

FullyConnectedNets

November 4, 2022

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

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'cs231n/assignments/assignment1/'
FOLDERNAME = 'enpm809K/assignments/assignment2/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

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

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

Drive already mounted at /content/drive; to attempt to forcibly remount, call `drive.mount("/content/drive", force_remount=True)`.

/content/drive/My Drive/enpm809K/assignments/assignment2/cs231n/datasets
/content/drive/My Drive/enpm809K/assignments/assignment2

1 Multi-Layer Fully Connected Network

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

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

Implement the network initialization, forward pass, and backward pass. Throughout this assignment, you will be implementing layers in `cs231n/layers.py`. You can re-use your implementations for `affine_forward`, `affine_backward`, `relu_forward`, `relu_backward`, and `softmax_loss` from

Assignment 1. For right now, don't worry about implementing dropout or batch/layer normalization yet, as you will add those features later.

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

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

%load_ext autoreload
%autoreload 2

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

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

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

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

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

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

1.1 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. This is a good way to see if the initial losses seem reasonable.

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

```
[ ]: 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(f"{name} relative error: {rel_error(grad_num, grads[name])}")
```

```
Running check with reg = 0
[15, 20, 30, 10]
Initial loss: 2.3004790897684924
W1 relative error: 2.422780825313861e-07
W2 relative error: 0.00020570303548294835
W3 relative error: 7.875020476559293e-07
b1 relative error: 3.5730097461013185e-09
b2 relative error: 2.085654276112763e-09
b3 relative error: 1.893955861655706e-10
Running check with reg = 3.14
[15, 20, 30, 10]
Initial loss: 7.052114776533016
W1 relative error: 1.409028728052923e-08
```

W2 relative error: 6.86942277940646e-08
W3 relative error: 2.131129859578198e-08
b1 relative error: 1.4752427965311745e-08
b2 relative error: 1.8805361674745492e-09
b3 relative error: 2.378772438198909e-10

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

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

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

weight_scale = 1e-1 # Experiment with this!
learning_rate = 2e-4 # Experiment with this!
model = FullyConnectedNet(
    [100, 100],
    weight_scale=weight_scale,
    dtype=np.float64
)
solver = Solver(
    model,
    small_data,
    print_every=10,
    num_epochs=20,
    batch_size=25,
    update_rule="sgd",
    optim_config={"learning_rate": learning_rate},
)
solver.train()

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

```
[3072, 100, 100, 10]
```

/content/drive/My Drive/enpm809K/assignments/assignment2/cs231n/layers.py:168:

RuntimeWarning: divide by zero encountered in log

```
loss = -np.sum(np.log(sy) / np.sum(P, axis=1))
```

(Iteration 1 / 40) loss: inf

(Epoch 0 / 20) train acc: 0.080000; val_acc: 0.113000

(Epoch 1 / 20) train acc: 0.200000; val_acc: 0.117000

(Epoch 2 / 20) train acc: 0.340000; val_acc: 0.138000

(Epoch 3 / 20) train acc: 0.440000; val_acc: 0.151000

(Epoch 4 / 20) train acc: 0.520000; val_acc: 0.148000

(Epoch 5 / 20) train acc: 0.640000; val_acc: 0.164000

(Iteration 11 / 40) loss: 13.361355

(Epoch 6 / 20) train acc: 0.700000; val_acc: 0.159000

(Epoch 7 / 20) train acc: 0.840000; val_acc: 0.159000

(Epoch 8 / 20) train acc: 0.880000; val_acc: 0.158000

(Epoch 9 / 20) train acc: 0.920000; val_acc: 0.155000

(Epoch 10 / 20) train acc: 0.940000; val_acc: 0.158000

(Iteration 21 / 40) loss: 4.491945

(Epoch 11 / 20) train acc: 0.960000; val_acc: 0.158000

(Epoch 12 / 20) train acc: 0.980000; val_acc: 0.163000

(Epoch 13 / 20) train acc: 1.000000; val_acc: 0.158000

(Epoch 14 / 20) train acc: 1.000000; val_acc: 0.158000

(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.158000

(Iteration 31 / 40) loss: 0.000208

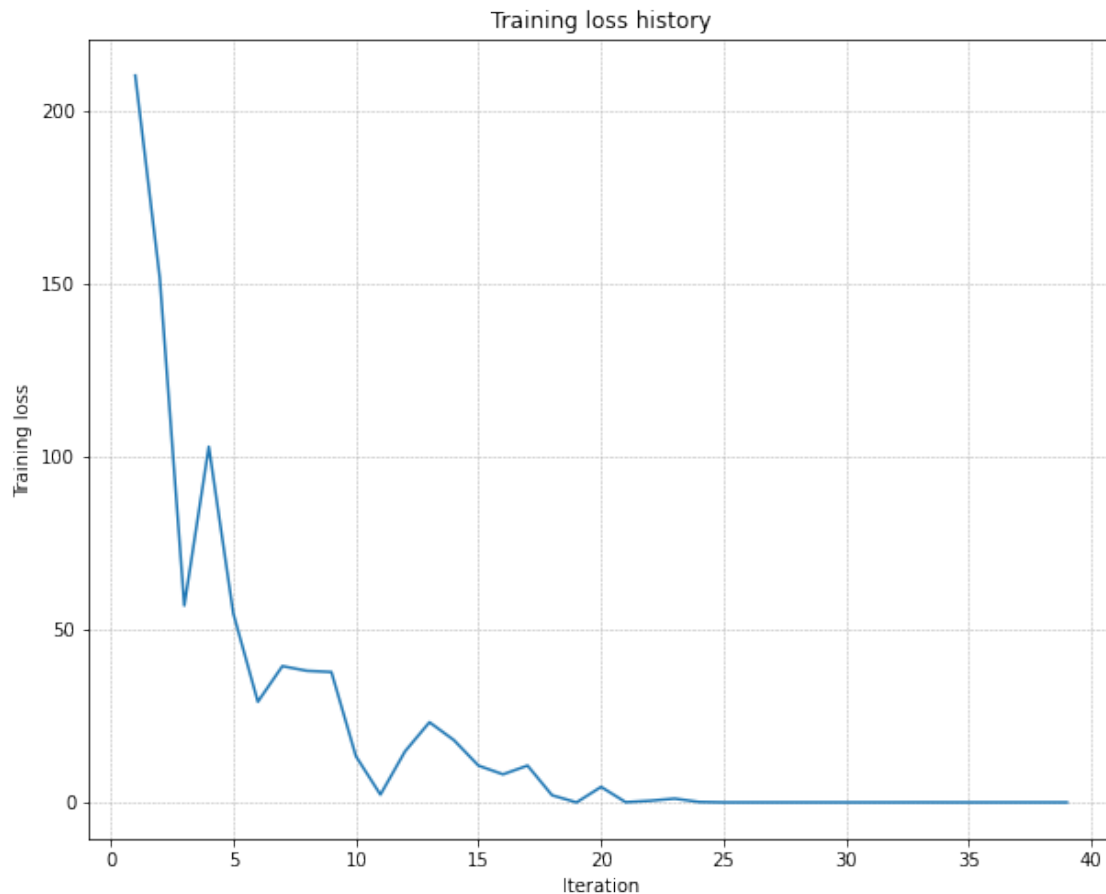
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.158000

