

# FullyConnectedNets

May 10, 2020

```
[346]: # 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. 'cs231n/assignments/assignment3/'
FOLDERNAME = "cs231n/assignments/assignment2/"
assert FOLDERNAME is not None, "[!] Enter the foldername."

# symlink to make it easier to load your files
!ln -s "/content/drive/My Drive/$FOLDERNAME" "/content/assignment2"

# 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/cs231n/datasets/
!bash get_datasets.sh
%cd /content
```

Mounted at /content/drive

ln: failed to create symbolic link '/content/assignment2/assignment2': Operation not supported

/content/drive/My Drive/cs231n/assignments/assignment2/cs231n/datasets  
/content

## 1 Fully-Connected Neural Nets

In the previous homework you implemented a fully-connected two-layer neural network on CIFAR-10. The implementation was simple but not very modular since the loss and gradient were computed in a single monolithic function. This is manageable for a simple two-layer network, but would become impractical as we move to bigger models. Ideally we want to build networks using

a more modular design so that we can implement different layer types in isolation and then snap them together into models with different architectures.

In this exercise we will implement fully-connected networks using a more modular approach. For each layer we will implement a forward and a backward function. The forward function will receive inputs, weights, and other parameters and will return both an output and a cache object storing data needed for the backward pass, like this:

```
def layer_forward(x, w):
    """ Receive inputs x and weights w """
    # Do some computations ...
    z = # ... some intermediate value
    # Do some more computations ...
    out = # the output

    cache = (x, w, z, out) # Values we need to compute gradients

    return out, cache
```

The backward pass will receive upstream derivatives and the cache object, and will return gradients with respect to the inputs and weights, like this:

```
def layer_backward(dout, cache):
    """
    Receive dout (derivative of loss with respect to outputs) and cache,
    and compute derivative with respect to inputs.
    """
    # Unpack cache values
    x, w, z, out = cache

    # Use values in cache to compute derivatives
    dx = # Derivative of loss with respect to x
    dw = # Derivative of loss with respect to w

    return dx, dw
```

After implementing a bunch of layers this way, we will be able to easily combine them to build classifiers with different architectures.

In addition to implementing fully-connected networks of arbitrary depth, we will also explore different update rules for optimization, and introduce Dropout as a regularizer and Batch/Layer Normalization as a tool to more efficiently optimize deep networks.

```
[347]: # As usual, a bit of setup
from __future__ import print_function
import time
import numpy as np
import matplotlib.pyplot as plt
from cs231n.classifiers.fc_net import *
from cs231n.data_utils import get_CIFAR10_data
```

```

from cs231n.gradient_check import eval_numerical_gradient, \
    eval_numerical_gradient_array
from cs231n.solver import Solver

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

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

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

```

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

```
%reload_ext autoreload
```

[348]: *# 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,))

```

## 2 Affine layer: forward

Open the file `cs231n/layers.py` and implement the `affine_forward` function.

Once you are done you can test your implementaion by running the following:

[349]: *# Test the affine\_forward function*

```

num_inputs = 2
input_shape = (4, 5, 6)
output_dim = 3

input_size = num_inputs * np.prod(input_shape)

```

```

weight_size = output_dim * np.prod(input_shape)

x = np.linspace(-0.1, 0.5, num=input_size).reshape(num_inputs, *input_shape)
w = np.linspace(-0.2, 0.3, num=weight_size).reshape(np.prod(input_shape),
→output_dim)
b = np.linspace(-0.3, 0.1, num=output_dim)

out, _ = affine_forward(x, w, b)
correct_out = np.array([[ 1.49834967,  1.70660132,  1.91485297],
                        [ 3.25553199,  3.5141327,  3.77273342]])

# Compare your output with ours. The error should be around e-9 or less.
print('Testing affine_forward function:')
print('difference: ', rel_error(out, correct_out))

```

Testing affine\_forward function:  
difference: 9.769849468192957e-10

### 3 Affine layer: backward

Now implement the `affine_backward` function and test your implementation using numeric gradient checking.

[350]:

```

# Test the affine_backward function
np.random.seed(231)
x = np.random.randn(10, 2, 3)
w = np.random.randn(6, 5)
b = np.random.randn(5)
dout = np.random.randn(10, 5)

dx_num = eval_numerical_gradient_array(lambda x: affine_forward(x, w, b)[0], x,
→dout)
dw_num = eval_numerical_gradient_array(lambda w: affine_forward(x, w, b)[0], w,
→dout)
db_num = eval_numerical_gradient_array(lambda b: affine_forward(x, w, b)[0], b,
→dout)

_, cache = affine_forward(x, w, b)
dx, dw, db = affine_backward(dout, cache)

# The error should be around e-10 or less
print('Testing affine_backward function:')
print('dx error: ', rel_error(dx_num, dx))
print('dw error: ', rel_error(dw_num, dw))
print('db error: ', rel_error(db_num, db))

```

Testing affine\_backward function:  
dx error: 5.399100368651805e-11  
dw error: 9.904211865398145e-11  
db error: 2.4122867568119087e-11

## 4 ReLU activation: forward

Implement the forward pass for the ReLU activation function in the `relu_forward` function and test your implementation using the following:

```
[351]: # Test the relu_forward function

x = np.linspace(-0.5, 0.5, num=12).reshape(3, 4)

out, _ = relu_forward(x)
correct_out = np.array([[ 0.,          0.,          0.,          0.],
                        [ 0.,          0.,          0.04545455, 0.13636364],
                        [ 0.22727273, 0.31818182, 0.40909091, 0.5]])

# Compare your output with ours. The error should be on the order of e-8
print('Testing relu_forward function:')
print('difference: ', rel_error(out, correct_out))
```

Testing relu\_forward function:  
difference: 4.999999798022158e-08

## 5 ReLU activation: backward

Now implement the backward pass for the ReLU activation function in the `relu_backward` function and test your implementation using numeric gradient checking:

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

dx_num = eval_numerical_gradient_array(lambda x: relu_forward(x)[0], x, dout)

_, cache = relu_forward(x)
dx = relu_backward(dout, cache)

# The error should be on the order of e-12
print('Testing relu_backward function:')
print('dx error: ', rel_error(dx_num, dx))
```

Testing relu\_backward function:  
dx error: 3.2756349136310288e-12

## 5.1 Inline Question 1:

We've only asked you to implement ReLU, but there are a number of different activation functions that one could use in neural networks, each with its pros and cons. In particular, an issue commonly seen with activation functions is getting zero (or close to zero) gradient flow during backpropagation. Which of the following activation functions have this problem? If you consider these functions in the one dimensional case, what types of input would lead to this behaviour? 1. Sigmoid 2. ReLU 3. Leaky ReLU

## 5.2 Answer:

Close to zero gradient flow during backpropagation leads to the vanishing gradient problem.

1. Sigmoid function suffers from the vanishing gradient problem because the gradient is close to zero for very large positive and negative values of the input (in the tail portions of the function). A one dimensional example that can lead to saturation is to consider very large positive and negative values, e.g.,  $[-1e3, 1e3]$ .
2. One of ReLU's biggest advantages over Sigmoid is the fact that it is less susceptible to the vanishing gradient problem owing to its linear response to a positive input. ReLU's gradient is 0 (for negative inputs) or 1 (for positive inputs). ReLU can suffer from the vanishing gradient problem when all the input values are negative which is not a very probable scenario. When this occurs, some neurons fail to train further. This is known as the "dying ReLU problem". A one dimensional example that can lead to the vanishing gradient problem is to consider only negative values, e.g.,  $[-1, -2, -3]$ .
3. Leaky ReLU tries to solve the ReLU problem of "dead" neurons by applying a small negative slope for negative values, i.e., if  $x < 0$  then  $0.01x$  else  $x$ . Leaky ReLU is thus aimed at solving the vanishing gradient problem. However, since the  $\max(0.01x, x)$  function is not continuous at  $x = 0$ , the gradient at  $x = 0$  is undefined. So, if it is not explicitly handled in code, a one dimensional example that can lead to zero gradients is to consider all zero values, e.g.,  $[0, 0, 0]$ , which can only happen with a bad network initialization.

## 6 "Sandwich" layers

There are some common patterns of layers that are frequently used in neural nets. For example, affine layers are frequently followed by a ReLU nonlinearity. To make these common patterns easy, we define several convenience layers in the file `cs231n/layer_utils.py`.

For now take a look at the `affine_relu_forward` and `affine_relu_backward` functions, and run the following to numerically gradient check the backward pass:

```
[353]: from cs231n.layer_utils import affine_relu_forward, affine_relu_backward
np.random.seed(231)
x = np.random.randn(2, 3, 4)
w = np.random.randn(12, 10)
b = np.random.randn(10)
dout = np.random.randn(2, 10)

out, cache = affine_relu_forward(x, w, b)
dx, dw, db = affine_relu_backward(dout, cache)
```

```

dx_num = eval_numerical_gradient_array(lambda x: affine_relu_forward(x, w,
    ↪b)[0], x, dout)
dw_num = eval_numerical_gradient_array(lambda w: affine_relu_forward(x, w,
    ↪b)[0], w, dout)
db_num = eval_numerical_gradient_array(lambda b: affine_relu_forward(x, w,
    ↪b)[0], b, dout)

# Relative error should be around e-10 or less
print('Testing affine_relu_forward and affine_relu_backward:')
print('dx error: ', rel_error(dx_num, dx))
print('dw error: ', rel_error(dw_num, dw))
print('db error: ', rel_error(db_num, db))

```

Testing affine\_relu\_forward and affine\_relu\_backward:

dx error: 2.299579177309368e-11

dw error: 8.162011105764925e-11

db error: 7.826724021458994e-12

## 7 Loss layers: Softmax and SVM

You implemented these loss functions in the last assignment, so we'll give them to you for free here. You should still make sure you understand how they work by looking at the implementations in `cs231n/layers.py`.

You can make sure that the implementations are correct by running the following:

```

[354]: np.random.seed(231)
num_classes, num_inputs = 10, 50
x = 0.001 * np.random.randn(num_inputs, num_classes)
y = np.random.randint(num_classes, size=num_inputs)

dx_num = eval_numerical_gradient(lambda x: svm_loss(x, y)[0], x, verbose=False)
loss, dx = svm_loss(x, y)

# Test svm_loss function. Loss should be around 9 and dx error should be around
    ↪the order of e-9
print('Testing svm_loss:')
print('loss: ', loss)
print('dx error: ', rel_error(dx_num, dx))

dx_num = eval_numerical_gradient(lambda x: softmax_loss(x, y)[0], x,
    ↪verbose=False)
loss, dx = softmax_loss(x, y)

# Test softmax_loss function. Loss should be close to 2.3 and dx error should
    ↪be around e-8
print('\nTesting softmax_loss:')

```

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

Testing svm\_loss:

loss: 8.999602749096233

dx error: 1.4021566006651672e-09

Testing softmax\_loss:

loss: 2.302545844500738

dx error: 9.384673161989355e-09

```
[355]: np.random.seed(231)
num_classes, num_inputs = 10, 50
x = 0.001 * np.random.randn(num_inputs, num_classes)
y = np.random.randint(num_classes, size=num_inputs)

dx_num = eval_numerical_gradient(lambda x: svm_loss(x, y)[0], x, verbose=False)
loss, dx = svm_loss(x, y)

# Test svm_loss function. Loss should be around 9 and dx error should be around
→the order of e-9
print('Testing svm_loss:')
print('loss: ', loss)
print('dx error: ', rel_error(dx_num, dx))

dx_num = eval_numerical_gradient(lambda x: softmax_loss(x, y)[0], x,
    →verbose=False)
loss, dx = softmax_loss(x, y)

# Test softmax_loss function. Loss should be close to 2.3 and dx error should
→be around e-8
print('\nTesting softmax_loss:')
print('loss: ', loss)
print('dx error: ', rel_error(dx_num, dx))
```

Testing svm\_loss:

loss: 8.999602749096233

dx error: 1.4021566006651672e-09

Testing softmax\_loss:

loss: 2.302545844500738

dx error: 9.384673161989355e-09

```
[356]: np.random.seed(231)
num_classes, num_inputs = 10, 50
x = 0.001 * np.random.randn(num_inputs, num_classes)
```



```

y = np.random.randint(num_classes, size=num_inputs)

dx_num = eval_numerical_gradient(lambda x: svm_loss(x, y)[0], x, verbose=False)
loss, dx = svm_loss(x, y)

# Test svm_loss function. Loss should be around 9 and dx error should be around
→the order of e-9
print('Testing svm_loss:')
print('loss: ', loss)
print('dx error: ', rel_error(dx_num, dx))

dx_num = eval_numerical_gradient(lambda x: softmax_loss(x, y)[0], x,
→verbose=False)
loss, dx = softmax_loss(x, y)

# Test softmax_loss function. Loss should be close to 2.3 and dx error should
→be around e-8
print('\nTesting softmax_loss:')
print('loss: ', loss)
print('dx error: ', rel_error(dx_num, dx))

```

```

Testing svm_loss:
loss:  8.999602749096233
dx error:  1.4021566006651672e-09

```

```

Testing softmax_loss:
loss:  2.302545844500738
dx error:  9.384673161989355e-09

```

## 8 Two-layer network

In the previous assignment you implemented a two-layer neural network in a single monolithic class. Now that you have implemented modular versions of the necessary layers, you will reimplement the two layer network using these modular implementations.

Open the file `cs231n/classifiers/fc_net.py` and complete the implementation of the `TwoLayerNet` class. This class will serve as a model for the other networks you will implement in this assignment, so read through it to make sure you understand the API. You can run the cell below to test your implementation.

```

[357]: np.random.seed(231)
N, D, H, C = 3, 5, 50, 7
X = np.random.randn(N, D)
y = np.random.randint(C, size=N)

std = 1e-3
model = TwoLayerNet(input_dim=D, hidden_dim=H, num_classes=C, weight_scale=std)

```

```

print('Testing initialization ... ')
W1_std = abs(model.params['W1'].std() - std)
b1 = model.params['b1']
W2_std = abs(model.params['W2'].std() - std)
b2 = model.params['b2']
assert W1_std < std / 10, 'First layer weights do not seem right'
assert np.all(b1 == 0), 'First layer biases do not seem right'
assert W2_std < std / 10, 'Second layer weights do not seem right'
assert np.all(b2 == 0), 'Second layer biases do not seem right'

print('Testing test-time forward pass ... ')
model.params['W1'] = np.linspace(-0.7, 0.3, num=D*H).reshape(D, H)
model.params['b1'] = np.linspace(-0.1, 0.9, num=H)
model.params['W2'] = np.linspace(-0.3, 0.4, num=H*C).reshape(H, C)
model.params['b2'] = np.linspace(-0.9, 0.1, num=C)
X = np.linspace(-5.5, 4.5, num=N*D).reshape(D, N).T
scores = model.loss(X)
correct_scores = np.asarray(
    [[11.53165108, 12.2917344, 13.05181771, 13.81190102, 14.57198434, 15.
→33206765, 16.09215096],
    [12.05769098, 12.74614105, 13.43459113, 14.1230412, 14.81149128, 15.
→49994135, 16.18839143],
    [12.58373087, 13.20054771, 13.81736455, 14.43418138, 15.05099822, 15.
→66781506, 16.2846319 ]])
scores_diff = np.abs(scores - correct_scores).sum()
assert scores_diff < 1e-6, 'Problem with test-time forward pass'

print('Testing training loss (no regularization)')
y = np.asarray([0, 5, 1])
loss, grads = model.loss(X, y)
correct_loss = 3.4702243556
assert abs(loss - correct_loss) < 1e-10, 'Problem with training-time loss'

model.reg = 1.0
loss, grads = model.loss(X, y)
correct_loss = 26.5948426952
assert abs(loss - correct_loss) < 1e-10, 'Problem with regularization loss'

# Errors should be around e-7 or less
for reg in [0.0, 0.7]:
    print('Running numeric gradient check with reg = ', reg)
    model.reg = reg
    loss, grads = model.loss(X, y)

    for name in sorted(grads):
        f = lambda _: model.loss(X, y)[0]
        grad_num = eval_numerical_gradient(f, model.params[name], verbose=False)

```

```
print('%s relative error: %.2e' % (name, rel_error(grad_num, grads[name])))
```

```
Testing initialization ...
Testing test-time forward pass ...
Testing training loss (no regularization)
Running numeric gradient check with reg = 0.0
W1 relative error: 1.83e-08
W2 relative error: 3.12e-10
b1 relative error: 9.83e-09
b2 relative error: 4.33e-10
Running numeric gradient check with reg = 0.7
W1 relative error: 2.53e-07
W2 relative error: 2.85e-08
b1 relative error: 1.56e-08
b2 relative error: 7.76e-10
```

## 9 Solver

In the previous assignment, the logic for training models was coupled to the models themselves. Following a more modular design, for this assignment we have split the logic for training models into a separate class.

Open the file `cs231n/solver.py` and read through it to familiarize yourself with the API. After doing so, use a Solver instance to train a `TwoLayerNet` that achieves at least 50% accuracy on the validation set.

```
[358]: model = TwoLayerNet()
       solver = None

#####
# TODO: Use a Solver instance to train a TwoLayerNet that achieves at least #
# 50% accuracy on the validation set.                                     #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

# generate random combinations of hyperparameters given min-max ranges of
→hyperparams
# using a uniform probability distribution
# usage: lr, reg, hidden_dim, epochs = generate_random_hyperparams((-1, 0),
→(-7, -4), (10, 500), (10, 20))
def generate_random_hyperparams(lr, reg, hidden_size, epoch_values):
    lr = 10*np.random.uniform(lr[0], lr[1])
    reg = 10*np.random.uniform(reg[0], reg[1])
    hidden_size = np.random.randint(hidden_size[0], hidden_size[1])
    epochs = epoch_values[np.random.randint(0, len(epoch_values))]

    return lr, reg, hidden_size, epochs
```

```

# number of hypercombinations combinations to look for
num_hyperparam_configs = 10

# form random combinations of learning_rate, regularization_strength,
→hidden_layer_size and num_training_epochs
grid_search = [generate_random_hyperparams((-2, -4), (-7, -3), (50, 100, 200),
→(10, 20))
                for count in range(num_hyperparam_configs)]

# get our data
# data = get_CIFAR10_data()

best_val = -1    # The highest validation accuracy that we have seen so far.
results = {}

for config_num, config in enumerate(grid_search):
    print("Hyperparam config #{} of {}: ".format(config_num+1,
→len(grid_search)), end='')

    lr, reg, hidden_size, epochs = config
    print("lr: {:.2e}, reg: {:.2e}, hidden_size: {:.2e}, epochs: {:.2e}".
→format(lr, reg, hidden_size, epochs))

    model = TwoLayerNet(hidden_dim=hidden_size, reg=reg)
    current_solver = Solver(model, data, update_rule='sgd',
→optim_config={'learning_rate': lr},
                    lr_decay=0.95, num_epochs=epochs, batch_size=100,
                    print_every=100, verbose=False)

    # train a 2-layer neural net on the training set
    current_solver.train()

    # store the best validation accuracy and the TwoLayerNet object
    if current_solver.best_val_acc > best_val:
        best_val = current_solver.best_val_acc
        solver = current_solver

    # print results in-line to avoid saving them
    print('Validation accuracy: %.4f' % (solver.best_val_acc,))
    print() # new line after every hyperparam config

print('Best validation accuracy achieved: %.4f' % best_val)

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

```

Hyperparam config #1 of #10: lr: 1.80e-04, reg: 1.37e-06, hidden\_size: 7.40e+01,  
epochs: 1.00e+01  
Validation accuracy: 0.4910

Hyperparam config #2 of #10: lr: 2.15e-03, reg: 8.08e-04, hidden\_size: 9.30e+01,  
epochs: 1.00e+01  
Validation accuracy: 0.4910

Hyperparam config #3 of #10: lr: 2.31e-04, reg: 1.26e-05, hidden\_size: 7.70e+01,  
epochs: 2.00e+01  
Validation accuracy: 0.5230

Hyperparam config #4 of #10: lr: 1.30e-03, reg: 1.33e-05, hidden\_size: 7.20e+01,  
epochs: 1.00e+01  
Validation accuracy: 0.5230

Hyperparam config #5 of #10: lr: 3.53e-04, reg: 1.69e-05, hidden\_size: 6.80e+01,  
epochs: 1.00e+01  
Validation accuracy: 0.5230

Hyperparam config #6 of #10: lr: 3.40e-03, reg: 4.12e-05, hidden\_size: 5.80e+01,  
epochs: 2.00e+01  
Validation accuracy: 0.5230

Hyperparam config #7 of #10: lr: 2.93e-04, reg: 5.47e-05, hidden\_size: 8.20e+01,  
epochs: 1.00e+01  
Validation accuracy: 0.5230

Hyperparam config #8 of #10: lr: 6.33e-04, reg: 7.90e-06, hidden\_size: 7.40e+01,  
epochs: 1.00e+01  
Validation accuracy: 0.5260

Hyperparam config #9 of #10: lr: 3.99e-04, reg: 7.09e-06, hidden\_size: 8.90e+01,  
epochs: 2.00e+01  
Validation accuracy: 0.5360

Hyperparam config #10 of #10: lr: 1.83e-04, reg: 2.24e-06, hidden\_size:  
5.60e+01, epochs: 2.00e+01  
Validation accuracy: 0.5360

Best validation accuracy achieved: 0.5360

[359]: *# Run this cell to visualize training loss and train / val accuracy*

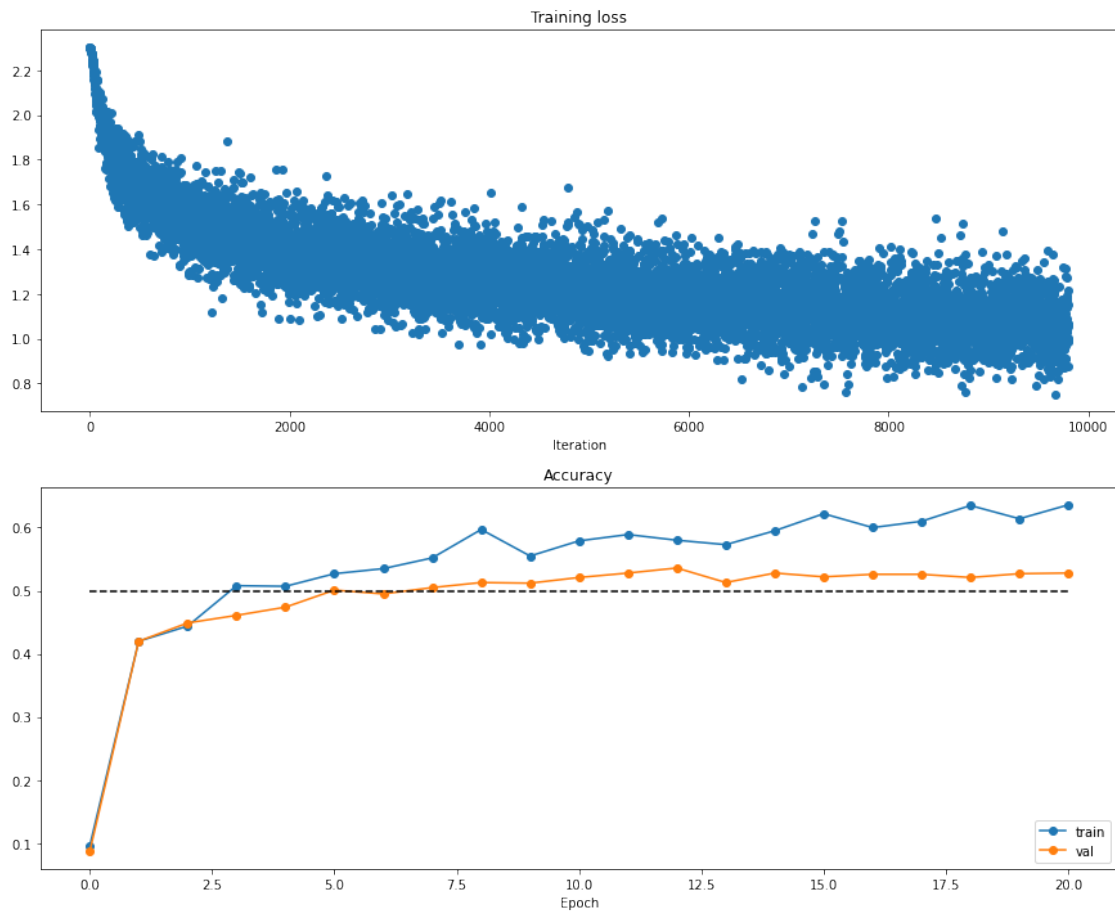
```
plt.subplot(2, 1, 1)
plt.title('Training loss')
plt.plot(solver.loss_history, 'o')
```

```

plt.xlabel('Iteration')

plt.subplot(2, 1, 2)
plt.title('Accuracy')
plt.plot(solver.train_acc_history, '-o', label='train')
plt.plot(solver.val_acc_history, '-o', label='val')
plt.plot([0.5] * len(solver.val_acc_history), 'k--')
plt.xlabel('Epoch')
plt.legend(loc='lower right')
plt.gcf().set_size_inches(15, 12)
plt.show()

```



## 10 Multilayer network

Next you will implement a fully-connected network with an arbitrary number of hidden layers. Read through the `FullyConnectedNet` class in the file `cs231n/classifiers/fc_net.py`.

Implement the initialization, the forward pass, and the backward pass. For the moment don't worry about implementing dropout or batch/layer normalization; we will add those features soon.

## 10.1 Initial loss and gradient check

As a sanity check, run the following to check the initial loss and to gradient check the network both with and without regularization. Do the initial losses seem reasonable?

