network.

My Network Architecture:

The model architecture consists of three sets of **convolutional layers** followed by **batch normalization**, **ReLU activation**, and **max pooling** 3 times. After the convolutional layers, the output is flattened and passed through a **dropout** layer with a dropout rate of 0.2. Finally, a fully connected layer maps the flattened features to the output classes.

I used the **Adam** optimizer with a learning rate of 1e-3 for training the model. The Adam optimizer combines the concepts of adaptive learning rates and momentum, which helps in efficient gradient-based optimization.

3 * [CONV >> Batch Normalization >> Relu Activation >> Max Pooling] >> Dropout >> Fully Connected

To evaluate the performance of the network, I trained the model on the training dataset for a specified number of epochs. The training progress, including the loss value at each iteration, was printed during the training process. After training, the model was evaluated on the validation set using the `check_accuracy` function.

Those are the filters I tried:

Filters: 5x5, 3x3

Num of filters: 16, 32, 64

```
[42]: from torchsummary import summary

# Print the summary of the model
summary(model, input_size=(in_channels_1, 32, 32))
```

```
-----
Layer (type)           Output Shape          Param #
=====
```

Conv2d-1	[-1, 16, 32, 32]	1,216
BatchNorm2d-2	[-1, 16, 32, 32]	32
ReLU-3	[-1, 16, 32, 32]	0
MaxPool2d-4	[-1, 16, 16, 16]	0
Conv2d-5	[-1, 32, 16, 16]	4,640
BatchNorm2d-6	[-1, 32, 16, 16]	64
ReLU-7	[-1, 32, 16, 16]	0
MaxPool2d-8	[-1, 32, 8, 8]	0
Conv2d-9	[-1, 64, 8, 8]	18,496
BatchNorm2d-10	[-1, 64, 8, 8]	128
ReLU-11	[-1, 64, 8, 8]	0
MaxPool2d-12	[-1, 64, 4, 4]	0
Flatten-13	[-1, 1024]	0
Dropout-14	[-1, 1024]	0
Linear-15	[-1, 10]	10,250

```

=====
Total params: 34,826
Trainable params: 34,826
Non-trainable params: 0
-----
Input size (MB): 0.01
Forward/backward pass size (MB): 0.73
Params size (MB): 0.13
Estimated Total Size (MB): 0.87
-----

```

7.2 Test set – run this only once

Now that we've gotten a result we're happy with, we test our final model on the test set (which you should store in `best_model`). Think about how this compares to your validation set accuracy.

```
[40]: best_model = model
      check_accuracy_part34(loader_test, best_model)
```

```

Checking accuracy on test set
Got 7552 / 10000 correct (75.52)

```