(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.160000

(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.160000

(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.160000

(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.160000



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

```
[ ]: # 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 = 1e-1 # Experiment with this!  
model = FullyConnectedNet(  
    [100, 100, 100, 100],
```

```

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

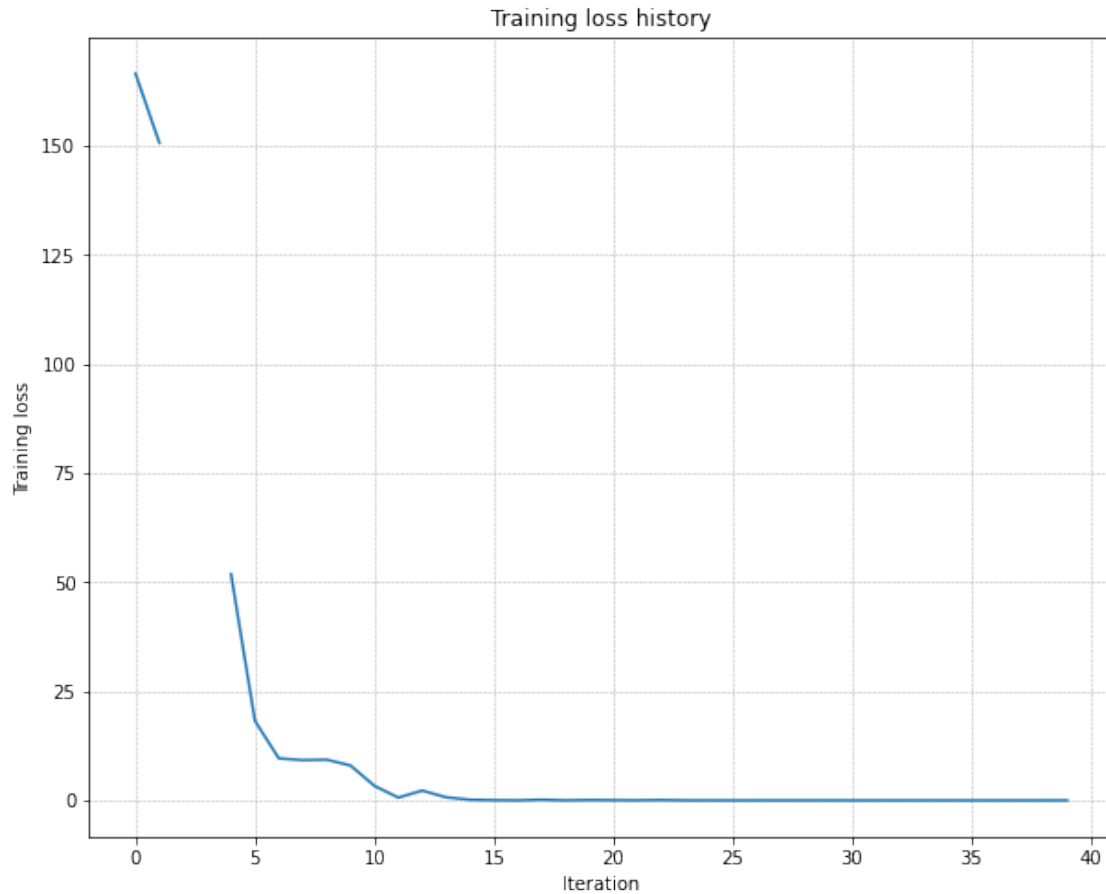
    plt.plot(solver.loss_history)
    plt.title('Training loss history')
    plt.xlabel('Iteration')
    plt.ylabel('Training loss')
    plt.grid(linestyle='--', linewidth=0.5)
    plt.show()

```

```

[3072, 100, 100, 100, 100, 10]
(Iteration 1 / 40) loss: 166.501707
(Epoch 0 / 20) train acc: 0.100000; val_acc: 0.107000
(Epoch 1 / 20) train acc: 0.320000; val_acc: 0.101000
(Epoch 2 / 20) train acc: 0.160000; val_acc: 0.122000
(Epoch 3 / 20) train acc: 0.380000; val_acc: 0.106000
(Epoch 4 / 20) train acc: 0.520000; val_acc: 0.111000
(Epoch 5 / 20) train acc: 0.760000; val_acc: 0.113000
(Iteration 11 / 40) loss: 3.343141
(Epoch 6 / 20) train acc: 0.840000; val_acc: 0.122000
(Epoch 7 / 20) train acc: 0.920000; val_acc: 0.113000
(Epoch 8 / 20) train acc: 0.940000; val_acc: 0.125000
(Epoch 9 / 20) train acc: 0.960000; val_acc: 0.125000
(Epoch 10 / 20) train acc: 0.980000; val_acc: 0.121000
(Iteration 21 / 40) loss: 0.039138
(Epoch 11 / 20) train acc: 0.980000; val_acc: 0.123000
(Epoch 12 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 13 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 14 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 15 / 20) train acc: 1.000000; val_acc: 0.121000
(Iteration 31 / 40) loss: 0.000644
(Epoch 16 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 17 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 18 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 19 / 20) train acc: 1.000000; val_acc: 0.121000
(Epoch 20 / 20) train acc: 1.000000; val_acc: 0.121000

```



1.2 Inline Question 1:

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

1.3 Answer:

Five layers is more difficult to train, and it's also more sensitive to the initialization scale. First of all, five layers have more parameters to tune, and it's easier to overfit so it's harder to train. The reason of why the deeper the architecture is the more sensitive to the initialization scale is because of the vanish gradient or exploding gradient. When doing backpropagation some of the nonlinear layer will get 0 gradient like too large or too small in sigmoid and negative part in relu, those will cause the networks learning slow or not converge to the minimum. So if the architecture is deeper than is have higher possibility for vanish gradient or exploding gradient happen.