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

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

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

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

    # Most of the errors should be on the order of e-7 or smaller.
    # NOTE: It is fine however to see an error for W2 on the order of e-5
    # for the check when reg = 0.0
    for name in sorted(grads):
        f = lambda _: model.loss(X, y)[0]
        grad_num = eval_numerical_gradient(f, model.params[name], verbose=False,
        →h=1e-5)
        print('%s relative error: %.2e' % (name, rel_error(grad_num, grads[name])))
```

```
Running check with reg = 0
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 reg = 3.14
Initial loss: 7.052114776533016
W1 relative error: 3.90e-09
W2 relative error: 6.87e-08
W3 relative error: 2.13e-08
b1 relative error: 1.48e-08
b2 relative error: 1.72e-09
b3 relative error: 1.57e-10
```

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

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

```
num_train = 50
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

# obtained using random search
weight_scale = 2e-2 # Experiment with this!
learning_rate = 4e-3 # Experiment with this!

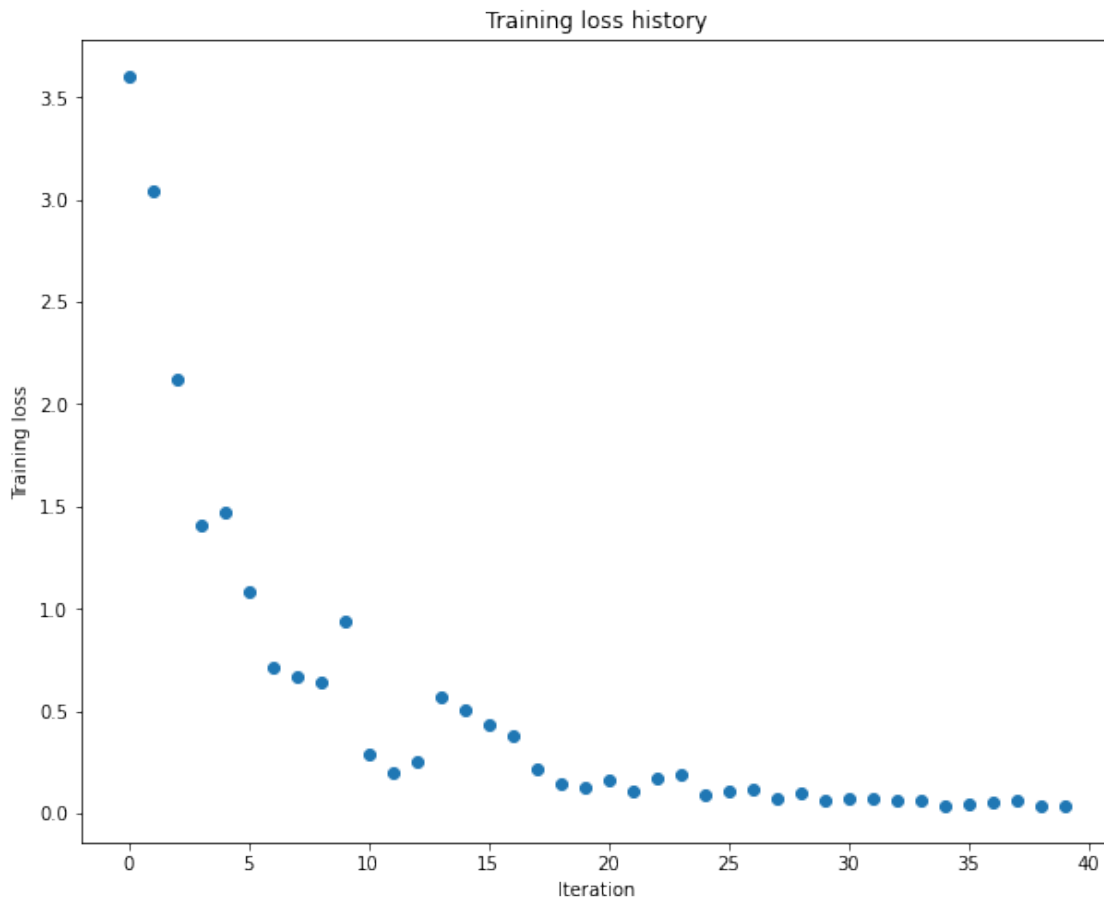
model = FullyConnectedNet([100, 100],
                           weight_scale=weight_scale, dtype=np.float64)
solver = Solver(model, small_data,
                print_every=10, num_epochs=20, batch_size=25,
                update_rule='sgd',
                optim_config={
                    'learning_rate': learning_rate,
                })
solver.train()

plt.plot(solver.loss_history, 'o')
plt.title('Training loss history')
plt.xlabel('Iteration')
plt.ylabel('Training loss')
plt.show()
```

```
(Iteration 1 / 40) loss: 3.604568
(Epoch 0 / 20) train acc: 0.200000; val_acc: 0.109000
(Epoch 1 / 20) train acc: 0.460000; val_acc: 0.141000
(Epoch 2 / 20) train acc: 0.700000; val_acc: 0.157000
(Epoch 3 / 20) train acc: 0.800000; val_acc: 0.175000
(Epoch 4 / 20) train acc: 0.840000; val_acc: 0.183000
(Epoch 5 / 20) train acc: 0.900000; val_acc: 0.199000
(Iteration 11 / 40) loss: 0.292878
(Epoch 6 / 20) train acc: 0.900000; val_acc: 0.197000
(Epoch 7 / 20) train acc: 0.940000; val_acc: 0.192000
(Epoch 8 / 20) train acc: 0.980000; val_acc: 0.206000
```



```
(Epoch 9 / 20) train acc: 1.000000; val_acc: 0.197000
(Epoch 10 / 20) train acc: 0.980000; val_acc: 0.190000
(Iteration 21 / 40) loss: 0.164482
(Epoch 11 / 20) train acc: 1.000000; val_acc: 0.195000
(Epoch 12 / 20) train acc: 1.000000; val_acc: 0.178000
(Epoch 13 / 20) train acc: 1.000000; val_acc: 0.185000
(Epoch 14 / 20) train acc: 1.000000; val_acc: 0.189000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.190000
(Iteration 31 / 40) loss: 0.070067
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.199000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.193000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.205000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.197000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.194000
```



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

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

```
num_train = 50
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

learning_rate = 2e-3 # Experiment with this!
weight_scale = 6e-2 # Experiment with this!
model = FullyConnectedNet([100, 100, 100, 100],
                           weight_scale=weight_scale, dtype=np.float64)
solver = Solver(model, small_data,
                print_every=10, num_epochs=20, batch_size=25,
                update_rule='sgd',
                optim_config={
                    'learning_rate': learning_rate,
                })
solver.train()

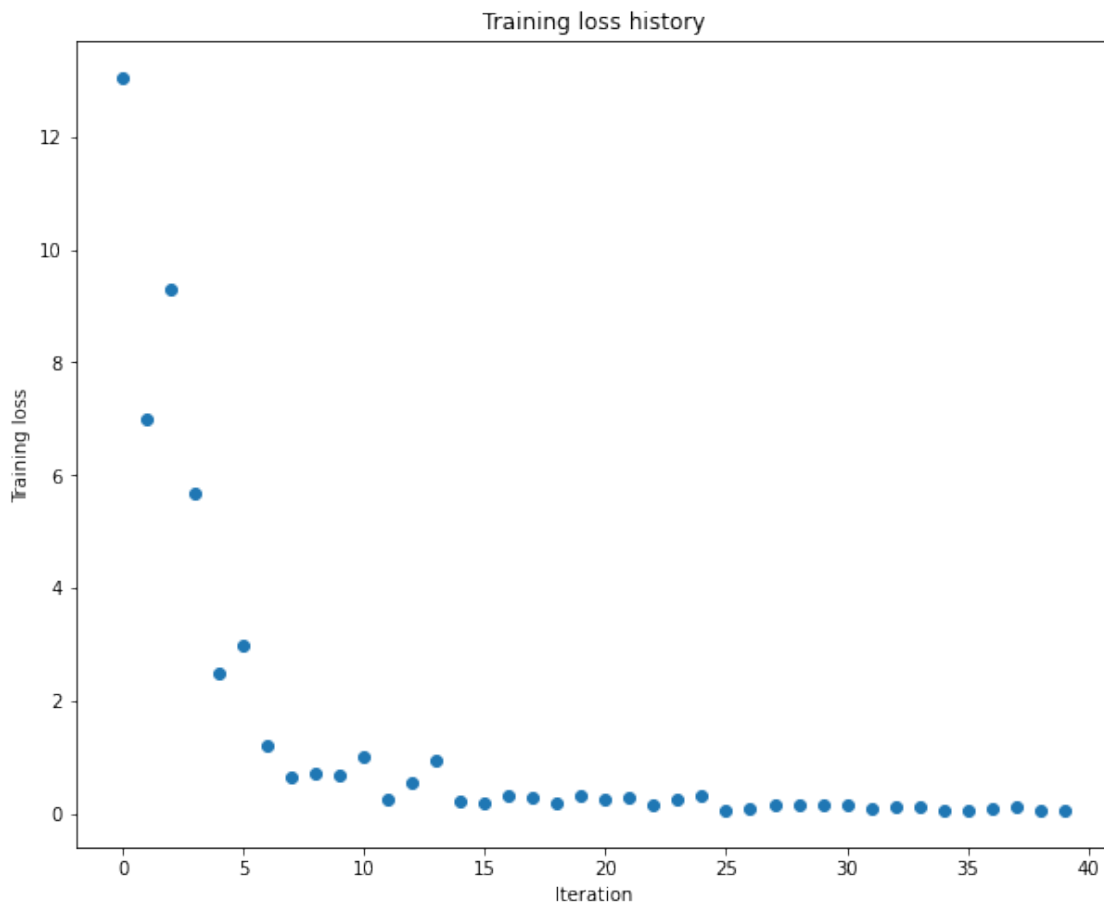
plt.plot(solver.loss_history, 'o')
plt.title('Training loss history')
plt.xlabel('Iteration')
plt.ylabel('Training loss')
plt.show()
```

```
(Iteration 1 / 40) loss: 13.054907
(Epoch 0 / 20) train acc: 0.260000; val_acc: 0.115000
(Epoch 1 / 20) train acc: 0.240000; val_acc: 0.088000
(Epoch 2 / 20) train acc: 0.340000; val_acc: 0.136000
(Epoch 3 / 20) train acc: 0.580000; val_acc: 0.131000
(Epoch 4 / 20) train acc: 0.740000; val_acc: 0.132000
(Epoch 5 / 20) train acc: 0.840000; val_acc: 0.128000
(Iteration 11 / 40) loss: 1.025352
(Epoch 6 / 20) train acc: 0.880000; val_acc: 0.133000
(Epoch 7 / 20) train acc: 0.880000; val_acc: 0.132000
(Epoch 8 / 20) train acc: 0.900000; val_acc: 0.128000
(Epoch 9 / 20) train acc: 0.960000; val_acc: 0.126000
(Epoch 10 / 20) train acc: 0.960000; val_acc: 0.122000
(Iteration 21 / 40) loss: 0.274423
(Epoch 11 / 20) train acc: 0.980000; val_acc: 0.125000
(Epoch 12 / 20) train acc: 0.980000; val_acc: 0.132000
(Epoch 13 / 20) train acc: 0.980000; val_acc: 0.133000
```

```

(Epoch 14 / 20) train acc: 0.980000; val_acc: 0.127000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.128000
(Iteration 31 / 40) loss: 0.165667
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.130000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.135000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.133000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.130000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.127000

```



## 10.2 Inline Question 2:

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

## 10.3 Answer:

- Comparing the difficulty of training the three-layer net vs training the five-layer net, the five-layer net was much more sensitive to the weight initialization scale. This is due to the

fact that the five-layer net is 66% more deeper than the three-layer net. With deep networks, a small weight scale leads to vanishing gradients (small products), while a large weight scale leads to exploding gradients (large products). The deeper the network, higher the probability of vanishing gradient issues. Finding the correct weight initialization scale for the five-layer net is thus harder than the three-layer net.

- In the below example, we can observe the sensitivity in the weight initialization scale with a five-layer net. Initially, with a small weight scale of  $1e-5$ , the mean and standard deviation of the parameters are observed to be close to zero (leading to vanishing gradients). Increasing the weight scale by a couple of orders of magnitude improves these values, leading to the best training performance with a weight scale of  $5e-2$ . However, nudging it further by an order of magnitude to  $1e-1$  leads to exploding gradients. While a shallower neural network such as one with three or two hidden layers also shows sensitivity to the weight scale, the magnitude is much lesser with shallower networks. The three-layer net's weight scale is relatively easier to find ( $2e-2$ ), unlike the five-layer net that requires a more specific weight scale ( $5e-2$ ). On a different note, we also noticed that the weight scale values depend on the learning rate.

```
[363]: # sensitivity to initialization scale

num_train = 50
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

# for deterministic/reproducible results
np.random.seed(1)

learning_rate = 1e-2
weight_scales = [1e-5, 1e-3, 1e-2, 5e-2, 1e-1, 5e-1]

for weight_scale in weight_scales:

    # train a 5-layer neural net on the training set
    model = FullyConnectedNet([100, 100, 100, 100],
                               weight_scale=weight_scale, dtype=np.float64)
    solver = Solver(model, small_data,
                     print_every=10, num_epochs=20, batch_size=25,
                     update_rule='sgd',
                     optim_config={
                         'learning_rate': learning_rate,
                     },
                     verbose=False
                    )

    solver.train()
```

```

train_accuracy = solver.best_train_acc # new field added to solver.py;
                                         # this holds the training accuracy
                                         # spotted before overfitting begins

val_accuracy = solver.best_val_acc
print("lr: {:.2e}, weight scale: {:.2e}, train accuracy: {:.2f}, val_
→accuracy: {:.2f}"\
      .format(learning_rate, weight_scale, train_accuracy, val_accuracy))

# print mean and std dev of weights
for k, v in sorted(model.params.items()):
    if k[0] == 'W': #only print weights
        print("%s -> mean: %lf, std: %lf" % (k, v.mean(), v.std()))
print()

```

```

lr: 1.00e-02, weight scale: 1.00e-05, train accuracy: 0.16, val accuracy: 0.11
W1 -> mean: 0.000000, std: 0.000010
W2 -> mean: -0.000000, std: 0.000010
W3 -> mean: 0.000000, std: 0.000010
W4 -> mean: 0.000000, std: 0.000010
W5 -> mean: -0.000000, std: 0.000010

```

```

lr: 1.00e-02, weight scale: 1.00e-03, train accuracy: 0.16, val accuracy: 0.11
W1 -> mean: -0.000001, std: 0.000999
W2 -> mean: -0.000008, std: 0.000997
W3 -> mean: -0.000009, std: 0.000999
W4 -> mean: 0.000014, std: 0.000995
W5 -> mean: 0.000061, std: 0.000975

```

```

lr: 1.00e-02, weight scale: 1.00e-02, train accuracy: 0.10, val accuracy: 0.11
W1 -> mean: -0.000005, std: 0.010000
W2 -> mean: 0.000012, std: 0.009918
W3 -> mean: 0.000009, std: 0.009958
W4 -> mean: 0.000079, std: 0.009939
W5 -> mean: 0.000143, std: 0.009803

```

```

lr: 1.00e-02, weight scale: 5.00e-02, train accuracy: 1.00, val accuracy: 0.17
W1 -> mean: 0.000027, std: 0.049990
W2 -> mean: -0.000561, std: 0.049903
W3 -> mean: 0.000143, std: 0.049845
W4 -> mean: -0.000881, std: 0.050095
W5 -> mean: -0.000547, std: 0.051637

```

```

lr: 1.00e-02, weight scale: 1.00e-01, train accuracy: 0.14, val accuracy: 0.13
W1 -> mean: -92904549857997046241368316112273408.000000, std:
1016972212849067896351167309930299392.000000
W2 -> mean: -71342956438476489197956625604168281751552.000000, std:

```

```

647222390529915769157952456889291792973824.000000
W3 -> mean: -76743173575393369925052441249306049742110720.000000, std:
1247698592842431146193283669261361189660655616.000000
W4 -> mean: -174214523618833681119023007966257167204352.000000, std:
3537329354785940224223184722588870264225792.000000
W5 -> mean: 9827244656481231962112.000000, std:
3203830765901411650654183640603220246528.000000

/content/drive/My Drive/cs231n/assignments/assignment2/cs231n/layers.py:1190:
RuntimeWarning: invalid value encountered in subtract
    shifted_logits = x - np.max(x, axis=1, keepdims=True)
/content/drive/My Drive/cs231n/assignments/assignment2/cs231n/layers.py:134:
RuntimeWarning: invalid value encountered in less
    dx[x < 0] = 0

lr: 1.00e-02, weight scale: 5.00e-01, train accuracy: 0.16, val accuracy: 0.11
W1 -> mean: -23744380149024319325891047981056.000000, std:
97116811295727163944483764043776.000000
W2 -> mean: -211847131294884436391421092015636480.000000, std:
2708588991607482183378979652146036736.000000
W3 -> mean: -2090359013607242829085789954024407040.000000, std:
138536107366855448803933158830757642240.000000
W4 -> mean: -893463820957497534924597082718208.000000, std:
24208472038810557266262360374378496.000000
W5 -> mean: -1855483183915597.750000, std:
1239582396156082958021079123099648.000000

```

## 11 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.

## 12 SGD+Momentum

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

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

```
[364]: from cs231n.optim import sgd_momentum
```

```
N, D = 4, 5
```

```

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

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

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

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

```

```

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

```

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

```

[365]: 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: 2.492729
(Epoch 0 / 5) train acc: 0.120000; val_acc: 0.118000
(Iteration 11 / 200) loss: 2.268942
(Iteration 21 / 200) loss: 2.172473
(Iteration 31 / 200) loss: 2.187313
(Epoch 1 / 5) train acc: 0.237000; val_acc: 0.231000
(Iteration 41 / 200) loss: 2.114944
(Iteration 51 / 200) loss: 2.108276
(Iteration 61 / 200) loss: 2.058074
(Iteration 71 / 200) loss: 1.980620
(Epoch 2 / 5) train acc: 0.286000; val_acc: 0.263000
(Iteration 81 / 200) loss: 2.124772
(Iteration 91 / 200) loss: 1.855776

```



```
(Iteration 101 / 200) loss: 1.988221
(Iteration 111 / 200) loss: 2.029581
(Epoch 3 / 5) train acc: 0.326000; val_acc: 0.297000
(Iteration 121 / 200) loss: 1.807487
(Iteration 131 / 200) loss: 2.025019
(Iteration 141 / 200) loss: 1.940683
(Iteration 151 / 200) loss: 1.796607
(Epoch 4 / 5) train acc: 0.380000; val_acc: 0.311000
(Iteration 161 / 200) loss: 1.806488
(Iteration 171 / 200) loss: 1.803253
(Iteration 181 / 200) loss: 1.815219
(Iteration 191 / 200) loss: 1.713854
(Epoch 5 / 5) train acc: 0.370000; val_acc: 0.306000
```

running with sgd\_momentum

```
(Iteration 1 / 200) loss: 2.666380
(Epoch 0 / 5) train acc: 0.128000; val_acc: 0.108000
(Iteration 11 / 200) loss: 2.208883
(Iteration 21 / 200) loss: 2.053921
(Iteration 31 / 200) loss: 2.131006
(Epoch 1 / 5) train acc: 0.301000; val_acc: 0.257000
(Iteration 41 / 200) loss: 1.837410
(Iteration 51 / 200) loss: 1.861928
(Iteration 61 / 200) loss: 1.712057
(Iteration 71 / 200) loss: 1.777077
(Epoch 2 / 5) train acc: 0.385000; val_acc: 0.328000
(Iteration 81 / 200) loss: 1.782872
(Iteration 91 / 200) loss: 1.696865
(Iteration 101 / 200) loss: 1.550729
(Iteration 111 / 200) loss: 1.692089
(Epoch 3 / 5) train acc: 0.472000; val_acc: 0.337000
(Iteration 121 / 200) loss: 1.595567
(Iteration 131 / 200) loss: 1.638184
(Iteration 141 / 200) loss: 1.496519
(Iteration 151 / 200) loss: 1.294896
(Epoch 4 / 5) train acc: 0.480000; val_acc: 0.363000
(Iteration 161 / 200) loss: 1.435100
(Iteration 171 / 200) loss: 1.321631
(Iteration 181 / 200) loss: 1.280618
(Iteration 191 / 200) loss: 1.275304
(Epoch 5 / 5) train acc: 0.529000; val_acc: 0.366000
```

/usr/local/lib/python3.6/dist-packages/ipykernel\_launcher.py:39:

MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each

axes instance.

/usr/local/lib/python3.6/dist-packages/ipykernel\_launcher.py:42:

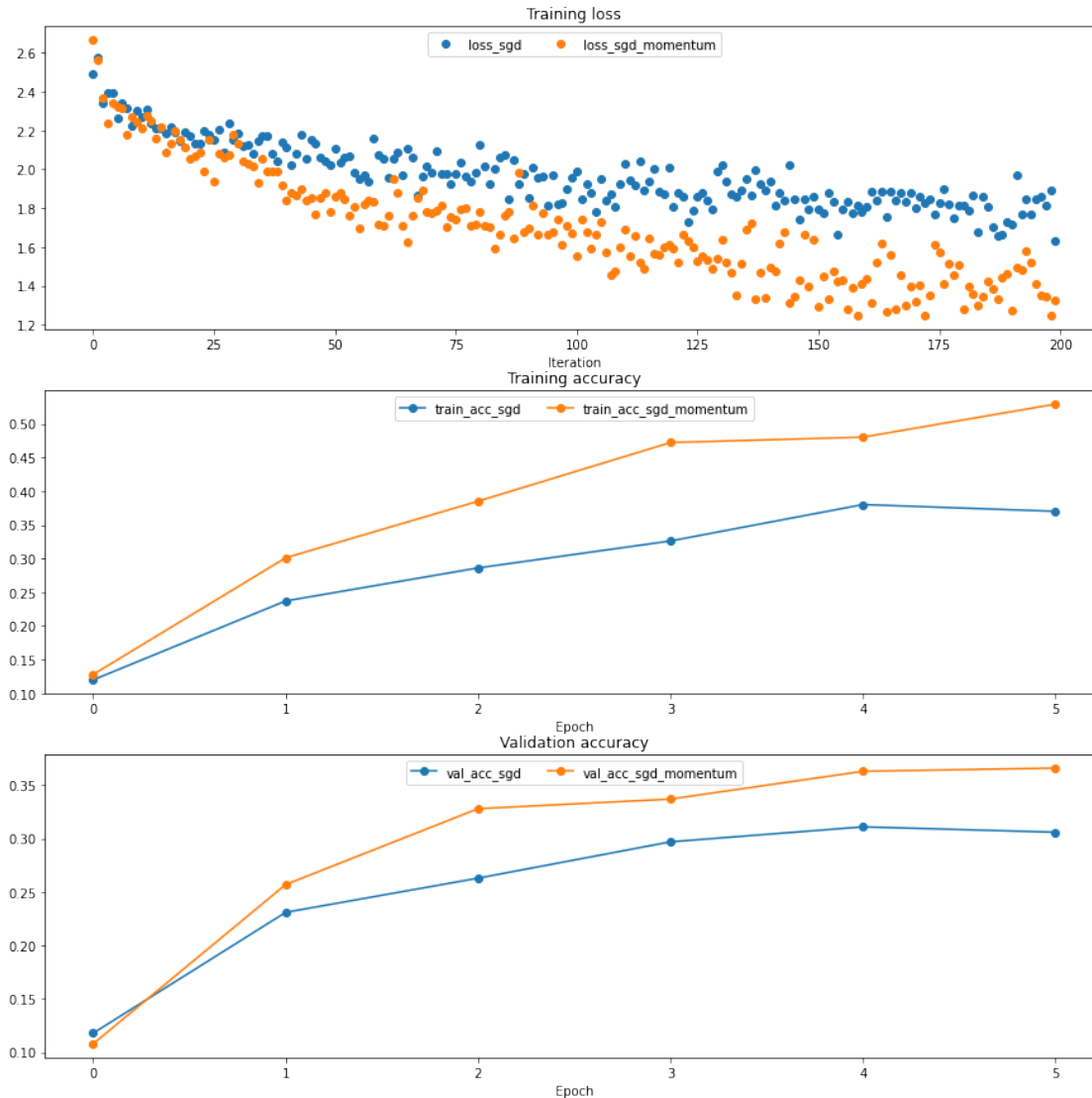
MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

/usr/local/lib/python3.6/dist-packages/ipykernel\_launcher.py:45:

MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

/usr/local/lib/python3.6/dist-packages/ipykernel\_launcher.py:49:

MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.



## 13 RMSProp and Adam

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

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

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

[1] Tijmen Tieleman and Geoffrey Hinton. "Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude." COURSE: Neural Networks for Machine Learning 4 (2012).

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

```
[366]: # Test RMSProp implementation
from cs231n.optim import rmsprop

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

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

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

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

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

```
[367]: # Test Adam implementation
from cs231n.optim import adam

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

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

expected_next_w = np.asarray([
    [-0.40094747, -0.34836187, -0.29577703, -0.24319299, -0.19060977],
    [-0.1380274, -0.08544591, -0.03286534, 0.01971428, 0.0722929],
    [ 0.1248705, 0.17744702, 0.23002243, 0.28259667, 0.33516969],
```

```

[ 0.38774145,  0.44031188,  0.49288093,  0.54544852,  0.59801459]])
expected_v = np.asarray([
[ 0.69966,      0.68908382,  0.67851319,  0.66794809,  0.65738853,],
[ 0.64683452,  0.63628604,  0.6257431,   0.61520571,  0.60467385,],
[ 0.59414753,  0.58362676,  0.57311152,  0.56260183,  0.55209767,],
[ 0.54159906,  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:

```

[368]: 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: 2.627778
(Epoch 0 / 5) train acc: 0.161000; val_acc: 0.142000
(Iteration 11 / 200) loss: 2.150370
(Iteration 21 / 200) loss: 2.030044
(Iteration 31 / 200) loss: 1.954323
(Epoch 1 / 5) train acc: 0.346000; val_acc: 0.311000
(Iteration 41 / 200) loss: 1.639610
(Iteration 51 / 200) loss: 1.751849
(Iteration 61 / 200) loss: 1.566358
(Iteration 71 / 200) loss: 1.743638
(Epoch 2 / 5) train acc: 0.468000; val_acc: 0.371000
(Iteration 81 / 200) loss: 1.605346
(Iteration 91 / 200) loss: 1.605137
(Iteration 101 / 200) loss: 1.687328
(Iteration 111 / 200) loss: 1.601762
(Epoch 3 / 5) train acc: 0.475000; val_acc: 0.348000
(Iteration 121 / 200) loss: 1.469905
(Iteration 131 / 200) loss: 1.335355
(Iteration 141 / 200) loss: 1.146708
(Iteration 151 / 200) loss: 1.381720
(Epoch 4 / 5) train acc: 0.488000; val_acc: 0.345000
(Iteration 161 / 200) loss: 1.332485
(Iteration 171 / 200) loss: 1.262432
(Iteration 181 / 200) loss: 1.155601
(Iteration 191 / 200) loss: 1.298816
(Epoch 5 / 5) train acc: 0.590000; val_acc: 0.369000

```

```
running with rmsprop
(Iteration 1 / 200) loss: 2.633719
(Epoch 0 / 5) train acc: 0.127000; val_acc: 0.127000
(Iteration 11 / 200) loss: 2.068507
(Iteration 21 / 200) loss: 1.864477
(Iteration 31 / 200) loss: 1.767447
(Epoch 1 / 5) train acc: 0.367000; val_acc: 0.306000
(Iteration 41 / 200) loss: 1.867340
(Iteration 51 / 200) loss: 1.641590
(Iteration 61 / 200) loss: 1.753767
(Iteration 71 / 200) loss: 1.606004
(Epoch 2 / 5) train acc: 0.428000; val_acc: 0.317000
(Iteration 81 / 200) loss: 1.639500
(Iteration 91 / 200) loss: 1.738378
(Iteration 101 / 200) loss: 1.641308
(Iteration 111 / 200) loss: 1.653115
(Epoch 3 / 5) train acc: 0.459000; val_acc: 0.331000
(Iteration 121 / 200) loss: 1.467239
(Iteration 131 / 200) loss: 1.622463
(Iteration 141 / 200) loss: 1.664754
(Iteration 151 / 200) loss: 1.530565
(Epoch 4 / 5) train acc: 0.500000; val_acc: 0.350000
(Iteration 161 / 200) loss: 1.437343
(Iteration 171 / 200) loss: 1.551063
(Iteration 181 / 200) loss: 1.483104
(Iteration 191 / 200) loss: 1.389258
(Epoch 5 / 5) train acc: 0.515000; val_acc: 0.356000
```

```
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:30:
```

MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

```
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:33:
```

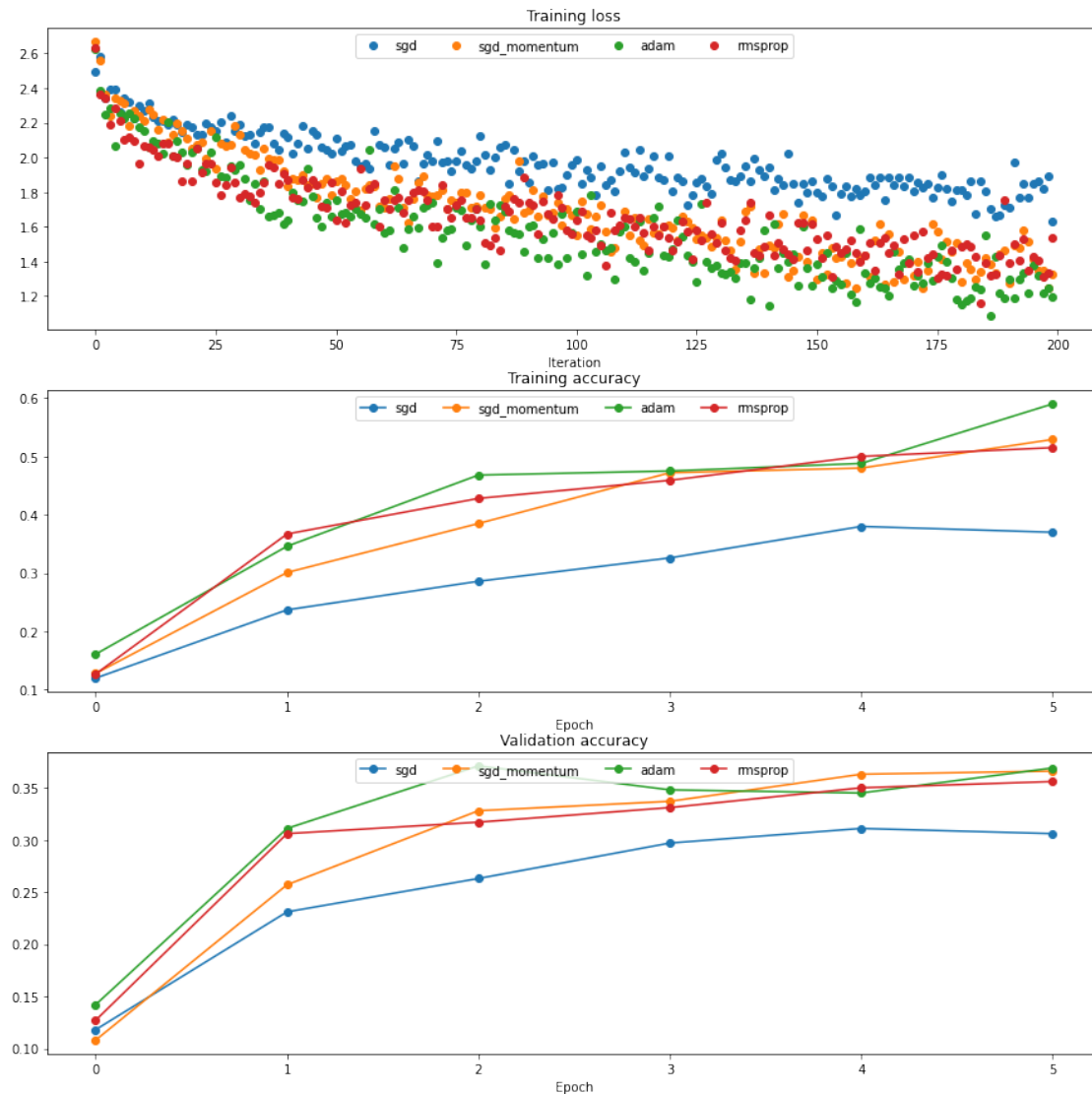
MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

```
/usr/local/lib/python3.6/dist-packages/ipykernel_launcher.py:36:
```

MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.

/usr/local/lib/python3.6/dist-packages/ipykernel\_launcher.py:40:

MatplotlibDeprecationWarning: Adding an axes using the same arguments as a previous axes currently reuses the earlier instance. In a future version, a new instance will always be created and returned. Meanwhile, this warning can be suppressed, and the future behavior ensured, by passing a unique label to each axes instance.



### 13.1 Inline Question 3:

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?

### 13.2 Answer:

- With AdaGrad, the updates tend to be very small because at every iteration, we are performing an update using the inverse of the accumulated value of squared gradients. For a large number of iterations, the cache value will be very large because we would have accumulated large positive values. Dividing the gradient by the square root of the cache value would lead to very small updates, and hence slow learning.
- On the other hand, Adam does not have this issue because it uses an exponentially weighted moving average (EWMA) of the gradients and their squares (inspired from RMSProp). In other words, it assigns larger weights to recent gradients while exponentially reducing the weights for the older gradients (older the datapoint, exponentially lower the weight). Adam computes the EWMA of the gradient and squared gradient (i.e., it utilizes the moment and second moment) rather than directly using the gradient and its square like AdaGrad. With Adam, the EWMA value of the squared gradient is much more limited in magnitude than the accumulation of the squared gradients. Specifically, Adam divides the EWMA of the gradients (first moment) by the square root of the EWMA of the squared gradients (second moment), instead of dividing the gradients with the square root of the accumulation of the gradient over all time-steps, as with AdaGrad. This helps Adam control large updates using the EWMA of the squared gradients.

## 14 Train a good model!

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

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

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.

```
[371]: 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)*****
```

```

# generate random combinations of hyperparameters given min-max ranges of
→hyperparams
# using a uniform probability distribution
# usage: lr, reg, hidden_dim, epochs = generate_random_hyperparams((-1, 0),
→(-7, -4), (10, 500), (10, 20))
def generate_random_hyperparams(lr, reg, hidden_size, epoch_values):
    lr = 10*np.random.uniform(lr[0], lr[1])
    reg = 10*np.random.uniform(reg[0], reg[1])
    hidden_size = np.random.randint(hidden_size[0], hidden_size[1])
    epochs = epoch_values[np.random.randint(0, len(epoch_values))]

    return lr, reg, hidden_size, epochs

# number of hypercombinations combinations to look for
num_hyperparam_configs = 10

# form random combinations of learning_rate, regularization_strength,
→hidden_layer_size and num_training_epochs
grid_search = [generate_random_hyperparams((-2, -4), (-7, -3), (50, 100, 200),
→(10, 20))
                for count in range(num_hyperparam_configs)]

# get our data
# data = get_CIFAR10_data()

best_val = -1    # The highest validation accuracy that we have seen so far.
results = {}

for config_num, config in enumerate(grid_search):
    print("Hyperparam config #{0} of #{0} -> ".format(config_num+1,
→len(grid_search)), end='')

    lr, reg, hidden_size, epochs = config
    print("lr: {:.2e}, reg: {:.2e}, hidden_size: {:.2e}, epochs: {:.2e}".
→format(lr, reg, hidden_size, epochs))

    model = TwoLayerNet(hidden_dim=hidden_size, reg=reg)
    current_solver = Solver(model, data, update_rule='sgd',
→optim_config={'learning_rate': lr},
                        lr_decay=0.95, num_epochs=epochs, batch_size=100,
                        print_every=100, verbose=False)

    # train a 2-layer neural net on the training set
    current_solver.train()

    # store the best validation accuracy and the TwoLayerNet object

```

```

if current_solver.best_val_acc > best_val:
    best_val = current_solver.best_val_acc
    best_model = model

# print results in-line to avoid saving them
print('Validation accuracy: %.4f' % (current_solver.best_val_acc,))
print() # new line after tending to every hyperparam config

print('Best validation accuracy achieved: %.4f' % best_val)

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

```

Hyperparam config #1 of #10 -> lr: 3.50e-04, reg: 5.29e-06, hidden\_size: 5.50e+01, epochs: 2.00e+01  
Validation accuracy: 0.5110

Hyperparam config #2 of #10 -> lr: 5.53e-04, reg: 5.89e-07, hidden\_size: 5.80e+01, epochs: 1.00e+01  
Validation accuracy: 0.5110

Hyperparam config #3 of #10 -> lr: 8.32e-04, reg: 2.17e-04, hidden\_size: 5.30e+01, epochs: 2.00e+01  
Validation accuracy: 0.5250

Hyperparam config #4 of #10 -> lr: 1.53e-03, reg: 1.80e-04, hidden\_size: 7.20e+01, epochs: 2.00e+01  
Validation accuracy: 0.5190

Hyperparam config #5 of #10 -> lr: 1.11e-04, reg: 1.60e-04, hidden\_size: 5.00e+01, epochs: 1.00e+01  
Validation accuracy: 0.4590

Hyperparam config #6 of #10 -> lr: 3.73e-03, reg: 7.12e-05, hidden\_size: 8.50e+01, epochs: 2.00e+01

```

/usr/local/lib/python3.6/dist-packages/numpy/core/fromnumeric.py:90:
RuntimeWarning: overflow encountered in reduce
    return ufunc.reduce(obj, axis, dtype, out, **passkwargs)
/content/drive/My Drive/cs231n/assignments/assignment2/cs231n/layers.py:1190:
RuntimeWarning: overflow encountered in subtract
    shifted_logits = x - np.max(x, axis=1, keepdims=True)
/content/drive/My Drive/cs231n/assignments/assignment2/cs231n/layers.py:1190:
RuntimeWarning: invalid value encountered in subtract
    shifted_logits = x - np.max(x, axis=1, keepdims=True)

```

```
/content/drive/My Drive/cs231n/assignments/assignment2/cs231n/layers.py:134:  
RuntimeWarning: invalid value encountered in less  
    dx[x < 0] = 0
```

Validation accuracy: 0.2040

Hyperparam config #7 of #10 -> lr: 6.64e-03, reg: 1.19e-05, hidden\_size:  
8.20e+01, epochs: 2.00e+01  
Validation accuracy: 0.1880

Hyperparam config #8 of #10 -> lr: 3.12e-03, reg: 2.70e-05, hidden\_size:  
5.20e+01, epochs: 1.00e+01  
Validation accuracy: 0.4530

Hyperparam config #9 of #10 -> lr: 1.32e-03, reg: 5.47e-05, hidden\_size:  
5.20e+01, epochs: 1.00e+01  
Validation accuracy: 0.5140

Hyperparam config #10 of #10 -> lr: 4.52e-03, reg: 6.36e-06, hidden\_size:  
6.40e+01, epochs: 2.00e+01  
Validation accuracy: 0.1490

Best validation accuracy achieved: 0.5250

## 15 Test your model!

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

```
[370]: 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: 0.517

Test set accuracy: 0.505

# BatchNormalization

May 10, 2020

```
[86]: # 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. 'cs231n/assignments/assignment3/'
FOLDERNAME = "cs231n/assignments/assignment2/"
assert FOLDERNAME is not None, "[!] Enter the foldername."

# symlink to make it easier to load your files
!ln -s "/content/drive/My Drive/$FOLDERNAME" "/content/assignment2"

# 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/cs231n/datasets/
!bash get_datasets.sh
%cd /content
```

Mounted at /content/drive

ln: failed to create symbolic link '/content/assignment2/assignment2': Operation not supported

/content/drive/My Drive/cs231n/assignments/assignment2/cs231n/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.](<https://arxiv.org/abs/1502.03167>)

```
[87]: # As usual, a bit of setup
import time
import numpy as np
import matplotlib.pyplot as plt
from cs231n.classifiers.fc_net import *
from cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient, \
    eval_numerical_gradient_array
from cs231n.solver import Solver

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

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

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

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

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

```
[88]: # 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 `cs231n/layers.py`, implement the batch normalization forward pass in the function `batchnorm_forward`. Once you have done so, run the following to test your implementation.

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

```
[89]: # 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:  [5.99520433e-17 6.93889390e-17 8.32667268e-19]
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]

```

[90]: *# 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.

```
[91]: # 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:  1.6674604875341426e-09
dgamma error:  7.417225040694815e-13
dbeta error:  2.379446949959628e-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  $\mu$  and variance  $v$ . With  $\mu$  and  $v$  calculated, we can calculate the standard deviation  $\sigma$  and normalized data  $Y$ . The equations and graph illustration below describe the computation ( $y_i$  is the  $i$ -th element of the vector  $Y$ ).

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

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

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.

```
[92]: 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: 1.2661271829168436e-12
dgamma difference: 1.3234426801152552e-14
dbeta difference: 0.0
speedup: 0.74x

```

## 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 cs231n/classifiers/fc\_net.py. Modify your implementation to add batch normalization.

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

HINT: You might find it useful to define an additional helper layer similar to those in the file cs231n/layer\_utils.py. If you decide to do so, do it in the file cs231n/classifiers/fc\_net.py.

```

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

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

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

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

```

```

Running check with reg = 0
Initial loss: 2.2611955101340957

```

```
W1 relative error: 1.10e-04
W2 relative error: 3.11e-06
W3 relative error: 4.05e-10
b1 relative error: 2.66e-07
b2 relative error: 2.55e-07
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: 2.41e-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: 1.91e-08
b2 relative error: 7.99e-07
b3 relative error: 1.42e-10
beta1 relative error: 6.65e-09
beta2 relative error: 3.48e-09
gamma1 relative error: 6.27e-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.

```
[94]: np.random.seed(231)
      # Try training a very deep net with batchnorm
      hidden_dims = [100, 100, 100, 100, 100]

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

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

```

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

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

```

Solver with batch norm:

```

(Iteration 1 / 200) loss: 2.340974
(Epoch 0 / 10) train acc: 0.107000; val_acc: 0.115000
(Epoch 1 / 10) train acc: 0.314000; val_acc: 0.266000
(Iteration 21 / 200) loss: 2.039365
(Epoch 2 / 10) train acc: 0.390000; val_acc: 0.279000
(Iteration 41 / 200) loss: 2.036710
(Epoch 3 / 10) train acc: 0.497000; val_acc: 0.316000
(Iteration 61 / 200) loss: 1.769764
(Epoch 4 / 10) train acc: 0.534000; val_acc: 0.308000
(Iteration 81 / 200) loss: 1.268882
(Epoch 5 / 10) train acc: 0.601000; val_acc: 0.320000
(Iteration 101 / 200) loss: 1.258349
(Epoch 6 / 10) train acc: 0.632000; val_acc: 0.317000
(Iteration 121 / 200) loss: 1.094234
(Epoch 7 / 10) train acc: 0.679000; val_acc: 0.332000
(Iteration 141 / 200) loss: 1.157111
(Epoch 8 / 10) train acc: 0.738000; val_acc: 0.324000
(Iteration 161 / 200) loss: 0.697783
(Epoch 9 / 10) train acc: 0.779000; val_acc: 0.338000
(Iteration 181 / 200) loss: 0.932644
(Epoch 10 / 10) train acc: 0.762000; val_acc: 0.321000

```

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: 1.696062
(Epoch 6 / 10) train acc: 0.536000; val_acc: 0.346000
(Iteration 121 / 200) loss: 1.550785
(Epoch 7 / 10) train acc: 0.530000; val_acc: 0.310000
(Iteration 141 / 200) loss: 1.436308
(Epoch 8 / 10) train acc: 0.622000; val_acc: 0.342000
(Iteration 161 / 200) loss: 1.000868
(Epoch 9 / 10) train acc: 0.654000; val_acc: 0.328000
(Iteration 181 / 200) loss: 0.925456
(Epoch 10 / 10) train acc: 0.726000; val_acc: 0.335000

```

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.

```

[95]: 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], \

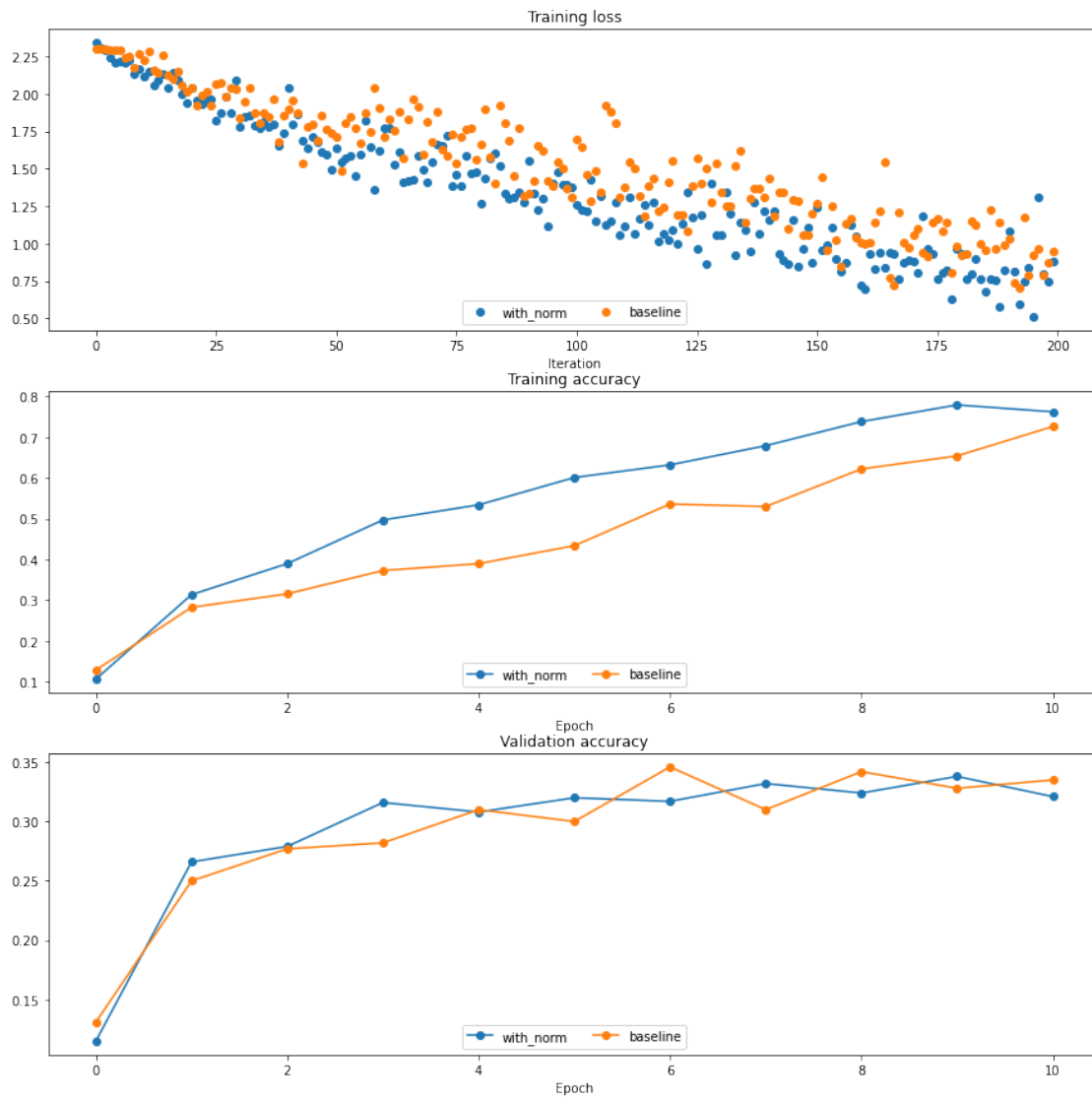
```

```

        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.

```
[96]: 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

```

```

[97]: # 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')

```

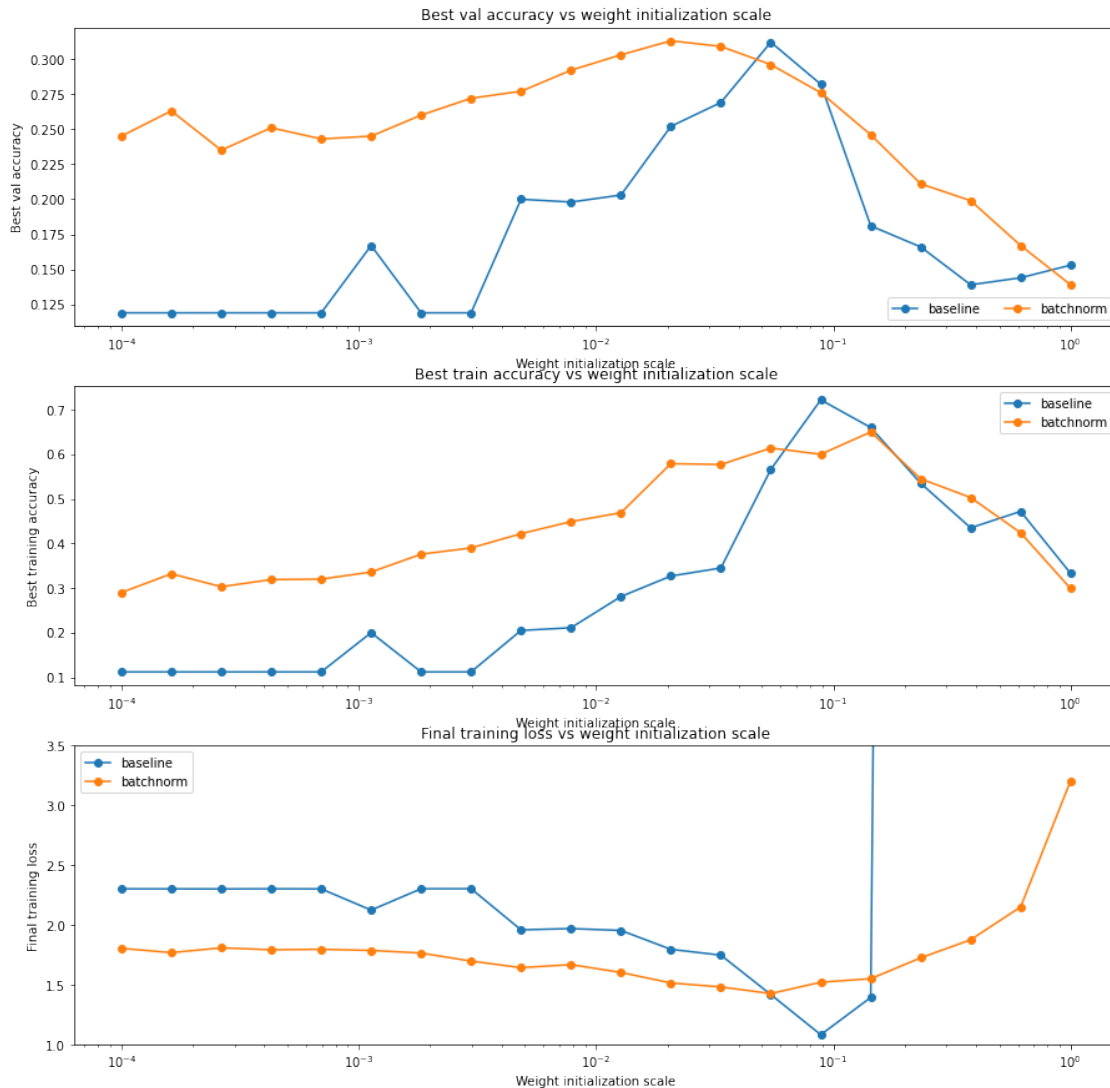
```

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:

- Batch normalization helps makes the model robust to poor initialization. The issue of vanishing gradients (small initial weights) can be clearly observed using the first and second plot. The baseline model is extremely sensitive to the weight initialization scale, and thus finding the right scale can be difficult. For this example, we can see that the baseline model shows poor performance between  $10^{-4}$  and  $10^{-3}$  and only obtains the best result when the weight scale is set to roughly  $2 \times 10^{-1}$ . On the other hand, we can see that the batchnorm model is much less sensitive to weight initialization because its accuracy is between 30% - 60% across the sweep of weight scales.
- The third plot depicts the problem of exploding gradients and it is glaringly evident in the baseline model for weight scale values greater than 1e-1. Batchnorm helps alleviate these issues and thus does not suffer from this problem.
- Batchnorm also helps with overfitting our model since it has regularization properties and thus, obtain better results than the baseline model.
- To summarize, with batch normalization we can avoid the problem of vanishing and exploding gradients because it normalizes every affine layer ( $Wx+b$ ) to have zero-centered data ( $\mu = 0$ ) with unit variance ( $\sigma = 1$ ), i.e., properties of a standard normal distribution, thus avoiding very large/small values. Moreover, Batchnorm's inherent regularization properties allow us to reduce overfitting.

## 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.

```
[98]: 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')

```

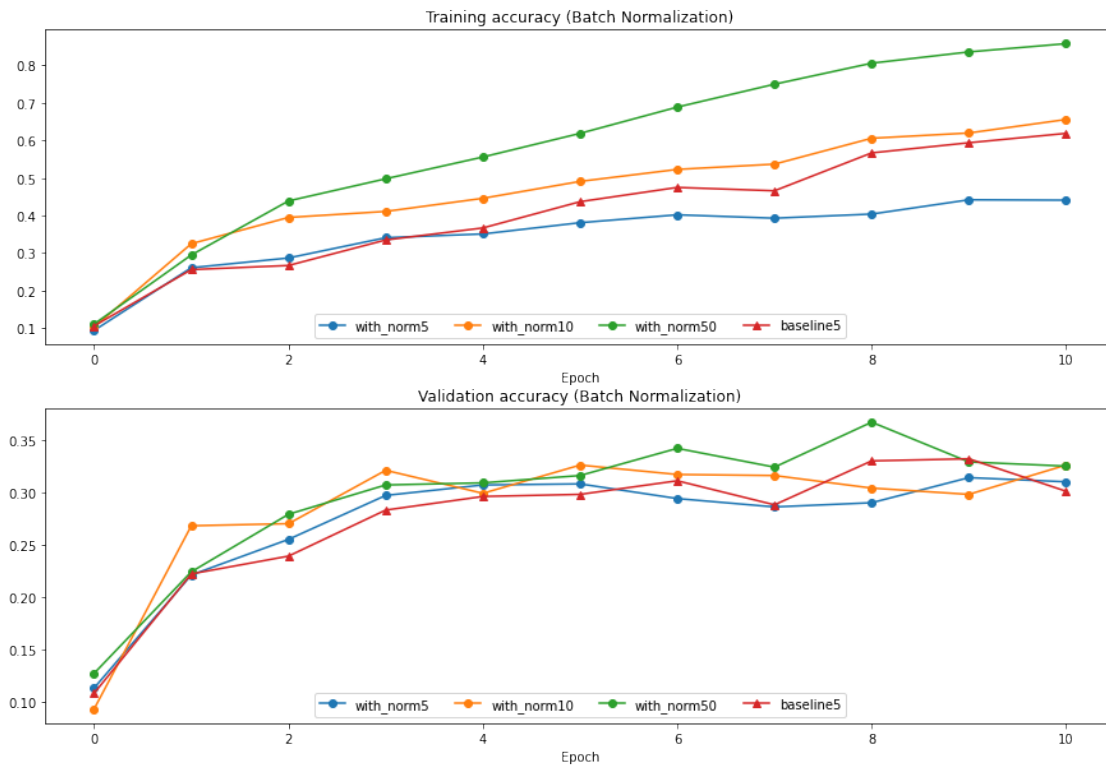
```

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

```

```
[99]: 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:

- According to the results, we can see that applying Batchnorm only makes sense when we have decently large batch sizes. The performance of batch normalization is a function of the

batch size. The smaller the batch size, worse the performance. Infact, for very small batch sizes ( $\leq 5$ ), the baseline model outperforms the model with Batchnorm.

- This issue occurs because our calculated batch statistics, i.e., the mean and variance, need to represent the batch as a whole. Thus, when we calculate these parameters, we are inherently trying to approximate the statistics of the entire dataset. With a small batch size, these statistics can be very noisy and thus, affect Batchnorm performance. On the other hand, with a large batch size we can obtain a much better approximation and thus better performance with Batchnorm.

## 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.](<https://arxiv.org/pdf/1607.06450.pdf>)

### 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:

2. This is analogous to layer normalization, since we are scaling over the dimensionality of the data (rather than over the training examples as is the case with batch normalization).

We can deduce this mathematically by considering  $\mu = 0$  and  $\beta = 0$ , i.e., zeroing out the shifting components since we are performing scaling here (and no shifting).

Per layernorm,

$$y = \gamma * x_{norm} + \beta$$

Since  $\mu = 0$  and  $\beta = 0$ ,

$$y = \gamma * x_{norm}$$

...where  $x_{norm} = \frac{x}{\sqrt{\sum x^2 + \epsilon}}$  and  $\gamma = \frac{x}{\sqrt{\sum x^2 + \epsilon}}$

Thus the result of layer normalization will be,

$$y = \frac{x^2}{\sum x^2 + \epsilon}$$

3. This is analogous to batch normalization since we are shifting over the training examples (rather than over the dimensionality of the data as is the case with layer normalization).

We can deduce this mathematically by considering  $\beta = 0$ ,  $\gamma = \sqrt{\sigma^2 + \epsilon}$  and mini-batch size = batch size (= size of the entire dataset).

Per batchnorm,

$$y = \gamma * x_{norm} + \beta$$

Since  $\beta = 0$ ,

$$y = \gamma * x_{norm}$$

...where  $x_{norm} = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}}$  and  $\gamma = \sqrt{\sigma^2 + \epsilon}$

Thus the result of layer normalization will be,

$$y = x - \mu$$

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

Run the cell below to check your results. \* In `cs231n/layers.py`, implement the backward pass for layer normalization in the function `layernorm_backward`.

Run the second cell below to check your results. \* Modify `cs231n/classifiers/fc_net.py` to add layer normalization to the `FullyConnectedNet`. When the normalization flag is set to "layernorm" in the constructor, you should insert a layer normalization layer before each ReLU nonlinearity.

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

```
[100]: # 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]

```

[101]: *# 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:  2.107277492956569e-09
dgamma error:  1.980045566295477e-12
dbeta error:  2.5842537629899423e-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!

```

[102]: 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)

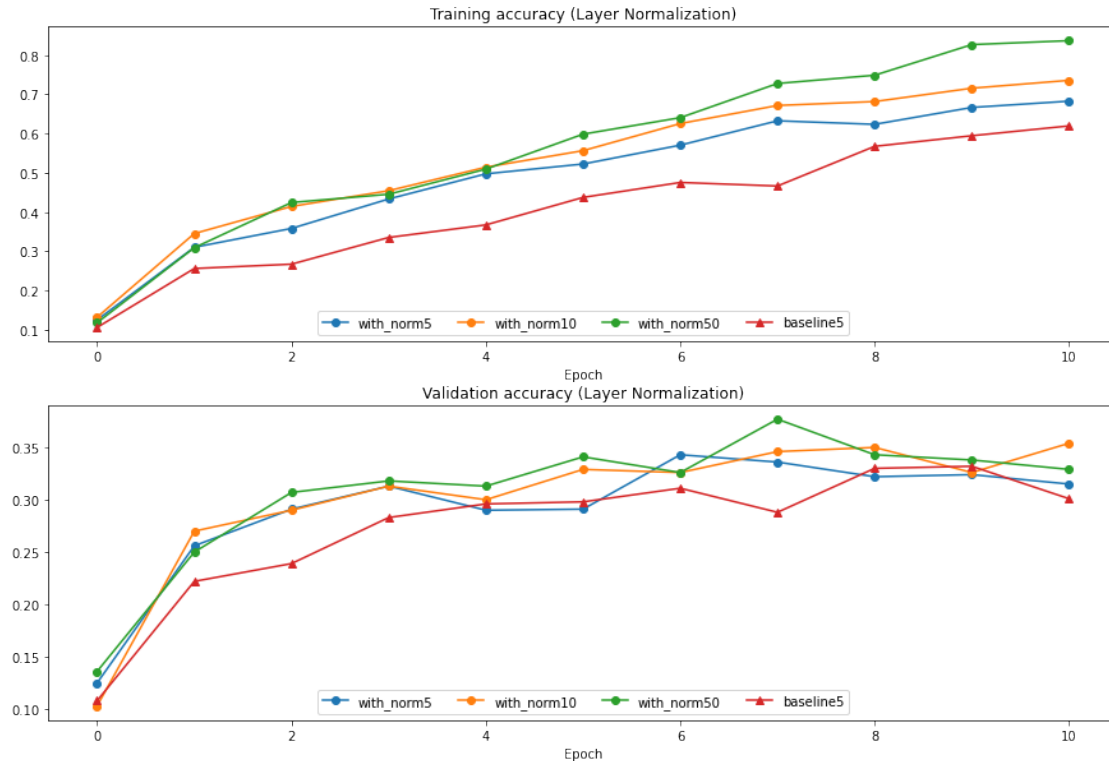
plt.gcf().set_size_inches(15, 10)
plt.show()

```

```

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. In the previous example, the network had five layers which can be considered a deep network. Using layer normalization in deep networks functions well. Also, in "[Layer Normalization](#)", Ba et al. offer examples of layernorm leading to better performance and faster training convergence with deep networks.
2. True. Having a small set of features affects the performance of layer normalization. The problem is analogous to that of batch normalization with a small batch size. In layer normalization, we calculate the statistics according to the number of hidden units, which represent the features that the network is learning, which in turn, is based on the dimensionality of data. Thus, smaller the feature spread, noisier the statistics computed during layer normalization.
3. True. Having a high regularization term affects the performance of layer normalization. In general, when the regularization term is very high, the model learns very simple functions

(underfitting). Code snippets showing the performance of layer normalization with different hidden sizes and regularization values can be seen in the next cells.