2 Update rules

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

2.1 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 $e-8$.

```
[ ]: 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.

```

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

solvers = {}

for update_rule in ['sgd', 'sgd_momentum']:
    print('Running with ', update_rule)
    model = FullyConnectedNet(
        [100, 100, 100, 100, 100],
        weight_scale=5e-2
    )

    solver = Solver(
        model,
        small_data,
        num_epochs=5,
        batch_size=100,
        update_rule=update_rule,
        optim_config={'learning_rate': 5e-3},
        verbose=True,
    )
    solvers[update_rule] = solver
    solver.train()

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

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

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

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

```

```
plt.show()
```

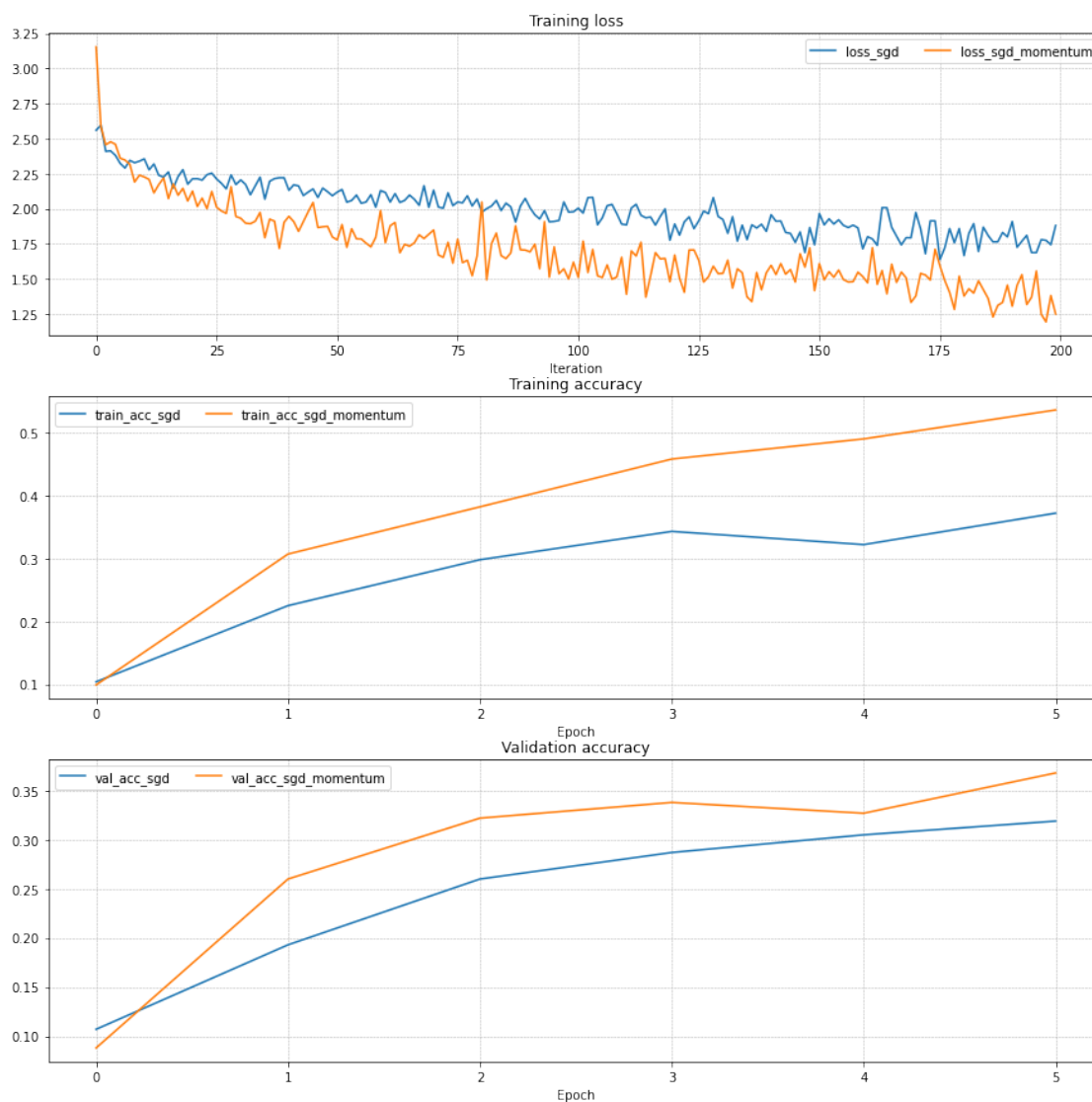
Running with sgd

```
[3072, 100, 100, 100, 100, 100, 10]
(Iteration 1 / 200) loss: 2.559978
(Epoch 0 / 5) train acc: 0.104000; val_acc: 0.107000
(Iteration 11 / 200) loss: 2.356069
(Iteration 21 / 200) loss: 2.214091
(Iteration 31 / 200) loss: 2.205928
(Epoch 1 / 5) train acc: 0.225000; val_acc: 0.193000
(Iteration 41 / 200) loss: 2.132095
(Iteration 51 / 200) loss: 2.118950
(Iteration 61 / 200) loss: 2.116443
(Iteration 71 / 200) loss: 2.132549
(Epoch 2 / 5) train acc: 0.298000; val_acc: 0.260000
(Iteration 81 / 200) loss: 1.977227
(Iteration 91 / 200) loss: 2.007528
(Iteration 101 / 200) loss: 2.004762
(Iteration 111 / 200) loss: 1.885342
(Epoch 3 / 5) train acc: 0.343000; val_acc: 0.287000
(Iteration 121 / 200) loss: 1.891516
(Iteration 131 / 200) loss: 1.923677
(Iteration 141 / 200) loss: 1.957743
(Iteration 151 / 200) loss: 1.966736
(Epoch 4 / 5) train acc: 0.322000; val_acc: 0.305000
(Iteration 161 / 200) loss: 1.801483
(Iteration 171 / 200) loss: 1.973780
(Iteration 181 / 200) loss: 1.666573
(Iteration 191 / 200) loss: 1.909494
(Epoch 5 / 5) train acc: 0.372000; val_acc: 0.319000
```

Running with sgd_momentum