## 8 Layer normalization and hidden size

We will now run experiments varying the hidden layer size with layer normalization.

```
[103]: def run_hiddensize_experiments(normalization_mode):
    np.random.seed(231)

    # train network with different hidden layer sizes
    hidden_size = [5,10,30,50,70,100]
    solver_hidden_dims = [10, 10, 10, 10]
    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
    lr = 10**(-3.5)

    print('No normalization: hidden_sizes = ', hidden_size[0])
    model = FullyConnectedNet(solver_hidden_dims, weight_scale=weight_scale,
    ↪normalization=None)
    solver = Solver(model, small_data,
                    num_epochs=n_epochs, batch_size=50,
                    update_rule='adam',
                    optim_config={
                        'learning_rate': lr,
                    },
                    verbose=False)
    solver.train()

    bn_solvers = []
    for i in range(len(hidden_size)):
        print('Normalization: hidden sizes = ', hidden_size[i])
        hidden_dims = [ hidden_size[i] for j in range(4)]
        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=50,
                          update_rule='adam',
                          optim_config={
                              'learning_rate': lr,
```

```

        },
        verbose=False)
    bn_solver.train()
    bn_solvers.append(bn_solver)

    return bn_solvers, solver, hidden_size

# Run model
ln_solvers_hsize, solver_hsize, hidden_size = \
    →run_hiddensize_experiments('layernorm')

```

```

No normalization: hidden_sizes = 5
Normalization: hidden sizes = 5
Normalization: hidden sizes = 10
Normalization: hidden sizes = 30
Normalization: hidden sizes = 50
Normalization: hidden sizes = 70
Normalization: hidden sizes = 100

```

```

[104]: # Plot results
def plot_training_history2(title, label, baseline, bn_solvers, plot_fn, \
    →bl_marker='.', bn_marker='.', labels=None, label_prefix= '-'):
    """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=label_prefix
        if labels is not None:
            label += str(labels[i]) #str("%.4lf" %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, bbox_to_anchor=(0.5, -0.3))

plt.subplot(2, 1, 1)
plot_training_history2('Training accuracy (Layer Normalization)', 'Epoch', \
    →solver_hsize, ln_solvers_hsize, \
        lambda x: x.train_acc_history, bl_marker='-^', \
    →bn_marker='-o', labels=hidden_size, label_prefix='hsize')
plt.subplots_adjust(hspace = 0.5)
plt.subplot(2, 1, 2)

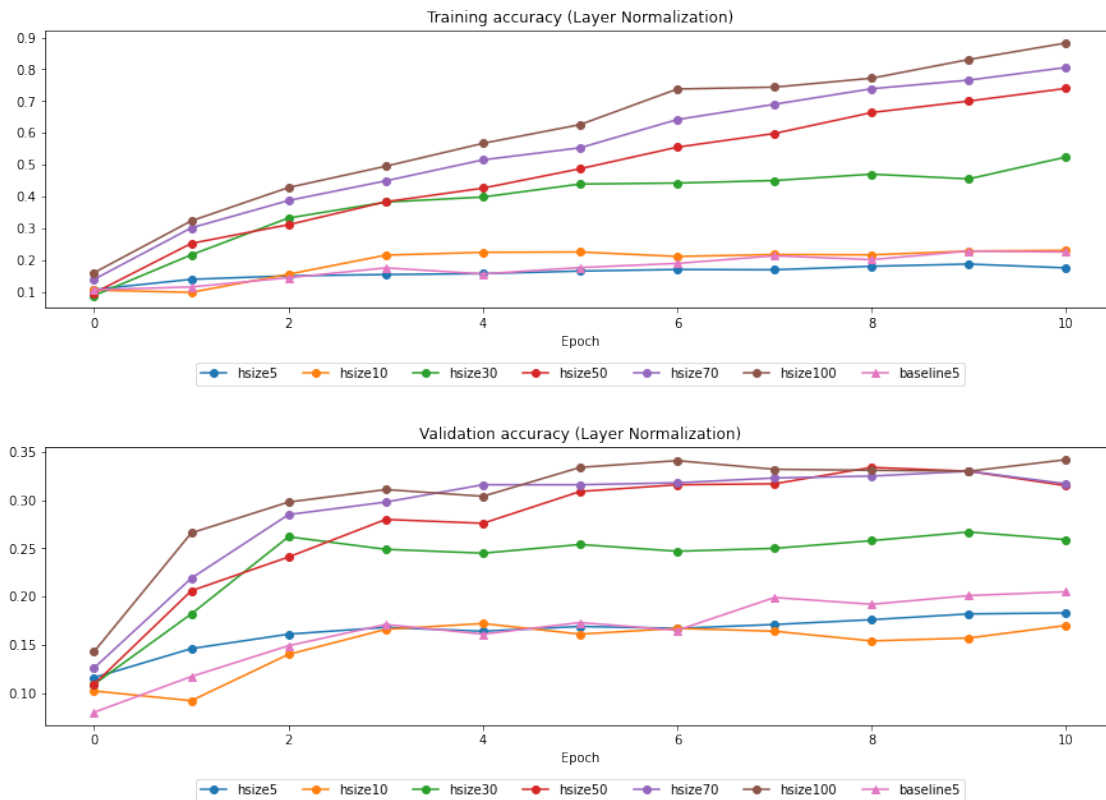
```

```

plot_training_history2('Validation accuracy (Layer Normalization)', 'Epoch', \
    →solver_hsize, ln_solvers_hsize, \
    lambda x: x.val_acc_history, bl_marker='~', \
    →bn_marker='-o', labels=hidden_size, label_prefix='hsize')

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

```



## 9 Layer Normalization and Regularization

We will now run experiments varying regularization values with layer normalization.

```

[105]: def run_regularization_experiments(normalization_mode):
    np.random.seed(231)

    # train a very deep net with regularization and layernorm
    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
    lr = 10**(-3.5)

    regularization = [0] # first try no regularization
    # vary regularization from 10-4 to 10-4
    regularization.extend(np.logspace(-4, 4, num=5))

    print('No normalization: regularization = ', regularization[0])
    model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale,
    ↪normalization=None, reg=regularization[0])
    solver = Solver(model, small_data,
                    num_epochs=n_epochs, batch_size=50,
                    update_rule='adam',
                    optim_config={
                        'learning_rate': lr,
                    },
                    verbose=False)
    solver.train()

    bn_solvers = []
    for i, reg in enumerate(regularization[1:]):
        print('Normalization: regularization = ', reg)
        bn_model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale,
    ↪normalization=normalization_mode, reg=reg)
        bn_solver = Solver(bn_model, small_data,
                        num_epochs=n_epochs, batch_size=50,
                        update_rule='adam',
                        optim_config={
                            'learning_rate': lr,
                        },
                        verbose=False)
        bn_solver.train()
        bn_solvers.append(bn_solver)

    return bn_solvers, solver, regularization

# Run model
ln_solvers_reg, solver_reg, regularization =
    ↪run_regularization_experiments('layernorm')

```

```

No normalization: regularization = 0
Normalization: regularization = 0.0001
Normalization: regularization = 0.01

```

```

Normalization: regularization = 1.0
Normalization: regularization = 100.0
Normalization: regularization = 10000.0

```

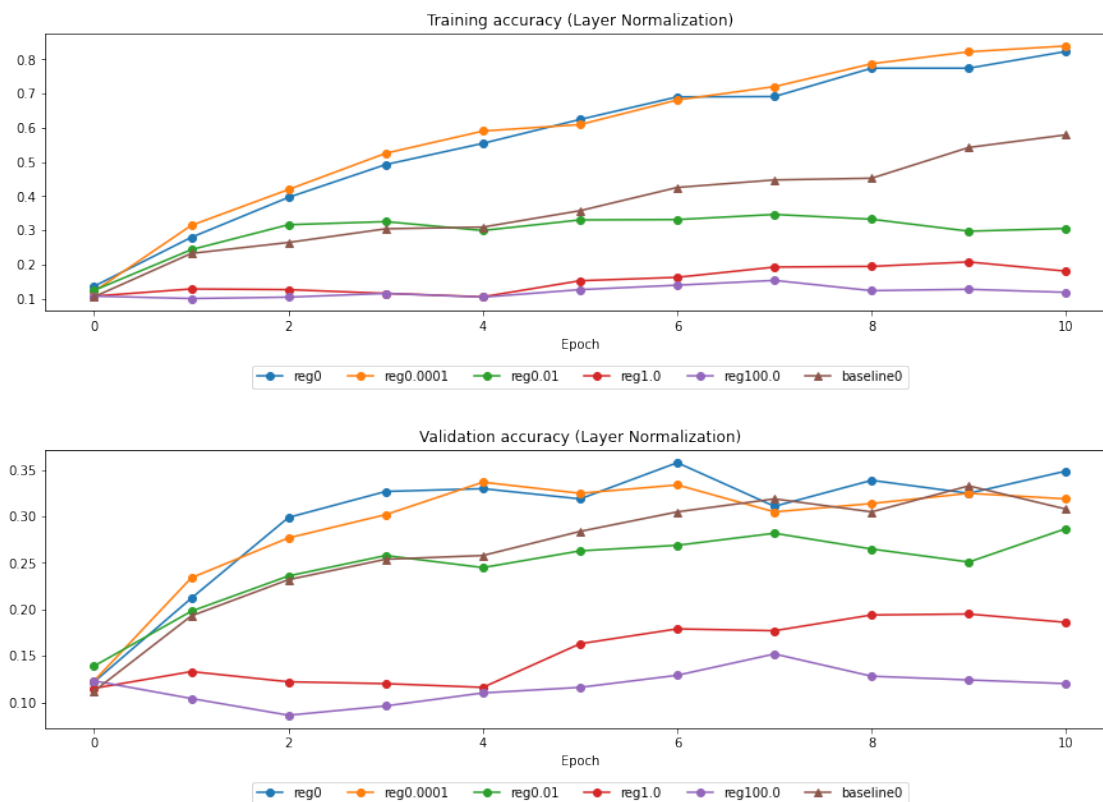
[106]: *# Plot results*

```

plt.subplot(2, 1, 1)
plot_training_history2('Training accuracy (Layer Normalization)', 'Epoch', \
    ↪ solver_reg, ln_solvers_reg, \
        lambda x: x.train_acc_history, bl_marker='-^', \
    ↪ bn_marker='-o', labels=regularization, \
        label_prefix='reg')
plt.subplots_adjust(hspace = 0.5)
plt.subplot(2, 1, 2)
plot_training_history2('Validation accuracy (Layer Normalization)', 'Epoch', \
    ↪ solver_reg, ln_solvers_reg, \
        lambda x: x.val_acc_history, bl_marker='-^', \
    ↪ bn_marker='-o', labels=regularization, \
        label_prefix='reg')

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

```



# Dropout

May 10, 2020

```
[14]: # 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. 'cs231n/assignments/assignment3/'
FOLDERNAME = "cs231n/assignments/assignment2/"
assert FOLDERNAME is not None, "[!] Enter the foldername."

# symlink to make it easier to load your files
!ln -s "/content/drive/My Drive/$FOLDERNAME" "/content/assignment2"

# 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/cs231n/datasets/
!bash get_datasets.sh
%cd /content
```

Mounted at /content/drive

```
ln: failed to create symbolic link '/content/assignment2/assignment2': Operation
not supported
/content/drive/My Drive/cs231n/assignments/assignment2/cs231n/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](<https://arxiv.org/abs/1207.0580>)

```
[15]: # As usual, a bit of setup
from __future__ import print_function
import time
import numpy as np
import matplotlib.pyplot as plt
from cs231n.classifiers.fc_net import *
from cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient, \
    eval_numerical_gradient_array
from cs231n.solver import Solver

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

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

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

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

```
[16]: # 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 `cs231n/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.

```
[17]: 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 `cs231n/layers.py`, implement the backward pass for dropout. After doing so, run the following cell to numerically gradient-check your implementation.

```
[18]: 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 by  $p$  (i.e., perform vanilla dropout instead of inverted dropout), then we would need to perform a scaling of the hidden layer outputs by  $p$  at test time. This is a crucial step because at test time all neurons see all their inputs (whereas with dropout, some neurons are randomly dropped during training), so we want the outputs of neurons at test time to be identical to their expected outputs at training time.
- This is necessary because at test time we will not be considering the (averaged) ensemble prediction of all the possible binary masks (and therefore all the exponentially many sub-networks). Without scaling the  $p$  at test time, we would only be considering the summation of all possible sub-networks which may lead to large values (and thus, exploding gradients).
- Since test-time performance is critical, it is always preferable to perform the scaling at train time and leave the forward pass at test time untouched, i.e., use inverted dropout. Additionally, this has the appealing property that the prediction code can remain untouched if you were to tweak the dropout (keep) probability or turn it off altogether.

## 4 Fully-connected nets with Dropout

In the file `cs231n/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.

```
[19]: 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.

```
[20]: # 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: 0.878000; val_acc: 0.269000
(Epoch 9 / 25) train acc: 0.902000; val_acc: 0.275000
(Epoch 10 / 25) train acc: 0.888000; val_acc: 0.261000
(Epoch 11 / 25) train acc: 0.926000; val_acc: 0.278000
(Epoch 12 / 25) train acc: 0.960000; val_acc: 0.302000
(Epoch 13 / 25) train acc: 0.964000; val_acc: 0.306000
```

```
(Epoch 14 / 25) train acc: 0.966000; val_acc: 0.309000
(Epoch 15 / 25) train acc: 0.976000; val_acc: 0.288000
(Epoch 16 / 25) train acc: 0.988000; val_acc: 0.301000
(Epoch 17 / 25) train acc: 0.988000; val_acc: 0.310000
(Epoch 18 / 25) train acc: 0.990000; val_acc: 0.311000
(Epoch 19 / 25) train acc: 0.990000; val_acc: 0.310000
(Epoch 20 / 25) train acc: 0.988000; val_acc: 0.312000
(Iteration 101 / 125) loss: 0.084611
(Epoch 21 / 25) train acc: 0.990000; val_acc: 0.302000
(Epoch 22 / 25) train acc: 0.978000; val_acc: 0.299000
(Epoch 23 / 25) train acc: 0.986000; val_acc: 0.291000
(Epoch 24 / 25) train acc: 0.994000; val_acc: 0.302000
(Epoch 25 / 25) train acc: 0.994000; val_acc: 0.293000
```

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: 0.297000
(Epoch 6 / 25) train acc: 0.626000; val_acc: 0.290000
(Epoch 7 / 25) train acc: 0.628000; val_acc: 0.298000
(Epoch 8 / 25) train acc: 0.686000; val_acc: 0.310000
(Epoch 9 / 25) train acc: 0.722000; val_acc: 0.289000
(Epoch 10 / 25) train acc: 0.724000; val_acc: 0.300000
(Epoch 11 / 25) train acc: 0.760000; val_acc: 0.305000
(Epoch 12 / 25) train acc: 0.772000; val_acc: 0.278000
(Epoch 13 / 25) train acc: 0.818000; val_acc: 0.306000
(Epoch 14 / 25) train acc: 0.816000; val_acc: 0.339000
(Epoch 15 / 25) train acc: 0.854000; val_acc: 0.351000
(Epoch 16 / 25) train acc: 0.832000; val_acc: 0.296000
(Epoch 17 / 25) train acc: 0.854000; val_acc: 0.288000
(Epoch 18 / 25) train acc: 0.846000; val_acc: 0.320000
(Epoch 19 / 25) train acc: 0.872000; val_acc: 0.344000
(Epoch 20 / 25) train acc: 0.868000; val_acc: 0.305000
(Iteration 101 / 125) loss: 5.472347
(Epoch 21 / 25) train acc: 0.866000; val_acc: 0.329000
(Epoch 22 / 25) train acc: 0.902000; val_acc: 0.308000
(Epoch 23 / 25) train acc: 0.898000; val_acc: 0.313000
(Epoch 24 / 25) train acc: 0.912000; val_acc: 0.334000
(Epoch 25 / 25) train acc: 0.914000; val_acc: 0.322000
```

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

```

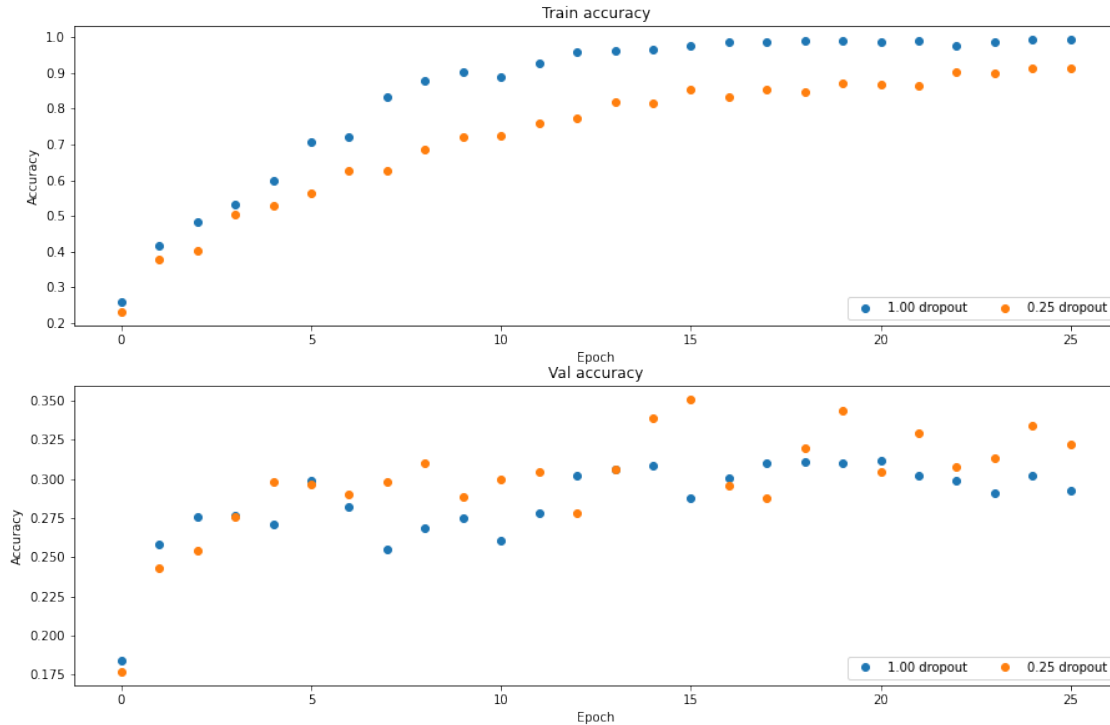
train_accs = []
val_accs = []
for dropout in dropout_choices:
    solver = solvers[dropout]
    train_accs.append(solver.train_acc_history[-1])
    val_accs.append(solver.val_acc_history[-1])

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 results show that the model is overfitting the training data.
- Without Dropout, there is a stark delta between the training (99%, at epoch 25) and validation accuracy (30%, at epoch 25), indicating overfitting. Dropout seeks to bridge the gap between the training accuracy and validation accuracy (that arises due to overfitting) by regularizing the model, which leads to an improvement in the validation accuracy and a slight degradation in training performance.
- With Dropout, the training accuracy reduces a bit (93%, at epoch 25) but the validation accuracy improves (32.5%, at epoch 25). Thus, Dropout helped the model's validation performance which would likely help the generalization capability of the model with unseen data. This suggests that with dropout, we are regularizing model complexity and learning a simpler model, effectively reducing overfitting.

### 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:

- Decreasing  $p$  would effectively impact the size of the hidden layer, since dropout leads to a reduced set of neurons being active. Thus, we do not need to modify the keep probability  $p$  when we vary the size of the hidden layers.
- Another way to interpret this is as follows. If we decide to decrease the size of the hidden layers, we are not required to modify  $p$  because the number of neurons, which will be dropped out will essentially be proportional to the size of the respective hidden layer. As an example, let's suppose we have  $n = 1024$  neurons in a hidden layer and we are using  $p = 0.5$ . The expected number of dropped neurons would be  $p * n = 0.5 * 1024 = 512$  which is essentially emulating a hidden layer of size 512.

# Convolutional Networks

May 10, 2020

```
[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. 'cs231n/assignments/assignment3/'
FOLDERNAME = "cs231n/assignments/assignment2/"
assert FOLDERNAME is not None, "[!] Enter the foldername."

# symlink to make it easier to load your files
!ln -s "/content/drive/My Drive/$FOLDERNAME" "/content/assignment2"

# 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/cs231n/datasets/
!bash get_datasets.sh
%cd /content
```

Mounted at /content/drive

ln: failed to create symbolic link '/content/assignment2/assignment2': Operation not supported

/content/drive/My Drive/cs231n/assignments/assignment2/cs231n/datasets  
/content

## 1 Convolutional Networks

So far we have worked with deep fully-connected networks, using them to explore different optimization strategies and network architectures. Fully-connected networks are a good testbed for experimentation because they are very computationally efficient, but in practice all state-of-the-art results use convolutional networks instead.

First you will implement several layer types that are used in convolutional networks. You will then use these layers to train a convolutional network on the CIFAR-10 dataset.

```
[0]: # As usual, a bit of setup
import numpy as np
import matplotlib.pyplot as plt
from cs231n.classifiers.cnn import *
from cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient_array,
    →eval_numerical_gradient
from cs231n.layers import *
from cs231n.fast_layers import *
from cs231n.solver import Solver

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

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

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

[0]: # 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 Convolution: Naive forward pass

The core of a convolutional network is the convolution operation. In the file `cs231n/layers.py`, implement the forward pass for the convolution layer in the function `conv_forward_naive`.

You don't have to worry too much about efficiency at this point; just write the code in whatever way you find most clear.

You can test your implementation by running the following:

```
[0]: x_shape = (2, 3, 4, 4)
w_shape = (3, 3, 4, 4)
x = np.linspace(-0.1, 0.5, num=np.prod(x_shape)).reshape(x_shape)
w = np.linspace(-0.2, 0.3, num=np.prod(w_shape)).reshape(w_shape)
b = np.linspace(-0.1, 0.2, num=3)

conv_param = {'stride': 2, 'pad': 1}
out, _ = conv_forward_naive(x, w, b, conv_param)
correct_out = np.array([[[[-0.08759809, -0.10987781],
                           [-0.18387192, -0.2109216 ]],
                          [[ 0.21027089,  0.21661097],
                           [ 0.22847626,  0.23004637]],
                          [[ 0.50813986,  0.54309974],
                           [ 0.64082444,  0.67101435]]],
                        [[[-0.98053589, -1.03143541],
                           [-1.19128892, -1.24695841]],
                          [[ 0.69108355,  0.66880383],
                           [ 0.59480972,  0.56776003]],
                          [[ 2.36270298,  2.36904306],
                           [ 2.38090835,  2.38247847]]]])

# Compare your output to ours; difference should be around e-8
print('Testing conv_forward_naive')
print('difference: ', rel_error(out, correct_out))
```

```
Testing conv_forward_naive
difference:  2.2121476417505994e-08
```

### 3 Aside: Image processing via convolutions

As fun way to both check your implementation and gain a better understanding of the type of operation that convolutional layers can perform, we will set up an input containing two images and manually set up filters that perform common image processing operations (grayscale conversion and edge detection). The convolution forward pass will apply these operations to each of the input images. We can then visualize the results as a sanity check.

#### 3.1 Colab Users Only

Please execute the below cell to copy two cat images to the Colab VM.

```
[0]: # Colab users only!
%mkdir -p cs231n/notebook_images
%cd drive/My\ Drive/$FOLDERNAME/cs231n
%cp -r notebook_images/ /content/cs231n/
%cd /content/
```

```
/content/drive/My Drive/cs231n/assignments/assignment2/cs231n
/content
```

```
[0]: from imageio import imread
      from PIL import Image

      kitten = imread('cs231n/notebook_images/kitten.jpg')
      puppy = imread('cs231n/notebook_images/puppy.jpg')
      # kitten is wide, and puppy is already square
      d = kitten.shape[1] - kitten.shape[0]
      kitten_cropped = kitten[:, d//2:-d//2, :]

      img_size = 200 # Make this smaller if it runs too slow
      resized_puppy = np.array(Image.fromarray(puppy).resize((img_size, img_size)))
      resized_kitten = np.array(Image.fromarray(kitten_cropped).resize((img_size,
      →img_size)))
      x = np.zeros((2, 3, img_size, img_size))
      x[0, :, :, :] = resized_puppy.transpose((2, 0, 1))
      x[1, :, :, :] = resized_kitten.transpose((2, 0, 1))

      # Set up a convolutional weights holding 2 filters, each 3x3
      w = np.zeros((2, 3, 3, 3))

      # The first filter converts the image to grayscale.
      # Set up the red, green, and blue channels of the filter.
      w[0, 0, :, :] = [[0, 0, 0], [0, 0.3, 0], [0, 0, 0]]
      w[0, 1, :, :] = [[0, 0, 0], [0, 0.6, 0], [0, 0, 0]]
      w[0, 2, :, :] = [[0, 0, 0], [0, 0.1, 0], [0, 0, 0]]

      # Second filter detects horizontal edges in the blue channel.
      w[1, 2, :, :] = [[1, 2, 1], [0, 0, 0], [-1, -2, -1]]

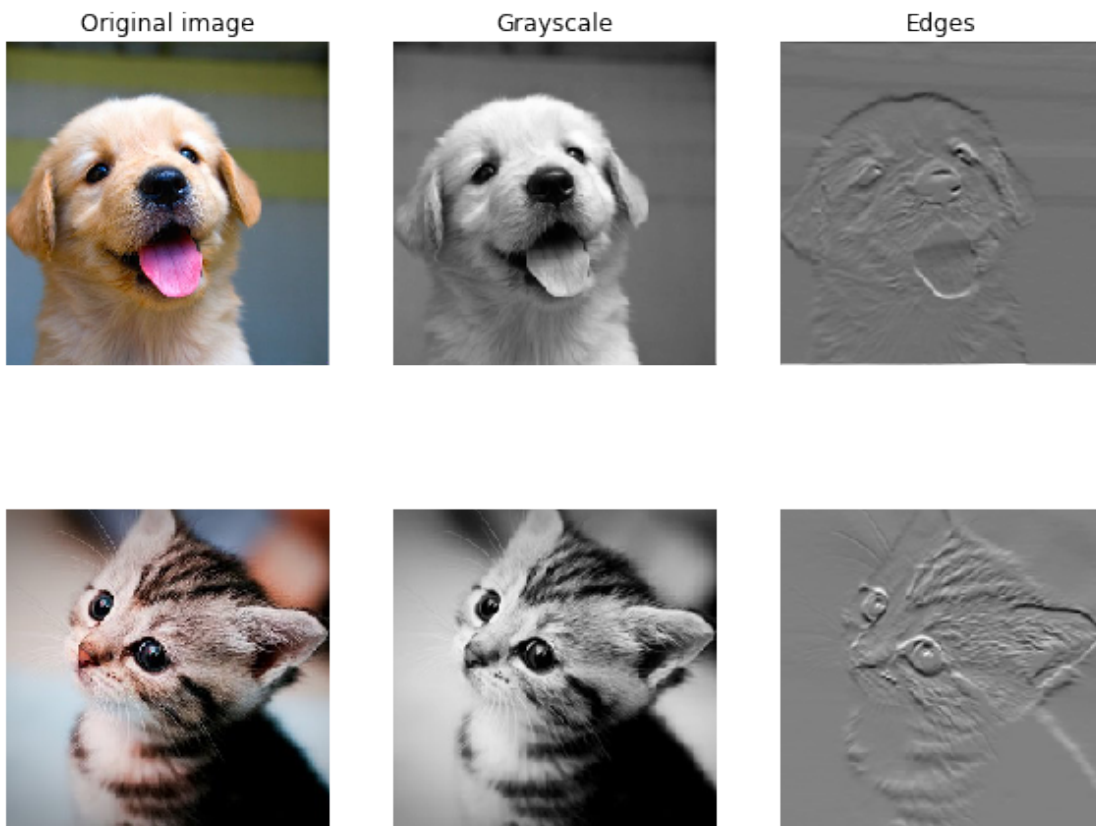
      # Vector of biases. We don't need any bias for the grayscale
      # filter, but for the edge detection filter we want to add 128
      # to each output so that nothing is negative.
      b = np.array([0, 128])

      # Compute the result of convolving each input in x with each filter in w,
      # offsetting by b, and storing the results in out.
      out, _ = conv_forward_naive(x, w, b, {'stride': 1, 'pad': 1})

      def imshow_no_ax(img, normalize=True):
          """ Tiny helper to show images as uint8 and remove axis labels """
          if normalize:
              img_max, img_min = np.max(img), np.min(img)
              img = 255.0 * (img - img_min) / (img_max - img_min)
          plt.imshow(img.astype('uint8'))
```

```
plt.gca().axis('off')

# Show the original images and the results of the conv operation
plt.subplot(2, 3, 1)
imshow_no_ax(puppy, normalize=False)
plt.title('Original image')
plt.subplot(2, 3, 2)
imshow_no_ax(out[0, 0])
plt.title('Grayscale')
plt.subplot(2, 3, 3)
imshow_no_ax(out[0, 1])
plt.title('Edges')
plt.subplot(2, 3, 4)
imshow_no_ax(kitten_cropped, normalize=False)
plt.subplot(2, 3, 5)
imshow_no_ax(out[1, 0])
plt.subplot(2, 3, 6)
imshow_no_ax(out[1, 1])
plt.show()
```



## 4 Convolution: Naive backward pass

Implement the backward pass for the convolution operation in the function `conv_backward_naive` in the file `cs231n/layers.py`. Again, you don't need to worry too much about computational efficiency.

When you are done, run the following to check your backward pass with a numeric gradient check.