```
[3072, 100, 100, 100, 100, 100, 10]
(Iteration 1 / 200) loss: 3.153778
(Epoch 0 / 5) train acc: 0.099000; val_acc: 0.088000
(Iteration 11 / 200) loss: 2.227203
(Iteration 21 / 200) loss: 2.125706
(Iteration 31 / 200) loss: 1.932695
(Epoch 1 / 5) train acc: 0.307000; val_acc: 0.260000
(Iteration 41 / 200) loss: 1.946488
(Iteration 51 / 200) loss: 1.778583
(Iteration 61 / 200) loss: 1.758119
(Iteration 71 / 200) loss: 1.849137
(Epoch 2 / 5) train acc: 0.382000; val_acc: 0.322000
(Iteration 81 / 200) loss: 2.048671
(Iteration 91 / 200) loss: 1.693223
(Iteration 101 / 200) loss: 1.511693
(Iteration 111 / 200) loss: 1.390754
```

(Epoch 3 / 5) train acc: 0.458000; val_acc: 0.338000
 (Iteration 121 / 200) loss: 1.670614
 (Iteration 131 / 200) loss: 1.540271
 (Iteration 141 / 200) loss: 1.597365
 (Iteration 151 / 200) loss: 1.609851
 (Epoch 4 / 5) train acc: 0.490000; val_acc: 0.327000
 (Iteration 161 / 200) loss: 1.472687
 (Iteration 171 / 200) loss: 1.378620
 (Iteration 181 / 200) loss: 1.378174
 (Iteration 191 / 200) loss: 1.305934
 (Epoch 5 / 5) train acc: 0.536000; val_acc: 0.368000



2.2 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.

```
[ ]: # 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

```
[ ]: # 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:

```

[ ]: learning_rates = {'rmsprop': 1e-4, 'adam': 1e-3}
for update_rule in ['adam', 'rmsprop']:
    print('Running with ', update_rule)
    model = FullyConnectedNet(
        [100, 100, 100, 100, 100],
        weight_scale=5e-2
    )
    solver = Solver(
        model,
        small_data,
        num_epochs=5,
        batch_size=100,

```

```

        update_rule=update_rule,
        optim_config={'learning_rate': learning_rates[update_rule]},
        verbose=True
    )
    solvers[update_rule] = solver
    solver.train()
    print()

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

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

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

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

plt.show()

```

Running with adam

```

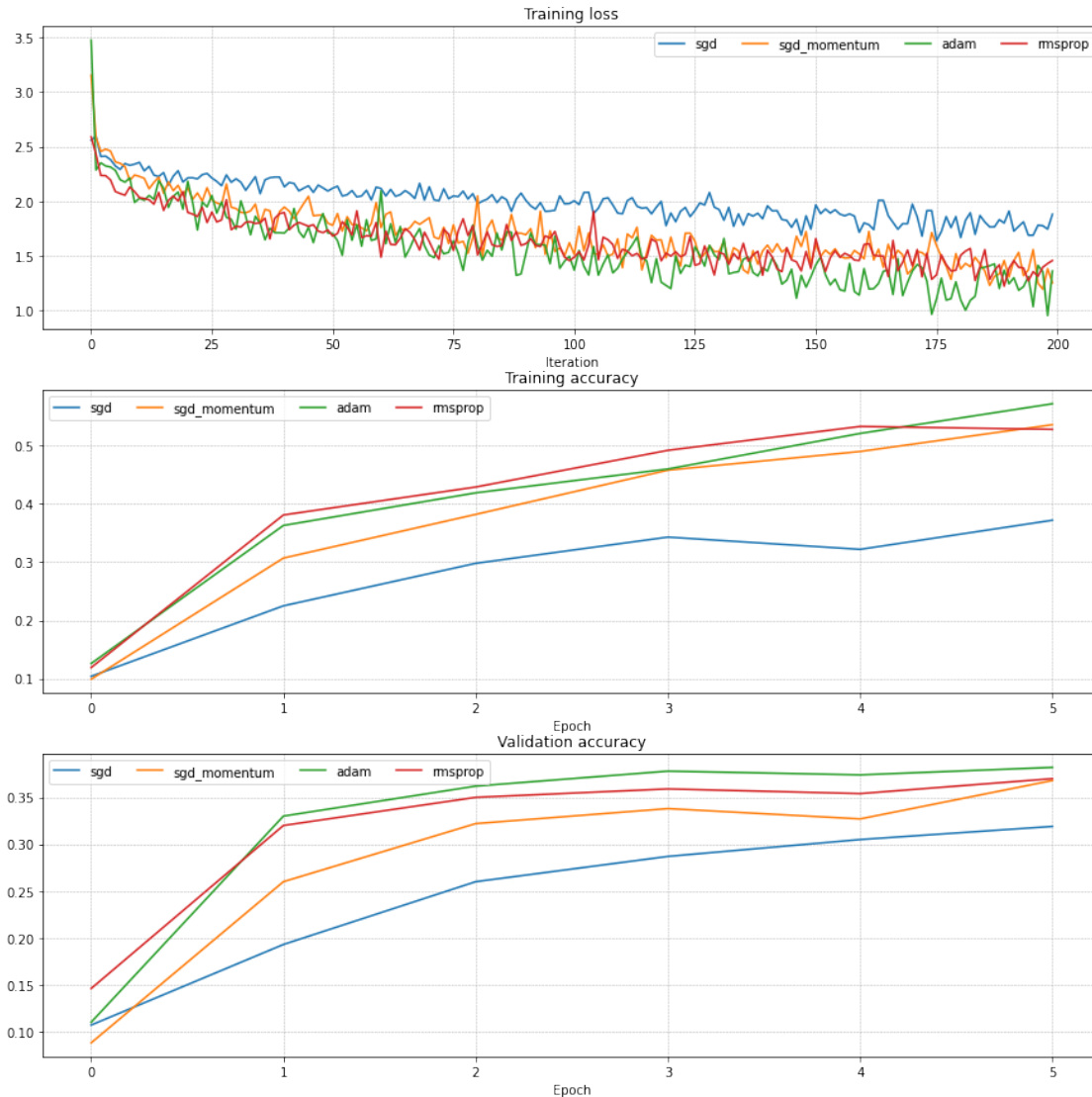
[3072, 100, 100, 100, 100, 100, 10]
(Iteration 1 / 200) loss: 3.476928
(Epoch 0 / 5) train acc: 0.126000; val_acc: 0.110000
(Iteration 11 / 200) loss: 2.027712
(Iteration 21 / 200) loss: 2.183358
(Iteration 31 / 200) loss: 1.744257
(Epoch 1 / 5) train acc: 0.363000; val_acc: 0.330000
(Iteration 41 / 200) loss: 1.707951
(Iteration 51 / 200) loss: 1.703835