```
[0]: np.random.seed(231)
x = np.random.randn(4, 3, 5, 5)
w = np.random.randn(2, 3, 3, 3)
b = np.random.randn(2,)
dout = np.random.randn(4, 2, 5, 5)
conv_param = {'stride': 1, 'pad': 1}

dx_num = eval_numerical_gradient_array(lambda x: conv_forward_naive(x, w, b,
    ↪conv_param)[0], x, dout)
dw_num = eval_numerical_gradient_array(lambda w: conv_forward_naive(x, w, b,
    ↪conv_param)[0], w, dout)
db_num = eval_numerical_gradient_array(lambda b: conv_forward_naive(x, w, b,
    ↪conv_param)[0], b, dout)

out, cache = conv_forward_naive(x, w, b, conv_param)
dx, dw, db = conv_backward_naive(dout, cache)

# Your errors should be around e-8 or less.
print('Testing conv_backward_naive function')
print('dx error: ', rel_error(dx, dx_num))
print('dw error: ', rel_error(dw, dw_num))
print('db error: ', rel_error(db, db_num))
```

```
Testing conv_backward_naive function
dx error:  1.159803161159293e-08
dw error:  2.2471264748452487e-10
db error:  3.37264006649648e-11
```

## 5 Max-Pooling: Naive forward

Implement the forward pass for the max-pooling operation in the function `max_pool_forward_naive` in the file `cs231n/layers.py`. Again, don't worry too much about computational efficiency.

Check your implementation by running the following:

```
[0]: x_shape = (2, 3, 4, 4)
x = np.linspace(-0.3, 0.4, num=np.prod(x_shape)).reshape(x_shape)
pool_param = {'pool_width': 2, 'pool_height': 2, 'stride': 2}

out, _ = max_pool_forward_naive(x, pool_param)
```

```

correct_out = np.array([[[[-0.26315789, -0.24842105],
                           [-0.20421053, -0.18947368]],
                          [[-0.14526316, -0.13052632],
                           [-0.08631579, -0.07157895]],
                          [[-0.02736842, -0.01263158],
                           [ 0.03157895,  0.04631579]]],
                        [[[ 0.09052632,  0.10526316],
                           [ 0.14947368,  0.16421053]],
                          [[ 0.20842105,  0.22315789],
                           [ 0.26736842,  0.28210526]],
                          [[ 0.32631579,  0.34105263],
                           [ 0.38526316,  0.4          ]]])

# Compare your output with ours. Difference should be on the order of e-8.
print('Testing max_pool_forward_naive function:')
print('difference: ', rel_error(out, correct_out))

```

Testing max\_pool\_forward\_naive function:  
difference: 4.1666665157267834e-08

## 6 Max-Pooling: Naive backward

Implement the backward pass for the max-pooling operation in the function `max_pool_backward_naive` in the file `cs231n/layers.py`. You don't need to worry about computational efficiency.

Check your implementation with numeric gradient checking by running the following:

```

[0]: np.random.seed(231)
x = np.random.randn(3, 2, 8, 8)
dout = np.random.randn(3, 2, 4, 4)
pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}

dx_num = eval_numerical_gradient_array(lambda x: max_pool_forward_naive(x,
    ↪pool_param)[0], x, dout)

out, cache = max_pool_forward_naive(x, pool_param)
dx = max_pool_backward_naive(dout, cache)

# Your error should be on the order of e-12
print('Testing max_pool_backward_naive function:')
print('dx error: ', rel_error(dx, dx_num))

```

Testing max\_pool\_backward\_naive function:  
dx error: 3.27562514223145e-12



## 7 Fast layers

Making convolution and pooling layers fast can be challenging. To spare you the pain, we've provided fast implementations of the forward and backward passes for convolution and pooling layers in the file `cs231n/fast_layers.py`.

The fast convolution implementation depends on a Cython extension; to compile it either execute the local development cell (option A) if you are developing locally, or the Colab cell (option B) if you are running this assignment in Colab.

---

**Very Important, Please Read.** For **both** option A and B, you have to **restart** the notebook after compiling the cython extension. In Colab, please save the notebook File -> Save, then click Runtime -> Restart Runtime -> Yes. This will restart the kernel which means local variables will be lost. Just re-execute the cells from top to bottom and skip the cell below as you only need to run it once for the compilation step.

---

### 7.1 Option A: Local Development

Go to the `cs231n` directory and execute the following in your terminal:

```
python setup.py build_ext --inplace
```

### 7.2 Option B: Colab

Execute the cell below only only **ONCE**.

```
[0]: %cd drive/My\ Drive/$FOLDERNAME/cs231n/  
!python setup.py build_ext --inplace
```

```
/content/drive/My Drive/cs231n/assignments/assignment2/cs231n  
running build_ext
```

The API for the fast versions of the convolution and pooling layers is exactly the same as the naive versions that you implemented above: the forward pass receives data, weights, and parameters and produces outputs and a cache object; the backward pass receives upstream derivatives and the cache object and produces gradients with respect to the data and weights.

**NOTE:** The fast implementation for pooling will only perform optimally if the pooling regions are non-overlapping and tile the input. If these conditions are not met then the fast pooling implementation will not be much faster than the naive implementation.

You can compare the performance of the naive and fast versions of these layers by running the following:

```
[0]: # Rel errors should be around e-9 or less  
from cs231n.fast_layers import conv_forward_fast, conv_backward_fast  
from time import time  
np.random.seed(231)  
x = np.random.randn(100, 3, 31, 31)
```

```

w = np.random.randn(25, 3, 3, 3)
b = np.random.randn(25,)
dout = np.random.randn(100, 25, 16, 16)
conv_param = {'stride': 2, 'pad': 1}

t0 = time()
out_naive, cache_naive = conv_forward_naive(x, w, b, conv_param)
t1 = time()
out_fast, cache_fast = conv_forward_fast(x, w, b, conv_param)
t2 = time()

print('Testing conv_forward_fast:')
print('Naive: %fs' % (t1 - t0))
print('Fast: %fs' % (t2 - t1))
print('Speedup: %fx' % ((t1 - t0) / (t2 - t1)))
print('Difference: ', rel_error(out_naive, out_fast))

t0 = time()
dx_naive, dw_naive, db_naive = conv_backward_naive(dout, cache_naive)
t1 = time()
dx_fast, dw_fast, db_fast = conv_backward_fast(dout, cache_fast)
t2 = time()

print('\nTesting conv_backward_fast:')
print('Naive: %fs' % (t1 - t0))
print('Fast: %fs' % (t2 - t1))
print('Speedup: %fx' % ((t1 - t0) / (t2 - t1)))
print('dx difference: ', rel_error(dx_naive, dx_fast))
print('dw difference: ', rel_error(dw_naive, dw_fast))
print('db difference: ', rel_error(db_naive, db_fast))

```

```

Testing conv_forward_fast:
Naive: 4.522038s
Fast: 0.008794s
Speedup: 514.242361x
Difference: 4.926407851494105e-11

```

```

Testing conv_backward_fast:
Naive: 7.278914s
Fast: 0.011704s
Speedup: 621.905828x
dx difference: 1.949764775345631e-11
dw difference: 3.681156828004736e-13
db difference: 0.0

```

```

[0]: # Relative errors should be close to 0.0
from cs231n.fast_layers import max_pool_forward_fast, max_pool_backward_fast

```

```

np.random.seed(231)
x = np.random.randn(100, 3, 32, 32)
dout = np.random.randn(100, 3, 16, 16)
pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}

t0 = time()
out_naive, cache_naive = max_pool_forward_naive(x, pool_param)
t1 = time()
out_fast, cache_fast = max_pool_forward_fast(x, pool_param)
t2 = time()

print('Testing pool_forward_fast:')
print('Naive: %fs' % (t1 - t0))
print('fast: %fs' % (t2 - t1))
print('speedup: %fx' % ((t1 - t0) / (t2 - t1)))
print('difference: ', rel_error(out_naive, out_fast))

t0 = time()
dx_naive = max_pool_backward_naive(dout, cache_naive)
t1 = time()
dx_fast = max_pool_backward_fast(dout, cache_fast)
t2 = time()

print('\nTesting pool_backward_fast:')
print('Naive: %fs' % (t1 - t0))
print('fast: %fs' % (t2 - t1))
print('speedup: %fx' % ((t1 - t0) / (t2 - t1)))
print('dx difference: ', rel_error(dx_naive, dx_fast))

```

```

Testing pool_forward_fast:
Naive: 0.192160s
fast: 0.003110s
speedup: 61.779703x
difference: 0.0

```

```

Testing pool_backward_fast:
Naive: 0.683393s
fast: 0.012554s
speedup: 54.434491x
dx difference: 0.0

```

## 8 Convolutional "sandwich" layers

Previously we introduced the concept of "sandwich" layers that combine multiple operations into commonly used patterns. In the file `cs231n/layer_utils.py` you will find sandwich layers that implement a few commonly used patterns for convolutional networks. Run the cells below to sanity check they're working.

```
[0]: from cs231n.layer_utils import conv_relu_pool_forward, conv_relu_pool_backward
np.random.seed(231)
x = np.random.randn(2, 3, 16, 16)
w = np.random.randn(3, 3, 3, 3)
b = np.random.randn(3,)
dout = np.random.randn(2, 3, 8, 8)
conv_param = {'stride': 1, 'pad': 1}
pool_param = {'pool_height': 2, 'pool_width': 2, 'stride': 2}

out, cache = conv_relu_pool_forward(x, w, b, conv_param, pool_param)
dx, dw, db = conv_relu_pool_backward(dout, cache)

dx_num = eval_numerical_gradient_array(lambda x: conv_relu_pool_forward(x, w,
    ↪b, conv_param, pool_param)[0], x, dout)
dw_num = eval_numerical_gradient_array(lambda w: conv_relu_pool_forward(x, w,
    ↪b, conv_param, pool_param)[0], w, dout)
db_num = eval_numerical_gradient_array(lambda b: conv_relu_pool_forward(x, w,
    ↪b, conv_param, pool_param)[0], b, dout)

# Relative errors should be around e-8 or less
print('Testing conv_relu_pool')
print('dx error: ', rel_error(dx_num, dx))
print('dw error: ', rel_error(dw_num, dw))
print('db error: ', rel_error(db_num, db))
```

```
Testing conv_relu_pool
dx error: 9.591132621921372e-09
dw error: 5.802391137330214e-09
db error: 1.0146343411762047e-09
```

```
[0]: from cs231n.layer_utils import conv_relu_forward, conv_relu_backward
np.random.seed(231)
x = np.random.randn(2, 3, 8, 8)
w = np.random.randn(3, 3, 3, 3)
b = np.random.randn(3,)
dout = np.random.randn(2, 3, 8, 8)
conv_param = {'stride': 1, 'pad': 1}

out, cache = conv_relu_forward(x, w, b, conv_param)
dx, dw, db = conv_relu_backward(dout, cache)

dx_num = eval_numerical_gradient_array(lambda x: conv_relu_forward(x, w, b,
    ↪conv_param)[0], x, dout)
dw_num = eval_numerical_gradient_array(lambda w: conv_relu_forward(x, w, b,
    ↪conv_param)[0], w, dout)
db_num = eval_numerical_gradient_array(lambda b: conv_relu_forward(x, w, b,
    ↪conv_param)[0], b, dout)
```

```
# Relative errors should be around e-8 or less
print('Testing conv_relu:')
print('dx error: ', rel_error(dx_num, dx))
print('dw error: ', rel_error(dw_num, dw))
print('db error: ', rel_error(db_num, db))
```

```
Testing conv_relu:
dx error:  1.5218619980349303e-09
dw error:  2.702022646099404e-10
db error:  1.451272393591721e-10
```

## 9 Three-layer ConvNet

Now that you have implemented all the necessary layers, we can put them together into a simple convolutional network.

Open the file `cs231n/classifiers/cnn.py` and complete the implementation of the `ThreeLayerConvNet` class. Remember you can use the fast/sandwich layers (already imported for you) in your implementation. Run the following cells to help you debug:

### 9.1 Sanity check loss

After you build a new network, one of the first things you should do is sanity check the loss. When we use the softmax loss, we expect the loss for random weights (and no regularization) to be about  $\log(C)$  for  $C$  classes. When we add regularization the loss should go up slightly.

```
[0]: model = ThreeLayerConvNet()

N = 50
X = np.random.randn(N, 3, 32, 32)
y = np.random.randint(10, size=N)

loss, grads = model.loss(X, y)
print('Initial loss (no regularization): ', loss)

model.reg = 0.5
loss, grads = model.loss(X, y)
print('Initial loss (with regularization): ', loss)
```

```
Initial loss (no regularization):  2.302586071243987
Initial loss (with regularization):  2.508255635671795
```

### 9.2 Gradient check

After the loss looks reasonable, use numeric gradient checking to make sure that your backward pass is correct. When you use numeric gradient checking you should use a small amount of artificial data and a small number of neurons at each layer. Note: correct implementations may still have relative errors up to the order of  $e^{-2}$ .

```
[0]: num_inputs = 2
input_dim = (3, 16, 16)
reg = 0.0
num_classes = 10
np.random.seed(231)
X = np.random.randn(num_inputs, *input_dim)
y = np.random.randint(num_classes, size=num_inputs)

model = ThreeLayerConvNet(num_filters=3, filter_size=3,
                           input_dim=input_dim, hidden_dim=7,
                           dtype=np.float64)
loss, grads = model.loss(X, y)
# Errors should be small, but correct implementations may have
# relative errors up to the order of e-2
for param_name in sorted(grads):
    f = lambda _: model.loss(X, y)[0]
    param_grad_num = eval_numerical_gradient(f, model.params[param_name],
    ↪ verbose=False, h=1e-6)
    e = rel_error(param_grad_num, grads[param_name])
    print('%s max relative error: %e' % (param_name, rel_error(param_grad_num,
    ↪ grads[param_name])))
```

```
W1 max relative error: 1.380104e-04
W2 max relative error: 1.822723e-02
W3 max relative error: 3.064049e-04
b1 max relative error: 3.477652e-05
b2 max relative error: 2.516375e-03
b3 max relative error: 7.945660e-10
```

### 9.3 Overfit small data

A nice trick is to train your model with just a few training samples. You should be able to overfit small datasets, which will result in very high training accuracy and comparatively low validation accuracy.

```
[0]: np.random.seed(231)

num_train = 100
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'],
}

model = ThreeLayerConvNet(weight_scale=1e-2)

solver = Solver(model, small_data,
```

```

        num_epochs=15, batch_size=50,
        update_rule='adam',
        optim_config={
            'learning_rate': 1e-3,
        },
        verbose=True, print_every=1)
solver.train()

```

```

(Iteration 1 / 30) loss: 2.414060
(Epoch 0 / 15) train acc: 0.200000; val_acc: 0.137000
(Iteration 2 / 30) loss: 3.102925
(Epoch 1 / 15) train acc: 0.140000; val_acc: 0.087000
(Iteration 3 / 30) loss: 2.270330
(Iteration 4 / 30) loss: 2.096705
(Epoch 2 / 15) train acc: 0.240000; val_acc: 0.094000
(Iteration 5 / 30) loss: 1.838880
(Iteration 6 / 30) loss: 1.934188
(Epoch 3 / 15) train acc: 0.510000; val_acc: 0.173000
(Iteration 7 / 30) loss: 1.827912
(Iteration 8 / 30) loss: 1.639574
(Epoch 4 / 15) train acc: 0.520000; val_acc: 0.188000
(Iteration 9 / 30) loss: 1.330082
(Iteration 10 / 30) loss: 1.756115
(Epoch 5 / 15) train acc: 0.630000; val_acc: 0.167000
(Iteration 11 / 30) loss: 1.024162
(Iteration 12 / 30) loss: 1.041826
(Epoch 6 / 15) train acc: 0.750000; val_acc: 0.229000
(Iteration 13 / 30) loss: 1.142777
(Iteration 14 / 30) loss: 0.835706
(Epoch 7 / 15) train acc: 0.790000; val_acc: 0.247000
(Iteration 15 / 30) loss: 0.587786
(Iteration 16 / 30) loss: 0.645509
(Epoch 8 / 15) train acc: 0.820000; val_acc: 0.252000
(Iteration 17 / 30) loss: 0.786844
(Iteration 18 / 30) loss: 0.467054
(Epoch 9 / 15) train acc: 0.820000; val_acc: 0.178000
(Iteration 19 / 30) loss: 0.429880
(Iteration 20 / 30) loss: 0.635498
(Epoch 10 / 15) train acc: 0.900000; val_acc: 0.206000
(Iteration 21 / 30) loss: 0.365807
(Iteration 22 / 30) loss: 0.284220
(Epoch 11 / 15) train acc: 0.820000; val_acc: 0.201000
(Iteration 23 / 30) loss: 0.469343
(Iteration 24 / 30) loss: 0.509369
(Epoch 12 / 15) train acc: 0.920000; val_acc: 0.211000
(Iteration 25 / 30) loss: 0.111638
(Iteration 26 / 30) loss: 0.145388

```

```
(Epoch 13 / 15) train acc: 0.930000; val_acc: 0.213000
(Iteration 27 / 30) loss: 0.155575
(Iteration 28 / 30) loss: 0.143398
(Epoch 14 / 15) train acc: 0.960000; val_acc: 0.212000
(Iteration 29 / 30) loss: 0.158160
(Iteration 30 / 30) loss: 0.118934
(Epoch 15 / 15) train acc: 0.990000; val_acc: 0.220000
```

```
[0]: # Print final training accuracy
print(
    "Small data training accuracy:",
    solver.check_accuracy(small_data['X_train'], small_data['y_train'])
)
```

Small data training accuracy: 0.82

```
[0]: # Print final validation accuracy
print(
    "Small data validation accuracy:",
    solver.check_accuracy(small_data['X_val'], small_data['y_val'])
)
```

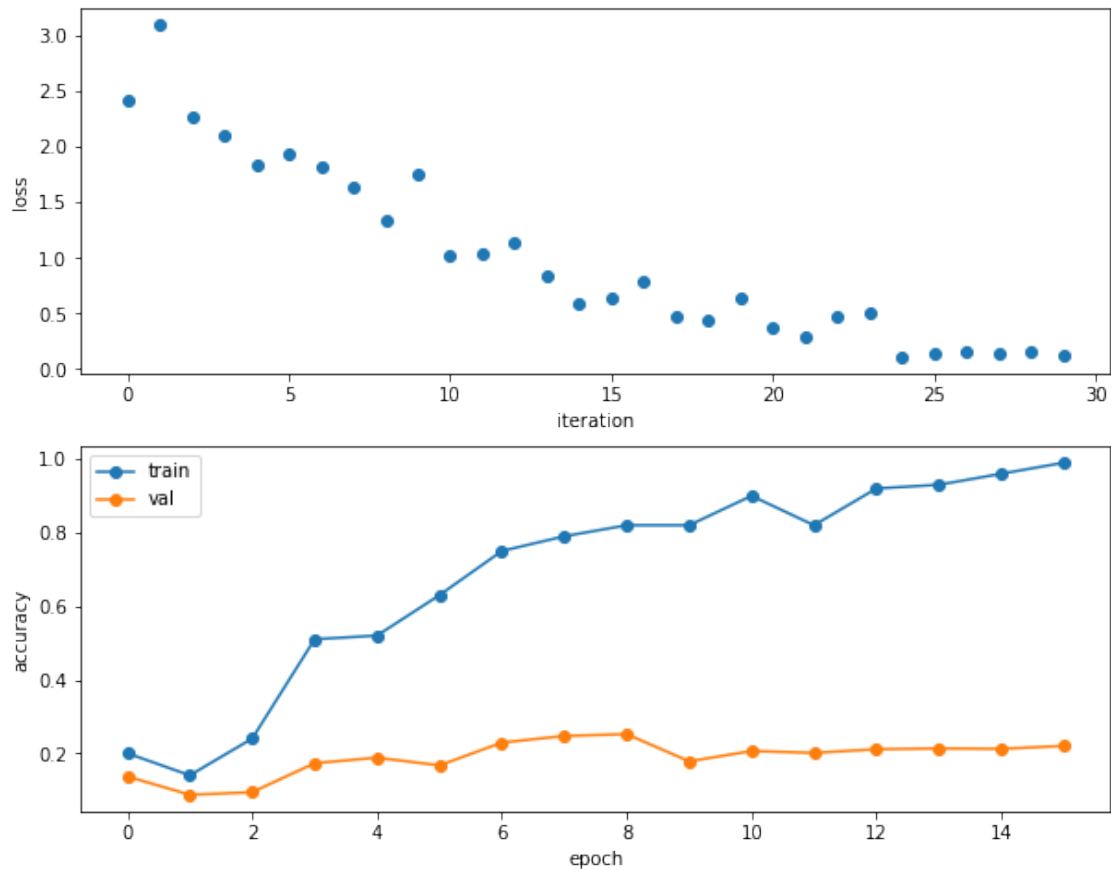
Small data validation accuracy: 0.252

Plotting the loss, training accuracy, and validation accuracy should show clear overfitting:

```
[0]: plt.subplot(2, 1, 1)
plt.plot(solver.loss_history, 'o')
plt.xlabel('iteration')
plt.ylabel('loss')

plt.subplot(2, 1, 2)
plt.plot(solver.train_acc_history, '-o')
plt.plot(solver.val_acc_history, '-o')
plt.legend(['train', 'val'], loc='upper left')
plt.xlabel('epoch')
plt.ylabel('accuracy')
plt.show()
```





## 9.4 Train the net

By training the three-layer convolutional network for one epoch, you should achieve greater than 40% accuracy on the training set:

```
[0]: model = ThreeLayerConvNet(weight_scale=0.001, hidden_dim=500, reg=0.001)

solver = Solver(model, data,
                 num_epochs=1, batch_size=50,
                 update_rule='adam',
                 optim_config={
                     'learning_rate': 1e-3,
                 },
                 verbose=True, print_every=20)
solver.train()
```

```
(Iteration 1 / 980) loss: 2.304647
(Epoch 0 / 1) train acc: 0.092000; val_acc: 0.087000
(Iteration 21 / 980) loss: 2.180548
(Iteration 41 / 980) loss: 1.755691
```

(Iteration 61 / 980) loss: 1.862384  
(Iteration 81 / 980) loss: 2.105984  
(Iteration 101 / 980) loss: 1.864433  
(Iteration 121 / 980) loss: 1.669615  
(Iteration 141 / 980) loss: 1.588268  
(Iteration 161 / 980) loss: 1.755229  
(Iteration 181 / 980) loss: 1.707818  
(Iteration 201 / 980) loss: 2.241003  
(Iteration 221 / 980) loss: 1.826519  
(Iteration 241 / 980) loss: 1.809789  
(Iteration 261 / 980) loss: 1.849963  
(Iteration 281 / 980) loss: 1.494044  
(Iteration 301 / 980) loss: 1.585512  
(Iteration 321 / 980) loss: 1.868825  
(Iteration 341 / 980) loss: 1.514661  
(Iteration 361 / 980) loss: 1.908325  
(Iteration 381 / 980) loss: 1.728278  
(Iteration 401 / 980) loss: 1.532411  
(Iteration 421 / 980) loss: 1.728782  
(Iteration 441 / 980) loss: 1.532376  
(Iteration 461 / 980) loss: 1.506589  
(Iteration 481 / 980) loss: 1.724665  
(Iteration 501 / 980) loss: 1.540936  
(Iteration 521 / 980) loss: 1.837002  
(Iteration 541 / 980) loss: 1.504125  
(Iteration 561 / 980) loss: 1.446879  
(Iteration 581 / 980) loss: 1.573095  
(Iteration 601 / 980) loss: 1.484758  
(Iteration 621 / 980) loss: 1.572361  
(Iteration 641 / 980) loss: 1.539987  
(Iteration 661 / 980) loss: 1.286299  
(Iteration 681 / 980) loss: 1.758689  
(Iteration 701 / 980) loss: 1.620575  
(Iteration 721 / 980) loss: 1.482950  
(Iteration 741 / 980) loss: 1.511648  
(Iteration 761 / 980) loss: 1.630812  
(Iteration 781 / 980) loss: 1.562391  
(Iteration 801 / 980) loss: 1.364693  
(Iteration 821 / 980) loss: 1.503164  
(Iteration 841 / 980) loss: 1.619294  
(Iteration 861 / 980) loss: 1.545362  
(Iteration 881 / 980) loss: 1.587725  
(Iteration 901 / 980) loss: 1.786351  
(Iteration 921 / 980) loss: 1.600237  
(Iteration 941 / 980) loss: 1.391016  
(Iteration 961 / 980) loss: 1.621583  
(Epoch 1 / 1) train acc: 0.494000; val\_acc: 0.480000

```
[0]: # Print final training accuracy
print(
    "Full data training accuracy:",
    solver.check_accuracy(data['X_train'], data['y_train'])
)
```

Full data training accuracy: 0.48316326530612247

```
[0]: # Print final validation accuracy
print(
    "Full data validation accuracy:",
    solver.check_accuracy(data['X_val'], data['y_val'])
)
```

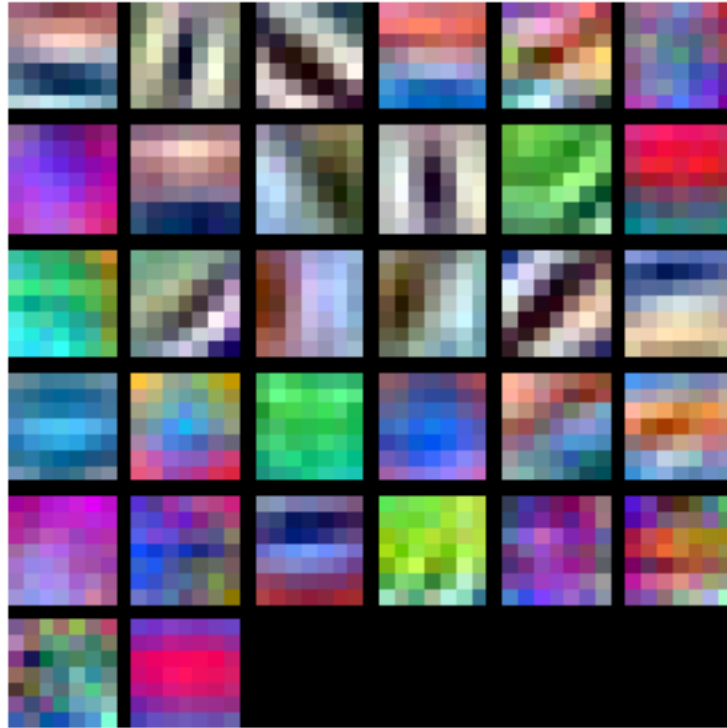
Full data validation accuracy: 0.48

## 9.5 Visualize Filters

You can visualize the first-layer convolutional filters from the trained network by running the following:

```
[0]: from cs231n.vis_utils import visualize_grid

grid = visualize_grid(model.params['W1'].transpose(0, 2, 3, 1))
plt.imshow(grid.astype('uint8'))
plt.axis('off')
plt.gcf().set_size_inches(5, 5)
plt.show()
```



## 10 Spatial Batch Normalization

We already saw that batch normalization is a very useful technique for training deep fully-connected networks. As proposed in the original paper (link in `BatchNormalization.ipynb`), batch normalization can also be used for convolutional networks, but we need to tweak it a bit; the modification will be called "spatial batch normalization."

Normally batch-normalization accepts inputs of shape  $(N, D)$  and produces outputs of shape  $(N, D)$ , where we normalize across the minibatch dimension  $N$ . For data coming from convolutional layers, batch normalization needs to accept inputs of shape  $(N, C, H, W)$  and produce outputs of shape  $(N, C, H, W)$  where the  $N$  dimension gives the minibatch size and the  $(H, W)$  dimensions give the spatial size of the feature map.

If the feature map was produced using convolutions, then we expect every feature channel's statistics e.g. mean, variance to be relatively consistent both between different images, and different locations within the same image -- after all, every feature channel is produced by the same convolutional filter! Therefore spatial batch normalization computes a mean and variance for each of the  $C$  feature channels by computing statistics over the minibatch dimension  $N$  as well the spatial dimensions  $H$  and  $W$ .

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

## 10.1 Spatial batch normalization: forward

In the file `cs231n/layers.py`, implement the forward pass for spatial batch normalization in the function `spatial_batchnorm_forward`. Check your implementation by running the following:

```
[0]: 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: [ 5.85642645e-16  5.93969318e-16 -8.88178420e-17]
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]
```

```
[0]: 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]
```

## 10.2 Spatial batch normalization: backward

In the file `cs231n/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:

```
[0]: 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:  3.083846825651097e-07
dgamma error:  7.09738489671469e-12
dbeta error:  3.275608725278405e-12
```

## 11 Group Normalization

In the previous notebook, we mentioned that Layer Normalization is an alternative normalization technique that mitigates the batch size limitations of Batch Normalization. However, as the authors of [2] observed, Layer Normalization does not perform as well as Batch Normalization when used with Convolutional Layers:

With fully connected layers, all the hidden units in a layer tend to make similar contributions to the final prediction, and re-centering and rescaling the summed inputs to a layer works well. However, the assumption of similar contributions is no longer true for convolutional neural networks. The large number of the hidden units whose receptive fields lie near the boundary of the image are rarely turned on and thus have very different statistics from the rest of the hidden units within the same layer.

The authors of [3] propose an intermediary technique. In contrast to Layer Normalization, where you normalize over the entire feature per-datapoint, they suggest a consistent splitting of each per-datapoint feature into  $G$  groups, and a per-group per-datapoint normalization instead.

Visual comparison of the normalization techniques discussed so far (image edited from [3])

Even though an assumption of equal contribution is still being made within each group, the authors hypothesize that this is not as problematic, as innate grouping arises within features for visual recognition. One example they use to illustrate this is that many high-performance hand-crafted features in traditional Computer Vision have terms that are explicitly grouped together. Take for example Histogram of Oriented Gradients [4]-- after computing histograms per spatially local block, each per-block histogram is normalized before being concatenated together to form the final feature vector.

You will now implement Group Normalization. Note that this normalization technique that you are to implement in the following cells was introduced and published to ECCV just in 2018 -- this truly is still an ongoing and excitingly active field of research!

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

[3] [Wu, Yuxin, and Kaiming He. "Group Normalization." arXiv preprint arXiv:1803.08494 (2018).](<https://arxiv.org/abs/1803.08494>)

[4] [N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In Computer Vision and Pattern Recognition (CVPR), 2005.](<https://ieeexplore.ieee.org/abstract/document/1467360/>)

## 11.1 Group normalization: forward

In the file `cs231n/layers.py`, implement the forward pass for group normalization in the function `spatial_groupnorm_forward`. Check your implementation by running the following:

```
[17]: 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, 6, 4, 5
G = 2
x = 4 * np.random.randn(N, C, H, W) + 10
x_g = x.reshape((N*G,-1))
print('Before spatial group normalization:')
print('  Shape: ', x.shape)
print('  Means: ', x_g.mean(axis=1))
print('  Stds: ', x_g.std(axis=1))

# Means should be close to zero and stds close to one
gamma, beta = np.ones((1,C,1,1)), np.zeros((1,C,1,1))
bn_param = {'mode': 'train'}

out, _ = spatial_groupnorm_forward(x, gamma, beta, G, bn_param)
out_g = out.reshape((N*G,-1))
print('After spatial group normalization:')
print('  Shape: ', out.shape)
print('  Means: ', out_g.mean(axis=1))
print('  Stds: ', out_g.std(axis=1))
```

Before spatial group normalization:

```
Shape: (2, 6, 4, 5)
Means: [9.72505327 8.51114185 8.9147544  9.43448077]
Stds:  [3.67070958 3.09892597 4.27043622 3.97521327]
```

After spatial group normalization:

```
Shape: (2, 6, 4, 5)
Means: [-2.14643118e-16  5.25505565e-16  2.58126853e-16 -3.62672855e-16]
Stds:  [0.99999963 0.99999948 0.99999973 0.99999968]
```

## 11.2 Spatial group normalization: backward

In the file `cs231n/layers.py`, implement the backward pass for spatial batch normalization in the function `spatial_groupnorm_backward`. Run the following to check your implementation using a numeric gradient check:

```
[18]: np.random.seed(231)
N, C, H, W = 2, 6, 4, 5
G = 2
x = 5 * np.random.randn(N, C, H, W) + 12
gamma = np.random.randn(1,C,1,1)
```



```

beta = np.random.randn(1,C,1,1)
dout = np.random.randn(N, C, H, W)

gn_param = {}
fx = lambda x: spatial_groupnorm_forward(x, gamma, beta, G, gn_param)[0]
fg = lambda a: spatial_groupnorm_forward(x, gamma, beta, G, gn_param)[0]
fb = lambda b: spatial_groupnorm_forward(x, gamma, beta, G, gn_param)[0]

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

_, cache = spatial_groupnorm_forward(x, gamma, beta, G, gn_param)
dx, dgamma, dbeta = spatial_groupnorm_backward(dout, cache)
#You should expect errors of magnitudes between 1e-12~1e-07
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:  6.345904562232329e-08
dgamma error:  1.0546047434202244e-11
dbeta error:  3.810857316122484e-12

```

# PyTorch

May 10, 2020

```
[0]: # 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. 'cs231n/assignments/assignment3/'
FOLDERNAME = "cs231n/assignments/assignment2/"
assert FOLDERNAME is not None, "[!] Enter the foldername."

# symlink to make it easier to load your files
!ln -s "/content/drive/My Drive/$FOLDERNAME" "/content/assignment2"

# 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/cs231n/datasets/
!bash get_datasets.sh
%cd /content
```

Mounted at /content/drive

```
ln: failed to create symbolic link '/content/assignment2/assignment2': Operation
not supported
/content/drive/My Drive/cs231n/assignments/assignment2/cs231n/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 code-base 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.

```
[0]: import torch
assert '.'.join(torch.__version__.split('.')[2]) == '1.5'
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
```

```
[0]: 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('./cs231n/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('./cs231n/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('./cs231n/datasets', train=False, download=True,
                             transform=transform)
loader_test = DataLoader(cifar10_test, batch_size=64)

```

Files already downloaded and verified  
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.

```

[0]: 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).

```
[0]: 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.

```
[0]: 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 = flatten(x) # shape: [batch\_size, C x H x W]*

*w1, w2 = params*

*# Forward pass: compute predicted y using operations on Tensors. Since w1*  
→*and*

*# w2 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(w1))*

*x = x.mm(w2)*

*return x*

**def** *two\_layer\_fc\_test*():

*hidden\_layer\_size = 42*

*x = torch.zeros((64, 50), dtype=dtype) # minibatch size 64, feature*  
→*dimension 50*

*w1 = torch.zeros((50, hidden\_layer\_size), dtype=dtype)*

*w2 = torch.zeros((hidden\_layer\_size, 10), dtype=dtype)*

*scores = two\_layer\_fc(x, [w1, w2])*

*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!

```
[0]: 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)*****

conv1 = F.conv2d(x, weight=conv_w1, bias=conv_b1, padding=2)
relu1 = F.relu(conv1)
conv2 = F.conv2d(relu1, weight=conv_w2, bias=conv_b2, padding=1)
relu2 = F.relu(conv2)
relu2_flat = flatten(relu2)
scores = relu2_flat.mm(fc_w) + fc_b

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

↳ #####
#
↳ #
#
↳ #####
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).

```

[0]: 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>

```
[0]: 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))
```

```
[0]: tensor([[ 0.9689, -0.0791,  0.4895,  0.0751, -0.0871],  
          [-0.0436, -0.7104,  1.3621, -0.5798, -0.8219],  
          [ 0.0617,  0.4377,  1.1820,  0.4716,  0.8677]], 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.

```
[0]: 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.

```
[0]: 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,  $w_1$  and  $w_2$ .

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  $w_1$ . The second dimension of  $w_1$  is the hidden layer size, which will also be the first dimension of  $w_2$ .

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.

```
[0]: 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 = 3.4991
Checking accuracy on the val set
Got 158 / 1000 correct (15.80%)
```

```
Iteration 100, loss = 2.2757
Checking accuracy on the val set
Got 276 / 1000 correct (27.60%)
```

```
Iteration 200, loss = 1.8914
Checking accuracy on the val set
Got 331 / 1000 correct (33.10%)
```

```
Iteration 300, loss = 1.6185
Checking accuracy on the val set
Got 394 / 1000 correct (39.40%)
```

```
Iteration 400, loss = 1.9287
Checking accuracy on the val set
Got 419 / 1000 correct (41.90%)
```

```
Iteration 500, loss = 2.0236
Checking accuracy on the val set
Got 437 / 1000 correct (43.70%)
```

```
Iteration 600, loss = 1.4150
Checking accuracy on the val set
Got 431 / 1000 correct (43.10%)
```

```
Iteration 700, loss = 1.8743
Checking accuracy on the val set
Got 432 / 1000 correct (43.20%)
```

#### 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.

```
[0]: 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)*****
#####
#                                     END OF YOUR CODE
→#
#####

params = [conv_w1, conv_b1, conv_w2, conv_b2, fc_w, fc_b]
train_part2(three_layer_convnet, params, learning_rate)
```

Iteration 0, loss = 2.6195

Checking accuracy on the val set

Got 119 / 1000 correct (11.90%)

```
Iteration 100, loss = 1.8748
Checking accuracy on the val set
Got 347 / 1000 correct (34.70%)
```

```
Iteration 200, loss = 1.7426
Checking accuracy on the val set
Got 385 / 1000 correct (38.50%)
```

```
Iteration 300, loss = 1.6614
Checking accuracy on the val set
Got 423 / 1000 correct (42.30%)
```

```
Iteration 400, loss = 1.4883
Checking accuracy on the val set
Got 425 / 1000 correct (42.50%)
```

```
Iteration 500, loss = 1.6403
Checking accuracy on the val set
Got 454 / 1000 correct (45.40%)
```

```
Iteration 600, loss = 1.5540
Checking accuracy on the val set
Got 471 / 1000 correct (47.10%)
```

```
Iteration 700, loss = 1.6749
Checking accuracy on the val set
Got 451 / 1000 correct (45.10%)
```

## 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 RM-SProp, 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:

```
[0]: 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.

```
[0]: 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,
        bias=True)
        nn.init.kaiming_normal_(self.conv1.weight)
        nn.init.constant_(self.conv1.bias, 0)

        self.conv2 = nn.Conv2d(channel_1, channel_2, kernel_size=3, padding=1,
        bias=True)
        nn.init.kaiming_normal_(self.conv2.weight)
        nn.init.constant_(self.conv2.bias, 0)

        self.fc = nn.Linear(channel_2*32*32, num_classes)
        nn.init.kaiming_normal_(self.fc.weight)
        nn.init.constant_(self.fc.bias, 0)

        # *****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)*****

    relu1 = F.relu(self.conv1(x))
    relu2 = F.relu(self.conv2(relu1))
    scores = self.fc(flatten(relu2))

    # *****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.

```

[0]: 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.

```

[0]: 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.

```

[0]: 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 = 3.4954
Checking accuracy on validation set
Got 137 / 1000 correct (13.70)

```

```

Iteration 100, loss = 2.3577
Checking accuracy on validation set
Got 339 / 1000 correct (33.90)

```

```

Iteration 200, loss = 2.0446
Checking accuracy on validation set
Got 384 / 1000 correct (38.40)

```

```

Iteration 300, loss = 1.7398
Checking accuracy on validation set
Got 423 / 1000 correct (42.30)

```

```
Iteration 400, loss = 1.8418
Checking accuracy on validation set
Got 434 / 1000 correct (43.40)
```

```
Iteration 500, loss = 1.7542
Checking accuracy on validation set
Got 457 / 1000 correct (45.70)
```

```
Iteration 600, loss = 1.8005
Checking accuracy on validation set
Got 405 / 1000 correct (40.50)
```

```
Iteration 700, loss = 1.3685
Checking accuracy on validation set
Got 462 / 1000 correct (46.20)
```

### 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.

```
[0]: 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(3, channel_1, channel_2, 10)
optimizer = 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 = 3.3847
Checking accuracy on validation set
```

```
Got 91 / 1000 correct (9.10)

Iteration 100, loss = 1.7058
Checking accuracy on validation set
Got 336 / 1000 correct (33.60)

Iteration 200, loss = 1.6296
Checking accuracy on validation set
Got 387 / 1000 correct (38.70)

Iteration 300, loss = 1.7425
Checking accuracy on validation set
Got 398 / 1000 correct (39.80)

Iteration 400, loss = 1.6108
Checking accuracy on validation set
Got 436 / 1000 correct (43.60)

Iteration 500, loss = 1.5165
Checking accuracy on validation set
Got 452 / 1000 correct (45.20)

Iteration 600, loss = 1.5041
Checking accuracy on validation set
Got 458 / 1000 correct (45.80)

Iteration 700, loss = 1.3591
Checking accuracy on validation set
Got 462 / 1000 correct (46.20)
```

## 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.

```
[0]: # 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 = 2.3314  
Checking accuracy on validation set  
Got 161 / 1000 correct (16.10)
```

```
Iteration 100, loss = 1.8255  
Checking accuracy on validation set  
Got 400 / 1000 correct (40.00)
```

```
Iteration 200, loss = 1.6462  
Checking accuracy on validation set  
Got 414 / 1000 correct (41.40)
```

```
Iteration 300, loss = 1.4844  
Checking accuracy on validation set  
Got 426 / 1000 correct (42.60)
```

```
Iteration 400, loss = 1.7098  
Checking accuracy on validation set  
Got 437 / 1000 correct (43.70)
```

```
Iteration 500, loss = 1.3872  
Checking accuracy on validation set  
Got 420 / 1000 correct (42.00)
```



```
Iteration 600, loss = 2.0308
Checking accuracy on validation set
Got 447 / 1000 correct (44.70)
```

```
Iteration 700, loss = 1.4962
Checking accuracy on validation set
Got 423 / 1000 correct (42.30)
```

## 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.

```
[0]: 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.
      #→#
      #####
      model = nn.Sequential( nn.Conv2d(3, channel_1, 5, padding = 2),
                             nn.ReLU(),
                             nn.Conv2d(channel_1, channel_2, 3, padding = 1),
                             nn.ReLU(),
                             Flatten(),
                             nn.Linear(channel_2 * 32 * 32, 10)
                           )
```

```
optimizer = optim.SGD(model.parameters(), lr = learning_rate, momentum = 0.9,
    →nesterov = True)
#####
#                               END OF YOUR CODE
#####

train_part34(model, optimizer)
```

```
Iteration 0, loss = 2.3244
Checking accuracy on validation set
Got 98 / 1000 correct (9.80)
```

```
Iteration 100, loss = 1.4736
Checking accuracy on validation set
Got 452 / 1000 correct (45.20)
```

```
Iteration 200, loss = 1.3742
Checking accuracy on validation set
Got 475 / 1000 correct (47.50)
```

```
Iteration 300, loss = 1.6068
Checking accuracy on validation set
Got 496 / 1000 correct (49.60)
```

```
Iteration 400, loss = 1.4587
Checking accuracy on validation set
Got 530 / 1000 correct (53.00)
```

```
Iteration 500, loss = 1.5674
Checking accuracy on validation set
Got 554 / 1000 correct (55.40)
```

```
Iteration 600, loss = 1.3241
Checking accuracy on validation set
Got 578 / 1000 correct (57.80)
```

```
Iteration 700, loss = 1.4336
Checking accuracy on validation set
Got 567 / 1000 correct (56.70)
```

## 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!

```
[0]: #####
# 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)*****

# function signature:
# Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0,
#        dilation=1, groups=1, bias=True, padding_mode='zeros')

# A 4-layer convolutional network
```

```

# (conv -> batchnorm -> relu -> maxpool) * 3 -> fc

outchannel_1, outchannel_2, outchannel_3, num_classes = 16, 32, 64, 10
inchannel_1, inchannel_2, inchannel_3 = 3, 16, 32
filter_size1, filter_size2, filter_size3 = 5, 3, 3

# first set of CONV combo layers
# Input: 3 x 32 x 32 raw image
conv1 = nn.Sequential(
    nn.Conv2d(inchannel_1, outchannel_1, kernel_size=filter_size1, padding=2),
    nn.BatchNorm2d(outchannel_1),      # Output: 32 x 32 x 16
    nn.ReLU(),                        # Output: 32 x 32 x 16
    nn.MaxPool2d(2)                   # Output: 16 x 16 x 16
)

# second set of CONV combo layers,
# Input: 16 x 16 x 16
conv2 = nn.Sequential(
    nn.Conv2d(inchannel_2, outchannel_2, kernel_size=filter_size2, padding=1),
    nn.BatchNorm2d(outchannel_2),      # Output: 16 x 16 x 32
    nn.ReLU(),                        # Output: 16 x 16 x 32
    nn.MaxPool2d(2)                   # Output: 8 x 8 x 32
)

# Third set of CONV combo layers,
# Input: 8 x 8 x 32
conv3 = nn.Sequential(
    nn.Conv2d(inchannel_3, outchannel_3, kernel_size=filter_size3, padding=1),
    nn.BatchNorm2d(outchannel_3),      # Output: 8 x 8 x 64
    nn.ReLU(),                        # Output: 8 x 8 x 64
    nn.MaxPool2d(2)                   # Output: 4 x 4 x 64
)

# Input: 4 x 4 x 64
# move to fully connected layers
fc = nn.Sequential(
    nn.Dropout(0.2, inplace=True),
    nn.Linear(64*4*4, num_classes)
)

model = nn.Sequential(
    conv1,
    conv2,
    conv3,
    Flatten(),
    fc
)

```

```

learning_rate = 1e-3

optimizer = optim.Adam(model.parameters(), lr=learning_rate)

# Print training status every epoch: set print_every to a large number
print_every = 10000

# *****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)

# Experimental Model
# sample_model = nn.Sequential(
#     # first convolutional combo layers
#     # Input: 3 x 32 x 32 raw image
#     nn.Conv2d(3, 32, kernel_size=7, stride=1),
#     nn.BatchNorm2d(32),
#     nn.ReLU(inplace=True),
#     # nn.MaxPool2d(kernel_size=2, stride=2),
#     # Output: 26 x 26 x 32

#     # second convolutional combo layers,
#     # Input: 26 x 26 x 32
#     nn.Conv2d(32, 24, kernel_size=5, stride=1),
#     nn.BatchNorm2d(24),
#     nn.ReLU(inplace=True),
#     # Output: 22 x 22 x 24
#     nn.MaxPool2d(kernel_size=2, stride=2),
#     # Output: 11 x 11 x 24

#     # Third convolutional combo layers,
#     # Input: 11 x 11 x 32
#     nn.Conv2d(24, 24, kernel_size=4, stride=1),
#     nn.BatchNorm2d(24),
#     nn.ReLU(inplace=True),
#     # Output: 8 x 8 x 24
#     nn.MaxPool2d(kernel_size=2, stride=2),
#     # Output: 4 x 4 x 24

#     # move to fully connected layers
#     nn.Dropout(0.2, inplace=True),
#     Flatten(),

```

```

#             nn.Linear(384, 256),
#             nn.ReLU(inplace=True),
#             nn.Linear(256, 10)
#         )
# #sample_model.type(gpu_dtype)

# # Print training status every epoch: set print_every to a large number
# print_every = 10000

# #loss_fn = nn.CrossEntropyLoss().type(gpu_dtype)
# optimizer = optim.RMSprop(sample_model.parameters(), lr=1e-3)
# # optimizer = optim.Adam(model.parameters(), lr=learning_rate)

# #train(sample_model, loss_fn, optimizer, num_epochs=10)
# train_part34(sample_model, optimizer, epochs=10)

```

Iteration 0, loss = 2.5262  
Checking accuracy on validation set  
Got 133 / 1000 correct (13.30)

Iteration 0, loss = 1.1148  
Checking accuracy on validation set  
Got 619 / 1000 correct (61.90)

Iteration 0, loss = 1.1139  
Checking accuracy on validation set  
Got 694 / 1000 correct (69.40)

Iteration 0, loss = 0.7234  
Checking accuracy on validation set  
Got 723 / 1000 correct (72.30)

Iteration 0, loss = 0.7457  
Checking accuracy on validation set  
Got 732 / 1000 correct (73.20)

Iteration 0, loss = 0.6344  
Checking accuracy on validation set  
Got 711 / 1000 correct (71.10)

Iteration 0, loss = 0.5841  
Checking accuracy on validation set  
Got 737 / 1000 correct (73.70)

Iteration 0, loss = 0.5105  
Checking accuracy on validation set  
Got 734 / 1000 correct (73.40)

```
Iteration 0, loss = 0.6077
Checking accuracy on validation set
Got 763 / 1000 correct (76.30)
```

```
Iteration 0, loss = 0.5303
Checking accuracy on validation set
Got 760 / 1000 correct (76.00)
```

## 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.

### Network Architecture:

$(CONV \rightarrow BatchNorm \rightarrow ReLu \rightarrow MaxPool) * 3 \rightarrow Dropout \rightarrow FC$

**Filter Sizes:** 5x5, 3x3, 3x3

**Number of Filters:** 16, 32, 64

**Pooling:** 2x2 Max-pooling

**Normalization:** Batch normalization

**Regularization:** Dropout

**Optimizer:** Adam

**Specifics:**

- The key here is that after the first CONV layer, we shrink the filter size in the next one. This enabled us to push the accuracy to 60% in the second epoch.
- We used Batchnorm because it enables training with larger learning rates, which leads to faster convergence and better generalization. Layernorm would probably have a similar effect; however, because it has been found to be less effective than batchnorm, the convergence wouldn't probably gain as much of a speedup.
- The final push was using Dropout right before the Fully Connected layer which got us an additional 1.6% (from 74.4% to 76.0%) on the model accuracy. Dropout helps prevent overfitting.

## 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.

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

```
Checking accuracy on test set
Got 7325 / 10000 correct (73.25)
```



# TensorFlow

May 10, 2020

```
[0]: # 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. 'cs231n/assignments/assignment3/'
FOLDERNAME = "cs231n/assignments/assignment2/"
assert FOLDERNAME is not None, "[!] Enter the foldername."

# symlink to make it easier to load your files
!ln -s "/content/drive/My Drive/$FOLDERNAME" "/content/assignment2"

# 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/cs231n/datasets/
!bash get_datasets.sh
%cd /content
```

Mounted at /content/drive

ln: failed to create symbolic link '/content/assignment2/assignment2': Operation not supported

/content/drive/My Drive/cs231n/assignments/assignment2/cs231n/datasets  
/content

## 1 What's this TensorFlow 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 code-base and instead migrate to one of two popular deep learning frameworks: in this instance, TensorFlow (or PyTorch, if you choose to work with that notebook).

**What is it?** TensorFlow is a system for executing computational graphs over Tensor objects, with native support for performing backpropagation for its Variables. In it, we work with Tensors which are n-dimensional arrays analogous to the numpy ndarray.

**Why?**

- Our code will now run on GPUs! Much faster training. Writing your own modules to run on GPUs is beyond the scope of this class, unfortunately.
- 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.

## 1.1 How will I learn TensorFlow?

TensorFlow has many excellent tutorials available, including those from [Google themselves](#).

Otherwise, this notebook will walk you through much of what you need to do to train models in TensorFlow. See the end of the notebook for some links to helpful tutorials if you want to learn more or need further clarification on topics that aren't fully explained here.

**NOTE: This notebook is meant to teach you the latest version of Tensorflow which is as of this homework version 2.2.0-rc3. Most examples on the web today are still in 1.x, so be careful not to confuse the two when looking up documentation.**

## 1.2 Install Tensorflow 2.0 (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 `tf_20_env`.
3. Run the command: `source activate tf_20_env`
4. Then pip install TF 2.0 as described here: <https://www.tensorflow.org/install>

## 2 Table of Contents

This notebook has 5 parts. We will walk through TensorFlow at **three different levels of abstraction**, which should help you better understand it and prepare you for working on your project.

1. Part I, Preparation: load the CIFAR-10 dataset.
2. Part II, Barebone TensorFlow: **Abstraction Level 1**, we will work directly with low-level TensorFlow graphs.
3. Part III, Keras Model API: **Abstraction Level 2**, we will use `tf.keras.Model` to define arbitrary neural network architecture.

4. Part IV, Keras Sequential + Functional API: **Abstraction Level 3**, we will use `tf.keras.Sequential` to define a linear feed-forward network very conveniently, and then explore the functional libraries for building unique and uncommon models that require more flexibility.
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.

We will discuss Keras in more detail later in the notebook.

Here is a table of comparison:

API	Flexibility	Convenience
Barebone	High	Low
<code>tf.keras.Model</code>	High	Medium
<code>tf.keras.Sequential</code>	Low	High

### 3 Part I: Preparation

First, we load the CIFAR-10 dataset. This might take a few minutes to download the first time you run it, but after that the files should be cached on disk and loading should be faster.

In previous parts of the assignment we used CS231N-specific code to download and read the CIFAR-10 dataset; however the `tf.keras.datasets` package in TensorFlow provides prebuilt utility functions for loading many common datasets.

For the purposes of this assignment we will still write our own code to preprocess the data and iterate through it in minibatches. The `tf.data` package in TensorFlow provides tools for automating this process, but working with this package adds extra complication and is beyond the scope of this notebook. However using `tf.data` can be much more efficient than the simple approach used in this notebook, so you should consider using it for your project.

```
[0]: import os
import tensorflow as tf
import numpy as np
import math
import timeit
import matplotlib.pyplot as plt

%matplotlib inline

[0]: def load_cifar10(num_training=49000, num_validation=1000, num_test=10000):
    """
    Fetch the CIFAR-10 dataset from the web and perform preprocessing to
    ↪prepare
    it for the two-layer neural net classifier. These are the same steps as
    we used for the SVM, but condensed to a single function.
    """
    # Load the raw CIFAR-10 dataset and use appropriate data types and shapes
    cifar10 = tf.keras.datasets.cifar10.load_data()
```

```

(X_train, y_train), (X_test, y_test) = cifar10
X_train = np.asarray(X_train, dtype=np.float32)
y_train = np.asarray(y_train, dtype=np.int32).flatten()
X_test = np.asarray(X_test, dtype=np.float32)
y_test = np.asarray(y_test, dtype=np.int32).flatten()

# Subsample the data
mask = range(num_training, num_training + num_validation)
X_val = X_train[mask]
y_val = y_train[mask]
mask = range(num_training)
X_train = X_train[mask]
y_train = y_train[mask]
mask = range(num_test)
X_test = X_test[mask]
y_test = y_test[mask]

# Normalize the data: subtract the mean pixel and divide by std
mean_pixel = X_train.mean(axis=(0, 1, 2), keepdims=True)
std_pixel = X_train.std(axis=(0, 1, 2), keepdims=True)
X_train = (X_train - mean_pixel) / std_pixel
X_val = (X_val - mean_pixel) / std_pixel
X_test = (X_test - mean_pixel) / std_pixel

return X_train, y_train, X_val, y_val, X_test, y_test

# If there are errors with SSL downloading involving self-signed certificates,
# it may be that your Python version was recently installed on the current
→machine.
# See: https://github.com/tensorflow/tensorflow/issues/10779
# To fix, run the command: /Applications/Python\ 3.7/Install\ Certificates.
→command
# ...replacing paths as necessary.

# Invoke the above function to get our data.
NHW = (0, 1, 2)
X_train, y_train, X_val, y_val, X_test, y_test = load_cifar10()
print('Train data shape: ', X_train.shape)
print('Train labels shape: ', y_train.shape, y_train.dtype)
print('Validation data shape: ', X_val.shape)
print('Validation labels shape: ', y_val.shape)
print('Test data shape: ', X_test.shape)
print('Test labels shape: ', y_test.shape)

```

```

Train data shape: (49000, 32, 32, 3)
Train labels shape: (49000,) int32
Validation data shape: (1000, 32, 32, 3)

```

```
Validation labels shape: (1000,)
Test data shape: (10000, 32, 32, 3)
Test labels shape: (10000,)
```

```
[0]: class Dataset(object):
    def __init__(self, X, y, batch_size, shuffle=False):
        """
        Construct a Dataset object to iterate over data X and labels y

        Inputs:
        - X: Numpy array of data, of any shape
        - y: Numpy array of labels, of any shape but with y.shape[0] == X.
        →shape[0]
        - batch_size: Integer giving number of elements per minibatch
        - shuffle: (optional) Boolean, whether to shuffle the data on each
        →epoch
        """
        assert X.shape[0] == y.shape[0], 'Got different numbers of data and
        →labels'
        self.X, self.y = X, y
        self.batch_size, self.shuffle = batch_size, shuffle

    def __iter__(self):
        N, B = self.X.shape[0], self.batch_size
        idxs = np.arange(N)
        if self.shuffle:
            np.random.shuffle(idxs)
        return iter((self.X[i:i+B], self.y[i:i+B]) for i in range(0, N, B))

train_dset = Dataset(X_train, y_train, batch_size=64, shuffle=True)
val_dset = Dataset(X_val, y_val, batch_size=64, shuffle=False)
test_dset = Dataset(X_test, y_test, batch_size=64)
```

```
[0]: # We can iterate through a dataset like this:
for t, (x, y) in enumerate(train_dset):
    print(t, x.shape, y.shape)
    if t > 5: break
```

```
0 (64, 32, 32, 3) (64,)
1 (64, 32, 32, 3) (64,)
2 (64, 32, 32, 3) (64,)
3 (64, 32, 32, 3) (64,)
4 (64, 32, 32, 3) (64,)
5 (64, 32, 32, 3) (64,)
6 (64, 32, 32, 3) (64,)
```

You can optionally use GPU by setting the flag to True below.