(Iteration 61 / 200) loss: 2.094758
(Iteration 71 / 200) loss: 1.505557
(Epoch 2 / 5) train acc: 0.419000; val_acc: 0.362000
(Iteration 81 / 200) loss: 1.594430
(Iteration 91 / 200) loss: 1.519017
(Iteration 101 / 200) loss: 1.368522
(Iteration 111 / 200) loss: 1.470400
(Epoch 3 / 5) train acc: 0.460000; val_acc: 0.378000
(Iteration 121 / 200) loss: 1.199064

```

(Iteration 131 / 200) loss: 1.464705
(Iteration 141 / 200) loss: 1.359863
(Iteration 151 / 200) loss: 1.415068
(Epoch 4 / 5) train acc: 0.521000; val_acc: 0.374000
(Iteration 161 / 200) loss: 1.382818
(Iteration 171 / 200) loss: 1.359900
(Iteration 181 / 200) loss: 1.095947
(Iteration 191 / 200) loss: 1.243088
(Epoch 5 / 5) train acc: 0.572000; val_acc: 0.382000

Running with rmsprop

[3072, 100, 100, 100, 100, 100, 10]
(Iteration 1 / 200) loss: 2.589166
(Epoch 0 / 5) train acc: 0.119000; val_acc: 0.146000
(Iteration 11 / 200) loss: 2.032921
(Iteration 21 / 200) loss: 1.897278
(Iteration 31 / 200) loss: 1.770793
(Epoch 1 / 5) train acc: 0.381000; val_acc: 0.320000
(Iteration 41 / 200) loss: 1.895731
(Iteration 51 / 200) loss: 1.681091
(Iteration 61 / 200) loss: 1.487204
(Iteration 71 / 200) loss: 1.629973
(Epoch 2 / 5) train acc: 0.429000; val_acc: 0.350000
(Iteration 81 / 200) loss: 1.506686
(Iteration 91 / 200) loss: 1.610742
(Iteration 101 / 200) loss: 1.486124
(Iteration 111 / 200) loss: 1.559454
(Epoch 3 / 5) train acc: 0.492000; val_acc: 0.359000
(Iteration 121 / 200) loss: 1.496860
(Iteration 131 / 200) loss: 1.531552
(Iteration 141 / 200) loss: 1.550195
(Iteration 151 / 200) loss: 1.657838
(Epoch 4 / 5) train acc: 0.533000; val_acc: 0.354000
(Iteration 161 / 200) loss: 1.603105
(Iteration 171 / 200) loss: 1.408064
(Iteration 181 / 200) loss: 1.504707
(Iteration 191 / 200) loss: 1.385212
(Epoch 5 / 5) train acc: 0.528000; val_acc: 0.370000



2.3 Inline Question 2:

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

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

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

2.4 Answer:

No, AdaGrad decrease the update if the direction is already update several times. In some case, the gradient should go to the lowest point which is the same direction, then the AdaGrad will not work well, however, the Adam have the speed (momentum) so it will still doing fine in this case.

3 Train a Good Model!

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

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

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

```
[ ]: 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)*****
# I increase the dataset a little bit for the better result
num_train = 20000
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"],
}
# finding good parameters
learning_rates = {'rmsprop': 1e-4, 'adam': 1e-3}
hss = [50, 75, 100]
lrs = [5e-3]
lr_decs = [1]
# define necessary parameters
input_size = 32 * 32 * 3
num_classes = 10
best = 0

weight_scale = 1e-1
```

```

# Find the best parameters by go through every pairs of parameters
for hs in hss:
    model = FullyConnectedNet(
        [hs],
        dtype=np.float64,
        reg=1
    )
    solver = Solver(
        model,
        data,
        num_epochs=10,
        batch_size=100,
        update_rule='sgd',
        optim_config={'learning_rate': 5e-4},
        verbose=True
    )

    # train model
    solver.train()
    # find best model
    if solver.best_val_acc > best:
        best = solver.best_val_acc
        best_model = model
print(best)
# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
#####
#                                     END OF YOUR CODE                                     #
#####

```

```

[3072, 50, 10]
(Iteration 1 / 4900) loss: 11.282419
(Epoch 0 / 10) train acc: 0.140000; val_acc: 0.135000
(Iteration 11 / 4900) loss: 9.968199
(Iteration 21 / 4900) loss: 9.633723
(Iteration 31 / 4900) loss: 9.476028
(Iteration 41 / 4900) loss: 9.364198
(Iteration 51 / 4900) loss: 9.266560
(Iteration 61 / 4900) loss: 9.233705
(Iteration 71 / 4900) loss: 9.433305
(Iteration 81 / 4900) loss: 8.930638
(Iteration 91 / 4900) loss: 8.958189
(Iteration 101 / 4900) loss: 8.762259
(Iteration 111 / 4900) loss: 8.876290
(Iteration 121 / 4900) loss: 8.804269
(Iteration 131 / 4900) loss: 8.602855
(Iteration 141 / 4900) loss: 8.583471

```

(Iteration 151 / 4900) loss: 8.538131
(Iteration 161 / 4900) loss: 8.390272
(Iteration 171 / 4900) loss: 8.250602
(Iteration 181 / 4900) loss: 8.086172
(Iteration 191 / 4900) loss: 8.135666
(Iteration 201 / 4900) loss: 8.235827
(Iteration 211 / 4900) loss: 8.006330
(Iteration 221 / 4900) loss: 7.865878
(Iteration 231 / 4900) loss: 7.807127
(Iteration 241 / 4900) loss: 7.890489
(Iteration 251 / 4900) loss: 7.673140
(Iteration 261 / 4900) loss: 7.808530
(Iteration 271 / 4900) loss: 7.598275
(Iteration 281 / 4900) loss: 7.629701
(Iteration 291 / 4900) loss: 7.596468
(Iteration 301 / 4900) loss: 7.276009
(Iteration 311 / 4900) loss: 7.382398
(Iteration 321 / 4900) loss: 7.353002
(Iteration 331 / 4900) loss: 7.359075
(Iteration 341 / 4900) loss: 7.157529
(Iteration 351 / 4900) loss: 7.081836
(Iteration 361 / 4900) loss: 7.103408
(Iteration 371 / 4900) loss: 7.010299
(Iteration 381 / 4900) loss: 6.912655
(Iteration 391 / 4900) loss: 6.825524
(Iteration 401 / 4900) loss: 6.921574
(Iteration 411 / 4900) loss: 6.720825
(Iteration 421 / 4900) loss: 6.864876
(Iteration 431 / 4900) loss: 6.816096
(Iteration 441 / 4900) loss: 6.575716
(Iteration 451 / 4900) loss: 6.532333
(Iteration 461 / 4900) loss: 6.531259
(Iteration 471 / 4900) loss: 6.451174
(Iteration 481 / 4900) loss: 6.437514
(Epoch 1 / 10) train acc: 0.413000; val_acc: 0.402000
(Iteration 491 / 4900) loss: 6.396889
(Iteration 501 / 4900) loss: 6.289858
(Iteration 511 / 4900) loss: 6.109087
(Iteration 521 / 4900) loss: 6.187563
(Iteration 531 / 4900) loss: 6.150017
(Iteration 541 / 4900) loss: 6.503099
(Iteration 551 / 4900) loss: 6.064406
(Iteration 561 / 4900) loss: 6.174423
(Iteration 571 / 4900) loss: 6.113128
(Iteration 581 / 4900) loss: 6.014652
(Iteration 591 / 4900) loss: 5.882647
(Iteration 601 / 4900) loss: 5.960659
(Iteration 611 / 4900) loss: 5.781807

(Iteration 621 / 4900) loss: 5.730725
(Iteration 631 / 4900) loss: 5.678347
(Iteration 641 / 4900) loss: 5.732500
(Iteration 651 / 4900) loss: 5.714536
(Iteration 661 / 4900) loss: 5.481376
(Iteration 671 / 4900) loss: 5.638440
(Iteration 681 / 4900) loss: 5.632932
(Iteration 691 / 4900) loss: 5.317564
(Iteration 701 / 4900) loss: 5.537261
(Iteration 711 / 4900) loss: 5.310406
(Iteration 721 / 4900) loss: 5.518395
(Iteration 731 / 4900) loss: 5.210270
(Iteration 741 / 4900) loss: 5.170378
(Iteration 751 / 4900) loss: 5.289611
(Iteration 761 / 4900) loss: 5.268250
(Iteration 771 / 4900) loss: 5.274011
(Iteration 781 / 4900) loss: 5.144196
(Iteration 791 / 4900) loss: 5.114886
(Iteration 801 / 4900) loss: 5.042549
(Iteration 811 / 4900) loss: 5.021859
(Iteration 821 / 4900) loss: 4.986270
(Iteration 831 / 4900) loss: 5.060661
(Iteration 841 / 4900) loss: 4.919314
(Iteration 851 / 4900) loss: 4.880406
(Iteration 861 / 4900) loss: 4.825929
(Iteration 871 / 4900) loss: 4.845911
(Iteration 881 / 4900) loss: 4.775297
(Iteration 891 / 4900) loss: 4.778613
(Iteration 901 / 4900) loss: 4.791790
(Iteration 911 / 4900) loss: 4.741757
(Iteration 921 / 4900) loss: 4.830344
(Iteration 931 / 4900) loss: 4.672534
(Iteration 941 / 4900) loss: 4.754941
(Iteration 951 / 4900) loss: 4.528821
(Iteration 961 / 4900) loss: 4.603891
(Iteration 971 / 4900) loss: 4.598834
(Epoch 2 / 10) train acc: 0.420000; val_acc: 0.440000
(Iteration 981 / 4900) loss: 4.534292
(Iteration 991 / 4900) loss: 4.426043
(Iteration 1001 / 4900) loss: 4.578464
(Iteration 1011 / 4900) loss: 4.549718
(Iteration 1021 / 4900) loss: 4.290295