### 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.

```
[0]: # Set up some global variables
USE_GPU = True

if USE_GPU:
    device = '/device:GPU:0'
else:
    device = '/cpu:0'

# Constant to control how often we print when training models
print_every = 100

print('Using device: ', device)
```

Using device: /device:GPU:0

## 4 Part II: Barebones TensorFlow

TensorFlow ships with various high-level APIs which make it very convenient to define and train neural networks; we will cover some of these constructs in Part III and Part IV of this notebook. In this section we will start by building a model with basic TensorFlow constructs to help you better understand what's going on under the hood of the higher-level APIs.

**"Barebones Tensorflow" is important to understanding the building blocks of TensorFlow, but much of it involves concepts from TensorFlow 1.x.** We will be working with legacy modules such as `tf.Variable`.

Therefore, please read and understand the differences between legacy (1.x) TF and the new (2.0) TF.

### 4.0.1 Historical background on TensorFlow 1.x

TensorFlow 1.x is primarily a framework for working with **static computational graphs**. Nodes in the computational graph are Tensors which will hold n-dimensional arrays when the graph is run; edges in the graph represent functions that will operate on Tensors when the graph is run to actually perform useful computation.

Before Tensorflow 2.0, we had to configure the graph into two phases. There are plenty of tutorials online that explain this two-step process. The process generally looks like the following for TF 1.x: 1. **Build a computational graph that describes the computation that you want to perform.** This stage doesn't actually perform any computation; it just builds up a symbolic representation of your computation. This stage will typically define one or more placeholder objects that represent inputs to the computational graph. 2. **Run the computational graph many times.** Each time the graph is run (e.g. for one gradient descent step) you will specify which parts of the graph you want to compute, and pass a `feed_dict` dictionary that will give concrete values to any placeholders in the graph.

## 4.0.2 The new paradigm in Tensorflow 2.0

Now, with Tensorflow 2.0, we can simply adopt a functional form that is more Pythonic and similar in spirit to PyTorch and direct Numpy operation. Instead of the 2-step paradigm with computation graphs, making it (among other things) easier to debug TF code. You can read more details at <https://www.tensorflow.org/guide/eager>.

The main difference between the TF 1.x and 2.0 approach is that the 2.0 approach doesn't make use of `tf.Session`, `tf.run`, `placeholder`, `feed_dict`. To get more details of what's different between the two version and how to convert between the two, check out the official migration guide: [https://www.tensorflow.org/alpha/guide/migration\\_guide](https://www.tensorflow.org/alpha/guide/migration_guide)

Later, in the rest of this notebook we'll focus on this new, simpler approach.

## 4.0.3 TensorFlow warmup: Flatten Function

We can see this in action by defining a simple `flatten` function that will reshape image data for use in a fully-connected network.

In TensorFlow, data for convolutional feature maps is typically stored in a Tensor of shape  $N \times H \times W \times C$  where:

- $N$  is the number of datapoints (minibatch size)
- $H$  is the height of the feature map
- $W$  is the width of the feature map
- $C$  is the number of channels in the feature map

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  $H \times W \times C$  values per representation into a single long vector.

Notice the `tf.reshape` call has the target shape as  $(N, -1)$ , meaning it will reshape/keep the first dimension to be  $N$ , and then infer as necessary what the second dimension is in the output, so we can collapse the remaining dimensions from the input properly.

**NOTE:** TensorFlow and PyTorch differ on the default Tensor layout; TensorFlow uses  $N \times H \times W \times C$  but PyTorch uses  $N \times C \times H \times W$ .

```
[0]: def flatten(x):  
    """  
    Input:  
    - TensorFlow Tensor of shape (N, D1, ..., DM)  
  
    Output:  
    - TensorFlow Tensor of shape (N, D1 * ... * DM)  
    """  
    N = tf.shape(x)[0]  
    return tf.reshape(x, (N, -1))
```

```
[0]: def test_flatten():  
    # Construct concrete values of the input data x using numpy  
    x_np = np.arange(24).reshape((2, 3, 4))
```

```

print('x_np:\n', x_np, '\n')
# Compute a concrete output value.
x_flat_np = flatten(x_np)
print('x_flat_np:\n', x_flat_np, '\n')

test_flatten()

```

```

x_np:
[[[ 0  1  2  3]
  [ 4  5  6  7]
  [ 8  9 10 11]]

 [[12 13 14 15]
  [16 17 18 19]
  [20 21 22 23]]]

x_flat_np:
tf.Tensor(
[[ 0  1  2  3  4  5  6  7  8  9 10 11]
 [12 13 14 15 16 17 18 19 20 21 22 23]], shape=(2, 12), dtype=int64)

```

#### 4.0.4 Barebones TensorFlow: Define a Two-Layer Network

We will now implement our first neural network with TensorFlow: a fully-connected ReLU network with two hidden layers and no biases on the CIFAR10 dataset. For now we will use only low-level TensorFlow operators to define the network; later we will see how to use the higher-level abstractions provided by `tf.keras` to simplify the process.

We will define the forward pass of the network in the function `two_layer_fc`; this will accept TensorFlow Tensors for the inputs and weights of the network, and return a TensorFlow Tensor for the scores.

After defining the network architecture in the `two_layer_fc` function, we will test the implementation by checking the shape of the output.

**It's important that you read and understand this implementation.**

```

[0]: def two_layer_fc(x, params):
      """
      A fully-connected neural network; the architecture is:
      fully-connected layer -> ReLU -> fully connected layer.
      Note that we only need to define the forward pass here; TensorFlow will
      ↪take
      care of computing the gradients 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 H1
      ↪units,
      and the output layer will produce scores for C classes.

```



*Inputs:*

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

*Returns:*

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

"""

```
w1, w2 = params                # Unpack the parameters
x = flatten(x)                 # Flatten the input; now x has shape (N, D)
h = tf.nn.relu(tf.matmul(x, w1)) # Hidden layer: h has shape (N, H)
scores = tf.matmul(h, w2)       # Compute scores of shape (N, C)
return scores
```

```
[0]: def two_layer_fc_test():
    hidden_layer_size = 42

    # Scoping our TF operations under a tf.device context manager
    # lets us tell TensorFlow where we want these Tensors to be
    # multiplied and/or operated on, e.g. on a CPU or a GPU.
    with tf.device(device):
        x = tf.zeros((64, 32, 32, 3))
        w1 = tf.zeros((32 * 32 * 3, hidden_layer_size))
        w2 = tf.zeros((hidden_layer_size, 10))

        # Call our two_layer_fc function for the forward pass of the network.
        scores = two_layer_fc(x, [w1, w2])

    print(scores.shape)

two_layer_fc_test()
```

(64, 10)

#### 4.0.5 Barebones TensorFlow: 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. The network should have the following architecture:

1. A convolutional layer (with bias) with `channel_1` filters, each with shape  $KW_1 \times KH_1$ , and zero-padding of two
2. ReLU nonlinearity
3. A convolutional layer (with bias) with `channel_2` filters, each with shape  $KW_2 \times KH_2$ , and zero-padding of one

4. ReLU nonlinearity
5. Fully-connected layer with bias, producing scores for C classes.

**HINT:** For convolutions: [https://www.tensorflow.org/versions/r2.0/api\\_docs/python/tf/nn/conv2d](https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/nn/conv2d); be careful with padding!

**HINT:** For biases: <https://www.tensorflow.org/performance/xla/broadcasting>

```
[0]: def three_layer_convnet(x, params):
    """
    A three-layer convolutional network with the architecture described above.

    Inputs:
    - x: A TensorFlow Tensor of shape (N, H, W, 3) giving a minibatch of images
    - params: A list of TensorFlow Tensors giving the weights and biases for
    → the
      network; should contain the following:
    - conv_w1: TensorFlow Tensor of shape (KH1, KW1, 3, channel_1) giving
      weights for the first convolutional layer.
    - conv_b1: TensorFlow Tensor of shape (channel_1,) giving biases for the
      first convolutional layer.
    - conv_w2: TensorFlow Tensor of shape (KH2, KW2, channel_1, channel_2)
      giving weights for the second convolutional layer
    - conv_b2: TensorFlow Tensor of shape (channel_2,) giving biases for the
      second convolutional layer.
    - fc_w: TensorFlow Tensor giving weights for the fully-connected layer.
      Can you figure out what the shape should be?
    - fc_b: TensorFlow Tensor giving biases for the fully-connected layer.
      Can you figure out what the shape should be?
    """
    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)*****

    paddings = tf.constant([[0,0], [2,2], [2,2], [0,0]])
    x = tf.pad(x, paddings, 'CONSTANT')
    conv1 = tf.nn.conv2d(x, conv_w1, strides=[1,1,1,1], padding="VALID")+conv_b1
    relu1 = tf.nn.relu(conv1)

    paddings = tf.constant([[0,0], [1,1], [1,1], [0,0]])
    conv1 = tf.pad(conv1, paddings, 'CONSTANT')
    conv2 = tf.nn.conv2d(conv1, conv_w2, strides=[1,1,1,1],
    →padding="VALID")+conv_b2
```

```

relu2 = tf.nn.relu(conv2)

relu2 = flatten(relu2)
scores = tf.matmul(relu2, 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 three-layer ConvNet above, run the following cell to test your implementation. Like the two-layer network, we run the graph on a batch of zeros just to make sure the function doesn't crash, and produces outputs of the correct shape.

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

```

[0]: def three_layer_convnet_test():

    with tf.device(device):
        x = tf.zeros((64, 32, 32, 3))
        conv_w1 = tf.zeros((5, 5, 3, 6))
        conv_b1 = tf.zeros((6,))
        conv_w2 = tf.zeros((3, 3, 6, 9))
        conv_b2 = tf.zeros((9,))
        fc_w = tf.zeros((32 * 32 * 9, 10))
        fc_b = tf.zeros((10,))
        params = [conv_w1, conv_b1, conv_w2, conv_b2, fc_w, fc_b]
        scores = three_layer_convnet(x, params)

    # Inputs to convolutional layers are 4-dimensional arrays with shape
    # [batch_size, height, width, channels]
    print('scores_np has shape: ', scores.shape)

three_layer_convnet_test()

```

scores\_np has shape: (64, 10)

#### 4.0.6 Barebones TensorFlow: Training Step

We now define the `training_step` function performs a single training step. This will take three basic steps:

1. Compute the loss
2. Compute the gradient of the loss with respect to all network weights
3. Make a weight update step using (stochastic) gradient descent.

We need to use a few new TensorFlow functions to do all of this: - For computing the cross-entropy loss we'll use `tf.nn.sparse_softmax_cross_entropy_with_logits`: [https://www.tensorflow.org/versions/r2.0/api\\_docs/python/tf/nn/sparse\\_softmax\\_cross\\_entropy\\_with\\_logits](https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/nn/sparse_softmax_cross_entropy_with_logits)

- For averaging the loss across a minibatch of data we'll use `tf.reduce_mean`: [https://www.tensorflow.org/versions/r2.0/api\\_docs/python/tf/reduce\\_mean](https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/reduce_mean)
- For computing gradients of the loss with respect to the weights we'll use `tf.GradientTape` (useful for Eager execution): [https://www.tensorflow.org/versions/r2.0/api\\_docs/python/tf/GradientTape](https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/GradientTape)
- We'll mutate the weight values stored in a TensorFlow Tensor using `tf.assign_sub` ("sub" is for subtraction): [https://www.tensorflow.org/versions/r2.0/api\\_docs/python/tf/assign\\_sub](https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/assign_sub)

```
[0]: def training_step(model_fn, x, y, params, learning_rate):  
    with tf.GradientTape() as tape:  
        scores = model_fn(x, params) # Forward pass of the model  
        loss = tf.nn.sparse_softmax_cross_entropy_with_logits(labels=y,   
->logits=scores)  
        total_loss = tf.reduce_mean(loss)  
        grad_params = tape.gradient(total_loss, params)  
  
        # Make a vanilla gradient descent step on all of the model parameters  
        # Manually update the weights using assign_sub()  
        for w, grad_w in zip(params, grad_params):  
            w.assign_sub(learning_rate * grad_w)  
  
    return total_loss
```

```
[0]: def train_part2(model_fn, init_fn, learning_rate):  
    """  
    Train a model on CIFAR-10.  
  
    Inputs:  
    - model_fn: A Python function that performs the forward pass of the model  
      using TensorFlow; it should have the following signature:  
      scores = model_fn(x, params) where x is a TensorFlow Tensor giving a  
      minibatch of image data, params is a list of TensorFlow Tensors holding  
      the model weights, and scores is a TensorFlow Tensor of shape (N, C)  
      giving scores for all elements of x.  
    - init_fn: A Python function that initializes the parameters of the model.  
      It should have the signature params = init_fn() where params is a list  
      of TensorFlow Tensors holding the (randomly initialized) weights of the  
      model.  
    - learning_rate: Python float giving the learning rate to use for SGD.  
    """  
  
    params = init_fn() # Initialize the model parameters
```

```

for t, (x_np, y_np) in enumerate(train_dset):
    # Run the graph on a batch of training data.
    loss = training_step(model_fn, x_np, y_np, params, learning_rate)

    # Periodically print the loss and check accuracy on the val set.
    if t % print_every == 0:
        print('Iteration %d, loss = %.4f' % (t, loss))
        check_accuracy(val_dset, x_np, model_fn, params)

```

```

[0]: def check_accuracy(dset, x, model_fn, params):
    """
    Check accuracy on a classification model, e.g. for validation.

    Inputs:
    - dset: A Dataset object against which to check accuracy
    - x: A TensorFlow placeholder Tensor where input images should be fed
    - model_fn: the Model we will be calling to make predictions on x
    - params: parameters for the model_fn to work with

    Returns: Nothing, but prints the accuracy of the model
    """
    num_correct, num_samples = 0, 0
    for x_batch, y_batch in dset:
        scores_np = model_fn(x_batch, params).numpy()
        y_pred = scores_np.argmax(axis=1)
        num_samples += x_batch.shape[0]
        num_correct += (y_pred == y_batch).sum()
    acc = float(num_correct) / num_samples
    print('Got %d / %d correct (%.2f%%)' % (num_correct, num_samples, 100 *
→acc))

```

#### 4.0.7 Barebones TensorFlow: Initialization

We'll use the following utility method to initialize the weight matrices for our models using Kaiming's normalization method.

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

```

[0]: def create_matrix_with_kaiming_normal(shape):
    if len(shape) == 2:
        fan_in, fan_out = shape[0], shape[1]
    elif len(shape) == 4:
        fan_in, fan_out = np.prod(shape[:3]), shape[3]
    return tf.keras.backend.random_normal(shape) * np.sqrt(2.0 / fan_in)

```

#### 4.0.8 Barebones TensorFlow: Train a Two-Layer Network

We are finally ready to use all of the pieces defined above to train a two-layer fully-connected network on CIFAR-10.

We just need to define a function to initialize the weights of the model, and call `train_part2`.

Defining the weights of the network introduces another important piece of TensorFlow API: `tf.Variable`. A TensorFlow Variable is a Tensor whose value is stored in the graph and persists across runs of the computational graph; however unlike constants defined with `tf.zeros` or `tf.random_normal`, the values of a Variable can be mutated as the graph runs; these mutations will persist across graph runs. Learnable parameters of the network are usually stored in Variables.

You don't need to tune any hyperparameters, but you should achieve validation accuracies above 40% after one epoch of training.

```
[0]: def two_layer_fc_init():
    """
    Initialize the weights of a two-layer network, for use with the
    two_layer_network function defined above.
    You can use the `create_matrix_with_kaiming_normal` helper!

    Inputs: None

    Returns: A list of:
    - w1: TensorFlow tf.Variable giving the weights for the first layer
    - w2: TensorFlow tf.Variable giving the weights for the second layer
    """
    hidden_layer_size = 4000
    w1 = tf.Variable(create_matrix_with_kaiming_normal((3 * 32 * 32, 4000)))
    w2 = tf.Variable(create_matrix_with_kaiming_normal((4000, 10)))
    return [w1, w2]

learning_rate = 1e-2
train_part2(two_layer_fc, two_layer_fc_init, learning_rate)
```

```
Iteration 0, loss = 2.8191
Got 150 / 1000 correct (15.00%)
Iteration 100, loss = 1.8836
Got 367 / 1000 correct (36.70%)
Iteration 200, loss = 1.4546
Got 401 / 1000 correct (40.10%)
Iteration 300, loss = 1.7125
Got 371 / 1000 correct (37.10%)
Iteration 400, loss = 1.8171
Got 430 / 1000 correct (43.00%)
Iteration 500, loss = 1.7354
Got 431 / 1000 correct (43.10%)
Iteration 600, loss = 1.8422
Got 434 / 1000 correct (43.40%)
Iteration 700, loss = 1.9712
Got 458 / 1000 correct (45.80%)
```

#### 4.0.9 Barebones TensorFlow: Train a three-layer ConvNet

We will now use TensorFlow to train a three-layer ConvNet on CIFAR-10.

You need to implement the `three_layer_convnet_init` function. Recall that the architecture of the network is:

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

You don't need to do any hyperparameter tuning, but you should see validation accuracies above 43% after one epoch of training.

```
[0]: def three_layer_convnet_init():
    """
    Initialize the weights of a Three-Layer ConvNet, for use with the
    three_layer_convnet function defined above.
    You can use the `create_matrix_with_kaiming_normal` helper!

    Inputs: None

    Returns a list containing:
    - conv_w1: TensorFlow tf.Variable giving weights for the first conv layer
    - conv_b1: TensorFlow tf.Variable giving biases for the first conv layer
    - conv_w2: TensorFlow tf.Variable giving weights for the second conv layer
    - conv_b2: TensorFlow tf.Variable giving biases for the second conv layer
    - fc_w: TensorFlow tf.Variable giving weights for the fully-connected layer
    - fc_b: TensorFlow tf.Variable giving biases for the fully-connected layer
    """
    params = None

    # TODO: Initialize the parameters of the three-layer network.

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

    conv_w1 = tf.Variable(create_matrix_with_kaiming_normal([5, 5, 3, 32]))
    conv_b1 = tf.Variable(np.zeros([32]), dtype=tf.float32)
    conv_w2 = tf.Variable(create_matrix_with_kaiming_normal([3, 3, 32, 16]))
    conv_b2 = tf.Variable(np.zeros([16]), dtype=tf.float32)
    fc_w = tf.Variable(create_matrix_with_kaiming_normal([32*32*16,10]))
    fc_b = tf.Variable(np.zeros([10]), dtype=tf.float32)
    params = (conv_w1, conv_b1, conv_w2, conv_b2, fc_w, fc_b)

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

```

    ↪#####
    #                                END OF YOUR CODE                                ↪
    ↪#
    ↪#####
    return params

learning_rate = 3e-3
train_part2(three_layer_convnet, three_layer_convnet_init, learning_rate)

```

```

Iteration 0, loss = 4.0863
Got 124 / 1000 correct (12.40%)
Iteration 100, loss = 1.7886
Got 397 / 1000 correct (39.70%)
Iteration 200, loss = 1.4878
Got 439 / 1000 correct (43.90%)
Iteration 300, loss = 1.5792
Got 413 / 1000 correct (41.30%)
Iteration 400, loss = 1.5824
Got 457 / 1000 correct (45.70%)
Iteration 500, loss = 1.6831
Got 495 / 1000 correct (49.50%)
Iteration 600, loss = 1.6586
Got 501 / 1000 correct (50.10%)
Iteration 700, loss = 1.5411
Got 497 / 1000 correct (49.70%)

```

## 5 Part III: Keras Model Subclassing API

Implementing a neural network using the low-level TensorFlow API is a good way to understand how TensorFlow works, but it's a little inconvenient - we had to manually keep track of all Tensors holding learnable parameters. This was fine for a small network, but could quickly become unweildy for a large complex model.

Fortunately TensorFlow 2.0 provides higher-level APIs such as `tf.keras` which make it easy to build models out of modular, object-oriented layers. Further, TensorFlow 2.0 uses eager execution that evaluates operations immediately, without explicitly constructing any computational graphs. This makes it easy to write and debug models, and reduces the boilerplate code.

In this part of the notebook we will define neural network models using the `tf.keras.Model` API. To implement your own model, you need to do the following:

1. Define a new class which subclasses `tf.keras.Model`. Give your class an intuitive name that describes it, like `TwoLayerFC` or `ThreeLayerConvNet`.
2. In the initializer `__init__()` for your new class, define all the layers you need as class attributes. The `tf.keras.layers` package provides many common neural-network layers, like `tf.keras.layers.Dense` for fully-connected layers and `tf.keras.layers.Conv2D`



for convolutional layers. Under the hood, these layers will construct Variable Tensors for any learnable parameters. **Warning:** Don't forget to call `super(YourModelName, self).__init__()` as the first line in your initializer!

3. Implement the `call()` method for your class; this implements the forward pass of your model, and defines the *connectivity* of your network. Layers defined in `__init__()` implement `__call__()` so they can be used as function objects that transform input Tensors into output Tensors. Don't define any new layers in `call()`; any layers you want to use in the forward pass should be defined in `__init__()`.

After you define your `tf.keras.Model` subclass, you can instantiate it and use it like the model functions from Part II.

### 5.0.1 Keras Model Subclassing API: Two-Layer Network

Here is a concrete example of using the `tf.keras.Model` API to define a two-layer network. There are a few new bits of API to be aware of here:

We use an `Initializer` object to set up the initial values of the learnable parameters of the layers; in particular `tf.initializers.VarianceScaling` gives behavior similar to the Kaiming initialization method we used in Part II. You can read more about it here: [https://www.tensorflow.org/versions/r2.0/api\\_docs/python/tf/initializers/VarianceScaling](https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/initializers/VarianceScaling)

We construct `tf.keras.layers.Dense` objects to represent the two fully-connected layers of the model. In addition to multiplying their input by a weight matrix and adding a bias vector, these layer can also apply a nonlinearity for you. For the first layer we specify a ReLU activation function by passing `activation='relu'` to the constructor; the second layer uses softmax activation function. Finally, we use `tf.keras.layers.Flatten` to flatten the output from the previous fully-connected layer.

```
[0]: class TwoLayerFC(tf.keras.Model):
    def __init__(self, hidden_size, num_classes):
        super(TwoLayerFC, self).__init__()
        initializer = tf.initializers.VarianceScaling(scale=2.0)
        self.fc1 = tf.keras.layers.Dense(hidden_size, activation='relu',
                                          kernel_initializer=initializer)
        self.fc2 = tf.keras.layers.Dense(num_classes, activation='softmax',
                                          kernel_initializer=initializer)
        self.flatten = tf.keras.layers.Flatten()

    def call(self, x, training=False):
        x = self.flatten(x)
        x = self.fc1(x)
        x = self.fc2(x)
        return x

def test_TwoLayerFC():
    """ A small unit test to exercise the TwoLayerFC model above. """
    input_size, hidden_size, num_classes = 50, 42, 10
    x = tf.zeros((64, input_size))
    model = TwoLayerFC(hidden_size, num_classes)
```

```

with tf.device(device):
    scores = model(x)
    print(scores.shape)

test_TwoLayerFC()

```

(64, 10)

## 5.0.2 Keras Model Subclassing API: Three-Layer ConvNet

Now it's your turn to implement a three-layer ConvNet using the `tf.keras.Model` API. Your model should have the same architecture used in Part II:

1. Convolutional layer with 5 x 5 kernels, with zero-padding of 2
2. ReLU nonlinearity
3. Convolutional layer with 3 x 3 kernels, with zero-padding of 1
4. ReLU nonlinearity
5. Fully-connected layer to give class scores
6. Softmax nonlinearity

You should initialize the weights of your network using the same initialization method as was used in the two-layer network above.

**Hint:** Refer to the documentation for `tf.keras.layers.Conv2D` and `tf.keras.layers.Dense`:  
[https://www.tensorflow.org/versions/r2.0/api\\_docs/python/tf/keras/layers/Conv2D](https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/keras/layers/Conv2D)  
[https://www.tensorflow.org/versions/r2.0/api\\_docs/python/tf/keras/layers/Dense](https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/keras/layers/Dense)

```

[0]: class ThreeLayerConvNet(tf.keras.Model):
    def __init__(self, channel_1, channel_2, num_classes):
        super(ThreeLayerConvNet, self).__init__()

        # TODO: Implement the __init__ method for a three-layer ConvNet. You
        # should instantiate layer objects to be used in the forward pass.

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

        initializer = tf.initializers.VarianceScaling(scale=2.0)
        self.conv1 = tf.keras.layers.Conv2D(channel_1, [5,5], [1,1],
        padding='valid',
                                           kernel_initializer=initializer,
                                           activation=tf.nn.relu)

        self.conv2 = tf.keras.layers.Conv2D(channel_2, [3,3], [1,1],
        padding='valid',
                                           kernel_initializer=initializer,
                                           activation=tf.nn.relu)

```

```

        self.fc = tf.keras.layers.Dense(num_classes,
→kernel_initializer=initializer)
        self.flatten = tf.keras.layers.Flatten()
        self.softmax = tf.keras.layers.Softmax()

        # *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
→#
        #
→#
→#
→#
        def call(self, x, training=False):
            scores = None
→#
→#
            # TODO: Implement the forward pass for a three-layer ConvNet. You
→#
            # should use the layer objects defined in the __init__ method.
→#
→#
→#
            # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

            padding = tf.constant([[0,0],[2,2],[2,2],[0,0]])
            x = tf.pad(x, padding, 'CONSTANT')
            x = self.conv1(x)
            padding = tf.constant([[0,0],[1,1],[1,1],[0,0]])
            x = tf.pad(x, padding, 'CONSTANT')
            x = self.conv2(x)
            x = self.flatten(x)
            x = self.fc(x)
            scores = self.softmax(x)

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

```

Once you complete the implementation of the ThreeLayerConvNet above you can run the following to ensure that your implementation does not crash and produces outputs of the expected

shape.

```
[0]: def test_ThreeLayerConvNet():
    channel_1, channel_2, num_classes = 12, 8, 10
    model = ThreeLayerConvNet(channel_1, channel_2, num_classes)
    with tf.device(device):
        x = tf.zeros((64, 3, 32, 32))
        scores = model(x)
        print(scores.shape)

test_ThreeLayerConvNet()
```

(64, 10)

### 5.0.3 Keras Model Subclassing API: Eager Training

While keras models have a builtin training loop (using the `model.fit`), sometimes you need more customization. Here's an example, of a training loop implemented with eager execution.

In particular, notice `tf.GradientTape`. Automatic differentiation is used in the backend for implementing backpropagation in frameworks like TensorFlow. During eager execution, `tf.GradientTape` is used to trace operations for computing gradients later. A particular `tf.GradientTape` can only compute one gradient; subsequent calls to `tape` will throw a runtime error.