(Iteration 1031 / 4900) loss: 4.288562
(Iteration 1041 / 4900) loss: 4.428186
(Iteration 1051 / 4900) loss: 4.288876
(Iteration 1061 / 4900) loss: 4.417127
(Iteration 1071 / 4900) loss: 4.334971
(Iteration 1081 / 4900) loss: 4.291742

(Iteration 1091 / 4900) loss: 4.209525
(Iteration 1101 / 4900) loss: 4.203670
(Iteration 1111 / 4900) loss: 4.218908
(Iteration 1121 / 4900) loss: 4.215084
(Iteration 1131 / 4900) loss: 3.961191
(Iteration 1141 / 4900) loss: 3.956223
(Iteration 1151 / 4900) loss: 4.262330
(Iteration 1161 / 4900) loss: 4.018383
(Iteration 1171 / 4900) loss: 3.970557
(Iteration 1181 / 4900) loss: 4.081530
(Iteration 1191 / 4900) loss: 3.878496
(Iteration 1201 / 4900) loss: 3.931389
(Iteration 1211 / 4900) loss: 3.863087
(Iteration 1221 / 4900) loss: 3.867409
(Iteration 1231 / 4900) loss: 3.771217
(Iteration 1241 / 4900) loss: 3.767381
(Iteration 1251 / 4900) loss: 3.758729
(Iteration 1261 / 4900) loss: 3.877628
(Iteration 1271 / 4900) loss: 3.753252
(Iteration 1281 / 4900) loss: 3.735858
(Iteration 1291 / 4900) loss: 3.782446
(Iteration 1301 / 4900) loss: 3.948104
(Iteration 1311 / 4900) loss: 3.560443
(Iteration 1321 / 4900) loss: 3.620041
(Iteration 1331 / 4900) loss: 3.650875
(Iteration 1341 / 4900) loss: 3.711539
(Iteration 1351 / 4900) loss: 3.547567
(Iteration 1361 / 4900) loss: 3.585062
(Iteration 1371 / 4900) loss: 3.550392
(Iteration 1381 / 4900) loss: 3.516006
(Iteration 1391 / 4900) loss: 3.459122
(Iteration 1401 / 4900) loss: 3.523197
(Iteration 1411 / 4900) loss: 3.540646
(Iteration 1421 / 4900) loss: 3.538179
(Iteration 1431 / 4900) loss: 3.389247
(Iteration 1441 / 4900) loss: 3.360791
(Iteration 1451 / 4900) loss: 3.501390
(Iteration 1461 / 4900) loss: 3.238175
(Epoch 3 / 10) train acc: 0.445000; val_acc: 0.450000
(Iteration 1471 / 4900) loss: 3.492374
(Iteration 1481 / 4900) loss: 3.410812
(Iteration 1491 / 4900) loss: 3.280988
(Iteration 1501 / 4900) loss: 3.314912
(Iteration 1511 / 4900) loss: 3.276452
(Iteration 1521 / 4900) loss: 3.189819
(Iteration 1531 / 4900) loss: 3.236500
(Iteration 1541 / 4900) loss: 3.150051
(Iteration 1551 / 4900) loss: 3.142013

(Iteration 1561 / 4900) loss: 3.168597
(Iteration 1571 / 4900) loss: 3.172030
(Iteration 1581 / 4900) loss: 3.384908
(Iteration 1591 / 4900) loss: 3.068518
(Iteration 1601 / 4900) loss: 3.248493
(Iteration 1611 / 4900) loss: 3.075060
(Iteration 1621 / 4900) loss: 3.305375
(Iteration 1631 / 4900) loss: 3.203135
(Iteration 1641 / 4900) loss: 2.917352
(Iteration 1651 / 4900) loss: 3.117749
(Iteration 1661 / 4900) loss: 3.046723
(Iteration 1671 / 4900) loss: 2.978084
(Iteration 1681 / 4900) loss: 2.940561
(Iteration 1691 / 4900) loss: 3.027775
(Iteration 1701 / 4900) loss: 2.981226
(Iteration 1711 / 4900) loss: 3.071202
(Iteration 1721 / 4900) loss: 3.055964
(Iteration 1731 / 4900) loss: 2.998065
(Iteration 1741 / 4900) loss: 2.795276
(Iteration 1751 / 4900) loss: 2.925977
(Iteration 1761 / 4900) loss: 2.994068
(Iteration 1771 / 4900) loss: 2.724434
(Iteration 1781 / 4900) loss: 2.896137
(Iteration 1791 / 4900) loss: 2.687204
(Iteration 1801 / 4900) loss: 2.892555
(Iteration 1811 / 4900) loss: 2.608178
(Iteration 1821 / 4900) loss: 2.707347
(Iteration 1831 / 4900) loss: 2.838930
(Iteration 1841 / 4900) loss: 2.802236
(Iteration 1851 / 4900) loss: 2.789340
(Iteration 1861 / 4900) loss: 2.775183
(Iteration 1871 / 4900) loss: 2.856689
(Iteration 1881 / 4900) loss: 2.676849
(Iteration 1891 / 4900) loss: 2.885970
(Iteration 1901 / 4900) loss: 2.744582
(Iteration 1911 / 4900) loss: 2.688391
(Iteration 1921 / 4900) loss: 2.648777
(Iteration 1931 / 4900) loss: 2.551921
(Iteration 1941 / 4900) loss: 2.636963
(Iteration 1951 / 4900) loss: 2.508312
(Epoch 4 / 10) train acc: 0.477000; val_acc: 0.479000
(Iteration 1961 / 4900) loss: 2.324823
(Iteration 1971 / 4900) loss: 2.563849
(Iteration 1981 / 4900) loss: 2.624141
(Iteration 1991 / 4900) loss: 2.492305
(Iteration 2001 / 4900) loss: 2.622084
(Iteration 2011 / 4900) loss: 2.600721
(Iteration 2021 / 4900) loss: 2.579433

(Iteration 2031 / 4900) loss: 2.590766
(Iteration 2041 / 4900) loss: 2.629797
(Iteration 2051 / 4900) loss: 2.487854
(Iteration 2061 / 4900) loss: 2.623735
(Iteration 2071 / 4900) loss: 2.567229
(Iteration 2081 / 4900) loss: 2.716104
(Iteration 2091 / 4900) loss: 2.484239
(Iteration 2101 / 4900) loss: 2.511346
(Iteration 2111 / 4900) loss: 2.559405
(Iteration 2121 / 4900) loss: 2.457385
(Iteration 2131 / 4900) loss: 2.428178
(Iteration 2141 / 4900) loss: 2.374234
(Iteration 2151 / 4900) loss: 2.462039
(Iteration 2161 / 4900) loss: 2.337800
(Iteration 2171 / 4900) loss: 2.344820
(Iteration 2181 / 4900) loss: 2.446388
(Iteration 2191 / 4900) loss: 2.403137
(Iteration 2201 / 4900) loss: 2.407101
(Iteration 2211 / 4900) loss: 2.401403
(Iteration 2221 / 4900) loss: 2.247599
(Iteration 2231 / 4900) loss: 2.527409
(Iteration 2241 / 4900) loss: 2.326273
(Iteration 2251 / 4900) loss: 2.192865
(Iteration 2261 / 4900) loss: 2.240725
(Iteration 2271 / 4900) loss: 2.345542
(Iteration 2281 / 4900) loss: 2.233930
(Iteration 2291 / 4900) loss: 2.312593
(Iteration 2301 / 4900) loss: 2.403011
(Iteration 2311 / 4900) loss: 2.345712
(Iteration 2321 / 4900) loss: 2.222351
(Iteration 2331 / 4900) loss: 2.371153