TensorFlow 2.0 ships with easy-to-use built-in metrics under `tf.keras.metrics` module. Each metric is an object, and we can use `update_state()` to add observations and `reset_state()` to clear all observations. We can get the current result of a metric by calling `result()` on the metric object.

```
[0]: def train_part34(model_init_fn, optimizer_init_fn, num_epochs=1,
    ↪is_training=False):
    """
    Simple training loop for use with models defined using tf.keras. It trains
    a model for one epoch on the CIFAR-10 training set and periodically checks
    accuracy on the CIFAR-10 validation set.

    Inputs:
    - model_init_fn: A function that takes no parameters; when called it
      constructs the model we want to train: model = model_init_fn()
    - optimizer_init_fn: A function which takes no parameters; when called it
      constructs the Optimizer object we will use to optimize the model:
      optimizer = optimizer_init_fn()
    - num_epochs: The number of epochs to train for

    Returns: Nothing, but prints progress during training
    """
    with tf.device(device):

        # Compute the loss like we did in Part II
        loss_fn = tf.keras.losses.SparseCategoricalCrossentropy()
```

```

model = model_init_fn()
optimizer = optimizer_init_fn()

train_loss = tf.keras.metrics.Mean(name='train_loss')
train_accuracy = tf.keras.metrics.
→SparseCategoricalAccuracy(name='train_accuracy')

val_loss = tf.keras.metrics.Mean(name='val_loss')
val_accuracy = tf.keras.metrics.
→SparseCategoricalAccuracy(name='val_accuracy')

t = 0
for epoch in range(num_epochs):

    # Reset the metrics - https://www.tensorflow.org/alpha/guide/
→migration_guide#new-style-metrics
    train_loss.reset_states()
    train_accuracy.reset_states()

    for x_np, y_np in train_dset:
        with tf.GradientTape() as tape:

            # Use the model function to build the forward pass.
            scores = model(x_np, training=is_training)
            loss = loss_fn(y_np, scores)

            gradients = tape.gradient(loss, model.trainable_variables)
            optimizer.apply_gradients(zip(gradients, model.
→trainable_variables))

            # Update the metrics
            train_loss.update_state(loss)
            train_accuracy.update_state(y_np, scores)

            if t % print_every == 0:
                val_loss.reset_states()
                val_accuracy.reset_states()
                for test_x, test_y in val_dset:
                    # During validation at end of epoch, training set
→to False

                    prediction = model(test_x, training=False)
                    t_loss = loss_fn(test_y, prediction)

                    val_loss.update_state(t_loss)
                    val_accuracy.update_state(test_y, prediction)

```

```

        template = 'Iteration {}, Epoch {}, Loss: {}, Accuracy: {}
        → {}, Val Loss: {}, Val Accuracy: {}'
        print (template.format(t, epoch+1,
                                train_loss.result(),
                                train_accuracy.result()*100,
                                val_loss.result(),
                                val_accuracy.result()*100))

        t += 1

```

#### 5.0.4 Keras Model Subclassing API: Train a Two-Layer Network

We can now use the tools defined above to train a two-layer network on CIFAR-10. We define the `model_init_fn` and `optimizer_init_fn` that construct the model and optimizer respectively when called. Here we want to train the model using stochastic gradient descent with no momentum, so we construct a `tf.keras.optimizers.SGD` function; you can [read about it here](#).

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

```

[0]: hidden_size, num_classes = 4000, 10
    learning_rate = 1e-2

    def model_init_fn():
        return TwoLayerFC(hidden_size, num_classes)

    def optimizer_init_fn():
        return tf.keras.optimizers.SGD(learning_rate=learning_rate)

    train_part34(model_init_fn, optimizer_init_fn)

```

```

Iteration 0, Epoch 1, Loss: 2.8891260623931885, Accuracy: 14.0625, Val Loss:
3.0216798782348633, Val Accuracy: 11.90000057220459
Iteration 100, Epoch 1, Loss: 2.2201895713806152, Accuracy: 28.434404373168945,
Val Loss: 1.8994783163070679, Val Accuracy: 38.0
Iteration 200, Epoch 1, Loss: 2.062427043914795, Accuracy: 32.40827178955078,
Val Loss: 1.8448724746704102, Val Accuracy: 39.099998474121094
Iteration 300, Epoch 1, Loss: 1.9858424663543701, Accuracy: 34.48920440673828,
Val Loss: 1.9095661640167236, Val Accuracy: 36.39999771118164
Iteration 400, Epoch 1, Loss: 1.9196314811706543, Accuracy: 36.218048095703125,
Val Loss: 1.7207499742507935, Val Accuracy: 42.599998474121094
Iteration 500, Epoch 1, Loss: 1.8778129816055298, Accuracy: 37.29104232788086,
Val Loss: 1.670654058456421, Val Accuracy: 42.89999771118164
Iteration 600, Epoch 1, Loss: 1.847782015800476, Accuracy: 38.316349029541016,
Val Loss: 1.678385853767395, Val Accuracy: 42.20000076293945
Iteration 700, Epoch 1, Loss: 1.8244125843048096, Accuracy: 38.88641357421875,
Val Loss: 1.6426002979278564, Val Accuracy: 44.29999923706055

```

### 5.0.5 Keras Model Subclassing API: Train a Three-Layer ConvNet

Here you should use the tools we've defined above to train a three-layer ConvNet on CIFAR-10. Your ConvNet should use 32 filters in the first convolutional layer and 16 filters in the second layer.

To train the model you should use gradient descent with Nesterov momentum 0.9.

**HINT:** [https://www.tensorflow.org/versions/r2.0/api\\_docs/python/tf/optimizers/SGD](https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/optimizers/SGD)

You don't need to perform any hyperparameter tuning, but you should achieve validation accuracies above 50% after training for one epoch.

```
[0]: learning_rate = 3e-3
channel_1, channel_2, num_classes = 32, 16, 10

def model_init_fn():
    model = None

    # TODO: Complete the implementation of model_fn.

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

    model = ThreeLayerConvNet(channel_1, channel_2, num_classes)

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

    # END OF YOUR CODE

    return model

def optimizer_init_fn():
    optimizer = None

    # TODO: Complete the implementation of model_fn.

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

    optimizer = tf.keras.optimizers.SGD(learning_rate, 0.9, nesterov=True)

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

```

└─
→#####
#                                END OF YOUR CODE                                └─
→#
└─
→#####
return optimizer

train_part34(model_init_fn, optimizer_init_fn)

```

```

Iteration 0, Epoch 1, Loss: 2.8280727863311768, Accuracy: 10.9375, Val Loss:
4.871011734008789, Val Accuracy: 8.600000381469727
Iteration 100, Epoch 1, Loss: 1.974219799041748, Accuracy: 33.41584014892578,
Val Loss: 1.6492583751678467, Val Accuracy: 43.20000076293945
Iteration 200, Epoch 1, Loss: 1.7758405208587646, Accuracy: 38.813743591308594,
Val Loss: 1.498644232749939, Val Accuracy: 47.20000076293945
Iteration 300, Epoch 1, Loss: 1.6737589836120605, Accuracy: 41.90199279785156,
Val Loss: 1.4523096084594727, Val Accuracy: 50.19999694824219
Iteration 400, Epoch 1, Loss: 1.597639799118042, Accuracy: 44.201995849609375,
Val Loss: 1.3690105676651, Val Accuracy: 50.80000305175781
Iteration 500, Epoch 1, Loss: 1.5436161756515503, Accuracy: 45.8925895690918,
Val Loss: 1.3240282535552979, Val Accuracy: 52.499996185302734
Iteration 600, Epoch 1, Loss: 1.507643461227417, Accuracy: 46.99979019165039,
Val Loss: 1.2929513454437256, Val Accuracy: 53.79999923706055
Iteration 700, Epoch 1, Loss: 1.4768180847167969, Accuracy: 48.00953674316406,
Val Loss: 1.2605247497558594, Val Accuracy: 56.5

```

## 6 Part IV: Keras Sequential API

In Part III we introduced the `tf.keras.Model` API, which allows you to define models with any number of learnable layers and with arbitrary connectivity between layers.

However for many models you don't need such flexibility - a lot of models can be expressed as a sequential stack of layers, with the output of each layer fed to the next layer as input. If your model fits this pattern, then there is an even easier way to define your model: using `tf.keras.Sequential`. You don't need to write any custom classes; you simply call the `tf.keras.Sequential` constructor with a list containing a sequence of layer objects.

One complication with `tf.keras.Sequential` is that you must define the shape of the input to the model by passing a value to the `input_shape` of the first layer in your model.

### 6.0.1 Keras Sequential API: Two-Layer Network

In this subsection, we will rewrite the two-layer fully-connected network using `tf.keras.Sequential`, and train it using the training loop defined above.

You don't need to perform any hyperparameter tuning here, but you should see validation accuracies above 40% after training for one epoch.



```
[0]: learning_rate = 1e-2

def model_init_fn():
    input_shape = (32, 32, 3)
    hidden_layer_size, num_classes = 4000, 10
    initializer = tf.initializers.VarianceScaling(scale=2.0)
    layers = [
        tf.keras.layers.Flatten(input_shape=input_shape),
        tf.keras.layers.Dense(hidden_layer_size, activation='relu',
                               kernel_initializer=initializer),
        tf.keras.layers.Dense(num_classes, activation='softmax',
                               kernel_initializer=initializer),
    ]
    model = tf.keras.Sequential(layers)
    return model

def optimizer_init_fn():
    return tf.keras.optimizers.SGD(learning_rate=learning_rate)

train_part34(model_init_fn, optimizer_init_fn)
```

```
Iteration 0, Epoch 1, Loss: 3.3227744102478027, Accuracy: 7.8125, Val Loss:
2.9715607166290283, Val Accuracy: 10.899999618530273
Iteration 100, Epoch 1, Loss: 2.2472212314605713, Accuracy: 28.18688201904297,
Val Loss: 1.8935211896896362, Val Accuracy: 38.29999923706055
Iteration 200, Epoch 1, Loss: 2.0815584659576416, Accuracy: 31.96517562866211,
Val Loss: 1.8076941967010498, Val Accuracy: 39.5
Iteration 300, Epoch 1, Loss: 2.0035083293914795, Accuracy: 34.01681900024414,
Val Loss: 1.867803692817688, Val Accuracy: 36.70000076293945
Iteration 400, Epoch 1, Loss: 1.9359807968139648, Accuracy: 35.64136505126953,
Val Loss: 1.713821530342102, Val Accuracy: 43.39999771118164
Iteration 500, Epoch 1, Loss: 1.8880256414413452, Accuracy: 36.82634735107422,
Val Loss: 1.6455496549606323, Val Accuracy: 43.099998474121094
Iteration 600, Epoch 1, Loss: 1.8585952520370483, Accuracy: 37.74958419799805,
Val Loss: 1.6740604639053345, Val Accuracy: 42.39999771118164
Iteration 700, Epoch 1, Loss: 1.831824779510498, Accuracy: 38.48297119140625,
Val Loss: 1.5924526453018188, Val Accuracy: 45.29999923706055
```

## 6.0.2 Abstracting Away the Training Loop

In the previous examples, we used a customised training loop to train models (e.g. `train_part34`). Writing your own training loop is only required if you need more flexibility and control during training your model. Alternately, you can also use built-in APIs like `tf.keras.Model.fit()` and `tf.keras.Model.evaluate` to train and evaluate a model. Also remember to configure your model for training by calling `tf.keras.Model.compile`.

You don't need to perform any hyperparameter tuning here, but you should see validation and test accuracies above 42% after training for one epoch.

```
[0]: model = model_init_fn()
model.compile(optimizer=tf.keras.optimizers.SGD(learning_rate=learning_rate),
              loss='sparse_categorical_crossentropy',
              metrics=[tf.keras.metrics.sparse_categorical_accuracy])
model.fit(X_train, y_train, batch_size=64, epochs=1, validation_data=(X_val,
→y_val))
model.evaluate(X_test, y_test)
```

```
766/766 [=====] - 3s 4ms/step - loss: 1.8157 -
sparse_categorical_accuracy: 0.3869 - val_loss: 1.6898 -
val_sparse_categorical_accuracy: 0.4300
313/313 [=====] - 1s 2ms/step - loss: 1.6684 -
sparse_categorical_accuracy: 0.4287
```

```
[0]: [1.6684269905090332, 0.4287000000476837]
```

### 6.0.3 Keras Sequential API: Three-Layer ConvNet

Here you should use `tf.keras.Sequential` to reimplement the same three-layer ConvNet architecture used in Part II and Part III. As a reminder, your model should have the following architecture:

1. Convolutional layer with 32 5x5 kernels, using zero padding of 2
2. ReLU nonlinearity
3. Convolutional layer with 16 3x3 kernels, using zero padding of 1
4. ReLU nonlinearity
5. Fully-connected layer giving class scores
6. Softmax nonlinearity

You should initialize the weights of the model using a `tf.initializers.VarianceScaling` as above.

You should train the model using Nesterov momentum 0.9.

You don't need to perform any hyperparameter search, but you should achieve accuracy above 45% after training for one epoch.

```
[0]: def model_init_fn():
    model = None

    →#####
    # TODO: Construct a three-layer ConvNet using tf.keras.Sequential.
    →#

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

    input_shape = (32,32,3)
    channel_1, channel_2, num_classes = 32, 16, 10
    initializer = tf.initializers.VarianceScaling(scale=2.0)
```

```

layers = [
    tf.keras.layers.InputLayer(input_shape=input_shape),
    tf.keras.layers.Conv2D(channel_1, [5,5], [1,1], padding='same',
                           kernel_initializer=initializer,
                           activation=tf.nn.relu),
    tf.keras.layers.Conv2D(channel_2, [3,3], [1,1], padding='same',
                           kernel_initializer=initializer,
                           activation=tf.nn.relu),
    tf.keras.layers.Flatten(),
    tf.keras.layers.Dense(num_classes, kernel_initializer=initializer),
    tf.keras.layers.Softmax()
]
model = tf.keras.Sequential(layers)

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

↳
↳#####
return model

learning_rate = 5e-4
def optimizer_init_fn():
    optimizer = None

    ↳#####
    # TODO: Complete the implementation of model_fn.                                     ↳
    ↳#

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

    optimizer = tf.keras.optimizers.SGD(learning_rate, momentum=0.9,↳
    ↳nesterov=True)

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

    ↳#####
    #                                     END OF YOUR CODE                                     ↳
    ↳#

    ↳
    ↳#####
    return optimizer

```

```
train_part34(model_init_fn, optimizer_init_fn)
```

```
Iteration 0, Epoch 1, Loss: 3.3954334259033203, Accuracy: 3.125, Val Loss:
2.895132541656494, Val Accuracy: 9.600000381469727
Iteration 100, Epoch 1, Loss: 1.9908231496810913, Accuracy: 29.563735961914062,
Val Loss: 1.7454599142074585, Val Accuracy: 40.099998474121094
Iteration 200, Epoch 1, Loss: 1.8512924909591675, Accuracy: 35.26896667480469,
Val Loss: 1.6331830024719238, Val Accuracy: 45.0
Iteration 300, Epoch 1, Loss: 1.779958724975586, Accuracy: 37.69725799560547,
Val Loss: 1.5917105674743652, Val Accuracy: 45.400001525878906
Iteration 400, Epoch 1, Loss: 1.7156802415847778, Accuracy: 39.736595153808594,
Val Loss: 1.5262717008590698, Val Accuracy: 48.89999771118164
Iteration 500, Epoch 1, Loss: 1.6727784872055054, Accuracy: 41.311126708984375,
Val Loss: 1.481580138206482, Val Accuracy: 49.900001525878906
Iteration 600, Epoch 1, Loss: 1.6445434093475342, Accuracy: 42.320091247558594,
Val Loss: 1.4551396369934082, Val Accuracy: 50.19999694824219
Iteration 700, Epoch 1, Loss: 1.6195082664489746, Accuracy: 43.20390701293945,
Val Loss: 1.4160679578781128, Val Accuracy: 51.400001525878906
```

We will also train this model with the built-in training loop APIs provided by TensorFlow.

```
[0]: model = model_init_fn()
model.compile(optimizer='sgd',
              loss='sparse_categorical_crossentropy',
              metrics=[tf.keras.metrics.sparse_categorical_accuracy])
model.fit(X_train, y_train, batch_size=64, epochs=1, validation_data=(X_val,
↪y_val))
model.evaluate(X_test, y_test)
```

```
766/766 [=====] - 3s 4ms/step - loss: 1.5475 -
sparse_categorical_accuracy: 0.4498 - val_loss: 1.4093 -
val_sparse_categorical_accuracy: 0.5020
313/313 [=====] - 1s 2ms/step - loss: 1.4161 -
sparse_categorical_accuracy: 0.5003
```

```
[0]: [1.4160820245742798, 0.5002999901771545]
```

## 6.1 Part IV: Functional API

### 6.1.1 Demonstration with a Two-Layer Network

In the previous section, we saw how we can use `tf.keras.Sequential` to stack layers to quickly build simple models. But this comes at the cost of losing flexibility.

Often we will have to write complex models that have non-sequential data flows: a layer can have **multiple inputs and/or outputs**, such as stacking the output of 2 previous layers together to feed as input to a third! (Some examples are residual connections and dense blocks.)

In such cases, we can use Keras functional API to write models with complex topologies such as:

1. Multi-input models
2. Multi-output models
3. Models with shared layers (the same layer called several times)
4. Models with non-sequential data flows (e.g. residual connections)

Writing a model with Functional API requires us to create a `tf.keras.Model` instance and explicitly write input tensors and output tensors for this model.

```
[0]: def two_layer_fc_functional(input_shape, hidden_size, num_classes):
    initializer = tf.initializers.VarianceScaling(scale=2.0)
    inputs = tf.keras.Input(shape=input_shape)
    flattened_inputs = tf.keras.layers.Flatten()(inputs)
    fc1_output = tf.keras.layers.Dense(hidden_size, activation='relu',
                                       ␣
                                       ↪kernel_initializer=initializer)(flattened_inputs)
    scores = tf.keras.layers.Dense(num_classes, activation='softmax',
                                   kernel_initializer=initializer)(fc1_output)

    # Instantiate the model given inputs and outputs.
    model = tf.keras.Model(inputs=inputs, outputs=scores)
    return model

def test_two_layer_fc_functional():
    """ A small unit test to exercise the TwoLayerFC model above. """
    input_size, hidden_size, num_classes = 50, 42, 10
    input_shape = (50,)

    x = tf.zeros((64, input_size))
    model = two_layer_fc_functional(input_shape, hidden_size, num_classes)

    with tf.device(device):
        scores = model(x)
        print(scores.shape)

test_two_layer_fc_functional()
```

(64, 10)

### 6.1.2 Keras Functional API: Train a Two-Layer Network

You can now train this two-layer network constructed using the functional API.

You don't need to perform any hyperparameter tuning here, but you should see validation accuracies above 40% after training for one epoch.

```
[0]: input_shape = (32, 32, 3)
    hidden_size, num_classes = 4000, 10
    learning_rate = 1e-2

    def model_init_fn():
```

```

    return two_layer_fc_functional(input_shape, hidden_size, num_classes)

def optimizer_init_fn():
    return tf.keras.optimizers.SGD(learning_rate=learning_rate)

train_part34(model_init_fn, optimizer_init_fn)

```

```

Iteration 0, Epoch 1, Loss: 3.147264003753662, Accuracy: 10.9375, Val Loss:
2.887949228286743, Val Accuracy: 12.399999618530273
Iteration 100, Epoch 1, Loss: 2.2481014728546143, Accuracy: 27.583539962768555,
Val Loss: 1.9231351613998413, Val Accuracy: 38.0
Iteration 200, Epoch 1, Loss: 2.0761499404907227, Accuracy: 32.12064743041992,
Val Loss: 1.8733704090118408, Val Accuracy: 41.29999923706055
Iteration 300, Epoch 1, Loss: 2.0036330223083496, Accuracy: 33.87146759033203,
Val Loss: 1.8841179609298706, Val Accuracy: 37.0
Iteration 400, Epoch 1, Loss: 1.9357917308807373, Accuracy: 35.582916259765625,
Val Loss: 1.7401769161224365, Val Accuracy: 42.79999923706055
Iteration 500, Epoch 1, Loss: 1.892579555114746, Accuracy: 36.682884216308594,
Val Loss: 1.6923105716705322, Val Accuracy: 42.20000076293945
Iteration 600, Epoch 1, Loss: 1.8612903356552124, Accuracy: 37.661190032958984,
Val Loss: 1.71518874168396, Val Accuracy: 42.0
Iteration 700, Epoch 1, Loss: 1.8347382545471191, Accuracy: 38.41610336303711,
Val Loss: 1.673095464706421, Val Accuracy: 43.70000076293945

```

## 7 Part V: CIFAR-10 open-ended challenge

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

You should experiment with architectures, hyperparameters, loss functions, regularization, or anything else you can think of to train a model that achieves **at least 70%** accuracy on the **validation** set within 10 epochs. You can use the built-in train function, the `train_part34` function from above, or implement your own training loop.

Describe what you did at the end of the notebook.

### 7.0.1 Some things you can try:

- **Filter size:** Above we used 5x5 and 3x3; is this optimal?
- **Number of filters:** Above we used 16 and 32 filters. Would more or fewer do better?
- **Pooling:** We didn't use any pooling above. Would this improve the model?
- **Normalization:** Would your model be improved with batch normalization, layer normalization, group normalization, or some other normalization strategy?
- **Network architecture:** The ConvNet above has only three layers of trainable parameters. Would a deeper model do better?
- **Global average pooling:** Instead of flattening after the final convolutional layer, would global average pooling do better? This strategy is used for example in Google's Inception network and in Residual Networks.
- **Regularization:** Would some kind of regularization improve performance? Maybe weight decay or dropout?

### 7.0.2 NOTE: Batch Normalization / Dropout

If you are using Batch Normalization and Dropout, remember to pass `is_training=True` if you use the `train_part34()` function. BatchNorm and Dropout layers have different behaviors at training and inference time. `training` is a specific keyword argument reserved for this purpose in any `tf.keras.Model`'s `call()` function. Read more about this here : [https://www.tensorflow.org/versions/r2.0/api\\_docs/python/tf/keras/layers/BatchNormalization#methods](https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/keras/layers/BatchNormalization#methods) [https://www.tensorflow.org/versions/r2.0/api\\_docs/python/tf/keras/layers/Dropout#methods](https://www.tensorflow.org/versions/r2.0/api_docs/python/tf/keras/layers/Dropout#methods)

### 7.0.3 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.

### 7.0.4 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](#)

### 7.0.5 Have fun and happy training!

```
[0]: class CustomConvNet(tf.keras.Model):
      def __init__(self):
          super(CustomConvNet, self).__init__()
          ␣
      → #####
          # TODO: Construct a model that performs well on CIFAR-10 ␣
      → #
```

```

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

    num_classes = 10
    input_shape = (32, 32, 3)
    outchannel_1, outchannel_2, outchannel_3, num_classes = 128, 256, 512, □
→10
    filter_size1, filter_size2, filter_size3 = 5, 3, 3
    initializer = tf.initializers.VarianceScaling(scale=1.0)

    #inputs = tf.keras.layers.InputLayer(input_shape=input_shape)
    self.conv1 = tf.keras.layers.Conv2D(outchannel_1, filter_size1, □
→padding='same', \
                                   kernel_initializer=initializer, \
                                   input_shape=input_shape)

    self.bn1 = tf.keras.layers.BatchNormalization()
    self.relu1 = tf.keras.layers.ReLU()
    self.drop1 = tf.keras.layers.Dropout(rate=0.2)
    self.pool1 = tf.keras.layers.MaxPool2D(2)

    self.conv2 = tf.keras.layers.Conv2D(outchannel_2, filter_size2, □
→padding='same', \
                                   kernel_initializer=initializer)

    self.bn2 = tf.keras.layers.BatchNormalization()
    self.relu2 = tf.keras.layers.ReLU()
    self.drop2 = tf.keras.layers.Dropout(rate=0.2)
    self.pool2 = tf.keras.layers.MaxPool2D(2)

    self.conv3 = tf.keras.layers.Conv2D(outchannel_3, filter_size3, □
→padding='same', \
                                   kernel_initializer=initializer)

    self.bn3 = tf.keras.layers.BatchNormalization()
    self.relu3 = tf.keras.layers.ReLU()
    self.drop3 = tf.keras.layers.Dropout(rate=0.2)
    #self.avg_pool = tf.keras.layers.AveragePooling2D(pool_size=(5, 5))
    self.flatten = tf.keras.layers.Flatten()

    self.fc1 = tf.keras.layers.Dense(512*4*4)
    self.fc2 = tf.keras.layers.Dense(10)
    self.scores = tf.keras.layers.Dense(num_classes, □
→kernel_initializer=initializer)
    self.softmax = tf.keras.layers.Softmax()

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

```



```

→ #####
#
#
→ #####

def call(self, input_tensor, training=False):
→ #####
# TODO: Construct a model that performs well on CIFAR-10
→ #
→ #####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

x = self.conv1(input_tensor)
x = self.bn1(x)
x = self.relu1(x)
#x = self.drop1(x)
x = self.pool1(x)

x = self.conv2(x)
x = self.bn2(x)
x = self.relu2(x)
#x = self.drop2(x)
x = self.pool2(x)

x = self.conv3(x)
x = self.bn3(x)
x = self.relu3(x)
#x = self.drop3(x)
#x = self.avg_pool(x)

x = self.drop3(x)
x = self.flatten(x)

x = self.fc1(x)
#x = self.fc2(x)

x = self.scores(x)
x = self.softmax(x)

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

```

```

#                                     END OF YOUR CODE
→ #
    □
→ #####
    return x

#device = '/device:GPU:0'
print_every = 700
num_epochs = 10

model = CustomConvNet()

def model_init_fn():
    return CustomConvNet()

def optimizer_init_fn():
    learning_rate = 1e-3
    return tf.keras.optimizers.SGD(learning_rate)

train_part34(model_init_fn, optimizer_init_fn, num_epochs=num_epochs, □
→ is_training=True)

```

Iteration 0, Epoch 1, Loss: 2.7738125324249268, Accuracy: 4.6875, Val Loss: 2.3566694259643555, Val Accuracy: 12.800000190734863

Iteration 700, Epoch 1, Loss: 1.5500390529632568, Accuracy: 44.96255111694336, Val Loss: 1.2871183156967163, Val Accuracy: 55.0

Iteration 1400, Epoch 2, Loss: 1.1807807683944702, Accuracy: 58.5014762878418, Val Loss: 1.131056547164917, Val Accuracy: 60.5

Iteration 2100, Epoch 3, Loss: 1.024848461151123, Accuracy: 64.32337188720703, Val Loss: 1.0728026628494263, Val Accuracy: 62.80000305175781

Iteration 2800, Epoch 4, Loss: 0.9252512454986572, Accuracy: 68.06970977783203, Val Loss: 1.0029585361480713, Val Accuracy: 64.60000610351562

Iteration 3500, Epoch 5, Loss: 0.8499141931533813, Accuracy: 70.73799133300781, Val Loss: 0.9862015247344971, Val Accuracy: 66.19999694824219

Iteration 4200, Epoch 6, Loss: 0.7865332365036011, Accuracy: 73.14689636230469, Val Loss: 0.9249181747436523, Val Accuracy: 69.5999984741211

Iteration 4900, Epoch 7, Loss: 0.7306836843490601, Accuracy: 75.0256118774414, Val Loss: 0.9044567346572876, Val Accuracy: 69.30000305175781

Iteration 5600, Epoch 8, Loss: 0.6855727434158325, Accuracy: 76.58211517333984, Val Loss: 0.8993682265281677, Val Accuracy: 68.69999694824219

Iteration 6300, Epoch 9, Loss: 0.6428729891777039, Accuracy: 78.06177520751953, Val Loss: 0.8766300082206726, Val Accuracy: 70.0

Iteration 7000, Epoch 10, Loss: 0.5973249673843384, Accuracy: 80.09637451171875, Val Loss: 0.8795017004013062, Val Accuracy: 70.20000457763672

## 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.

### Network Architecture:

$(CONV \rightarrow BatchNorm \rightarrow ReLu \rightarrow MaxPool) * 3 \rightarrow Dropout \rightarrow FC$

**Filter Sizes:** 5x5, 3x3, 3x3

**Number of Filters:** 128, 256, 512

**Pooling:** 2x2 Max-pooling

**Normalization:** Batch normalization

**Regularization:** Dropout

**Optimizer:** SGD

**Specifics:**

- The key here is that after the first CONV layer, we shrink the filter size in the next one. This enabled us to push the accuracy to 60% in the second epoch.
- We used Batchnorm because it enables training with larger learning rates, which leads to faster convergence and better generalization. Layernorm would probably have a similar effect; however, because it has been found to be less effective than batchnorm, the convergence wouldn't probably gain as much of a speedup.
- The final push was using Dropout right before the Fully Connected layer which got us an additional 2% (from 68.14% to 70.20%) on the model accuracy. Dropout helps prevent overfitting.
- We used Softmax loss as our loss function, which proved to be more effective than SVM.