(Iteration 2341 / 4900) loss: 2.383452
(Iteration 2351 / 4900) loss: 2.253364
(Iteration 2361 / 4900) loss: 2.152177
(Iteration 2371 / 4900) loss: 2.301509
(Iteration 2381 / 4900) loss: 2.268225
(Iteration 2391 / 4900) loss: 2.268050
(Iteration 2401 / 4900) loss: 2.149960
(Iteration 2411 / 4900) loss: 2.081528
(Iteration 2421 / 4900) loss: 2.230056
(Iteration 2431 / 4900) loss: 2.207303
(Iteration 2441 / 4900) loss: 2.299841
(Epoch 5 / 10) train acc: 0.491000; val_acc: 0.472000
(Iteration 2451 / 4900) loss: 2.168053
(Iteration 2461 / 4900) loss: 2.038197
(Iteration 2471 / 4900) loss: 2.228795
(Iteration 2481 / 4900) loss: 2.114399
(Iteration 2491 / 4900) loss: 1.960511

(Iteration 2501 / 4900) loss: 2.111447
(Iteration 2511 / 4900) loss: 2.260801
(Iteration 2521 / 4900) loss: 2.047407
(Iteration 2531 / 4900) loss: 2.148435
(Iteration 2541 / 4900) loss: 2.150414
(Iteration 2551 / 4900) loss: 2.227616
(Iteration 2561 / 4900) loss: 2.216553
(Iteration 2571 / 4900) loss: 2.172396
(Iteration 2581 / 4900) loss: 2.036021
(Iteration 2591 / 4900) loss: 2.075383
(Iteration 2601 / 4900) loss: 2.074709
(Iteration 2611 / 4900) loss: 1.943037
(Iteration 2621 / 4900) loss: 1.985107
(Iteration 2631 / 4900) loss: 2.078042
(Iteration 2641 / 4900) loss: 2.067537
(Iteration 2651 / 4900) loss: 2.072585
(Iteration 2661 / 4900) loss: 1.986338
(Iteration 2671 / 4900) loss: 2.109556
(Iteration 2681 / 4900) loss: 2.139933
(Iteration 2691 / 4900) loss: 1.971628
(Iteration 2701 / 4900) loss: 2.052781
(Iteration 2711 / 4900) loss: 2.091902
(Iteration 2721 / 4900) loss: 2.119129
(Iteration 2731 / 4900) loss: 2.050882
(Iteration 2741 / 4900) loss: 2.228422
(Iteration 2751 / 4900) loss: 2.029029
(Iteration 2761 / 4900) loss: 2.059679
(Iteration 2771 / 4900) loss: 2.201745
(Iteration 2781 / 4900) loss: 1.988102
(Iteration 2791 / 4900) loss: 1.858021
(Iteration 2801 / 4900) loss: 2.159015
(Iteration 2811 / 4900) loss: 2.066959
(Iteration 2821 / 4900) loss: 2.148929
(Iteration 2831 / 4900) loss: 1.999595
(Iteration 2841 / 4900) loss: 2.056111
(Iteration 2851 / 4900) loss: 2.113946
(Iteration 2861 / 4900) loss: 1.921687
(Iteration 2871 / 4900) loss: 1.938316
(Iteration 2881 / 4900) loss: 1.878329
(Iteration 2891 / 4900) loss: 1.875925
(Iteration 2901 / 4900) loss: 1.912492
(Iteration 2911 / 4900) loss: 2.010253
(Iteration 2921 / 4900) loss: 1.937808
(Iteration 2931 / 4900) loss: 1.783289
(Epoch 6 / 10) train acc: 0.520000; val_acc: 0.472000
(Iteration 2941 / 4900) loss: 2.066633
(Iteration 2951 / 4900) loss: 2.024844
(Iteration 2961 / 4900) loss: 2.094778

(Iteration 2971 / 4900) loss: 1.924314
(Iteration 2981 / 4900) loss: 1.918301
(Iteration 2991 / 4900) loss: 1.958458
(Iteration 3001 / 4900) loss: 2.009512
(Iteration 3011 / 4900) loss: 1.904492
(Iteration 3021 / 4900) loss: 1.889477
(Iteration 3031 / 4900) loss: 1.946074
(Iteration 3041 / 4900) loss: 1.940513
(Iteration 3051 / 4900) loss: 1.805045
(Iteration 3061 / 4900) loss: 1.713323
(Iteration 3071 / 4900) loss: 1.957798
(Iteration 3081 / 4900) loss: 1.986414
(Iteration 3091 / 4900) loss: 1.817472
(Iteration 3101 / 4900) loss: 2.006638
(Iteration 3111 / 4900) loss: 1.661014
(Iteration 3121 / 4900) loss: 1.986469
(Iteration 3131 / 4900) loss: 1.815904
(Iteration 3141 / 4900) loss: 1.794466
(Iteration 3151 / 4900) loss: 1.869385
(Iteration 3161 / 4900) loss: 1.708356
(Iteration 3171 / 4900) loss: 1.738672
(Iteration 3181 / 4900) loss: 1.690910
(Iteration 3191 / 4900) loss: 1.822841
(Iteration 3201 / 4900) loss: 1.805599
(Iteration 3211 / 4900) loss: 1.893313
(Iteration 3221 / 4900) loss: 1.782435
(Iteration 3231 / 4900) loss: 1.732201
(Iteration 3241 / 4900) loss: 1.753592
(Iteration 3251 / 4900) loss: 1.777553
(Iteration 3261 / 4900) loss: 1.868824
(Iteration 3271 / 4900) loss: 1.995637
(Iteration 3281 / 4900) loss: 1.856055
(Iteration 3291 / 4900) loss: 1.934562
(Iteration 3301 / 4900) loss: 1.853974
(Iteration 3311 / 4900) loss: 1.816999
(Iteration 3321 / 4900) loss: 1.889596
(Iteration 3331 / 4900) loss: 1.888565
(Iteration 3341 / 4900) loss: 1.694916
(Iteration 3351 / 4900) loss: 1.926625
(Iteration 3361 / 4900) loss: 1.882928
(Iteration 3371 / 4900) loss: 1.785786
(Iteration 3381 / 4900) loss: 1.780455
(Iteration 3391 / 4900) loss: 1.751042
(Iteration 3401 / 4900) loss: 1.669621
(Iteration 3411 / 4900) loss: 1.818193
(Iteration 3421 / 4900) loss: 1.873247
(Epoch 7 / 10) train acc: 0.489000; val_acc: 0.471000
(Iteration 3431 / 4900) loss: 1.920669

(Iteration 3441 / 4900) loss: 1.662033
(Iteration 3451 / 4900) loss: 1.968277
(Iteration 3461 / 4900) loss: 1.887043
(Iteration 3471 / 4900) loss: 1.644716
(Iteration 3481 / 4900) loss: 1.662648
(Iteration 3491 / 4900) loss: 1.782571
(Iteration 3501 / 4900) loss: 1.938963
(Iteration 3511 / 4900) loss: 1.854960
(Iteration 3521 / 4900) loss: 1.696779
(Iteration 3531 / 4900) loss: 1.599203
(Iteration 3541 / 4900) loss: 1.832434
(Iteration 3551 / 4900) loss: 1.575063
(Iteration 3561 / 4900) loss: 1.721467
(Iteration 3571 / 4900) loss: 1.857618
(Iteration 3581 / 4900) loss: 1.591602
(Iteration 3591 / 4900) loss: 1.870461
(Iteration 3601 / 4900) loss: 1.862983
(Iteration 3611 / 4900) loss: 1.803052
(Iteration 3621 / 4900) loss: 1.589775
(Iteration 3631 / 4900) loss: 1.867459
(Iteration 3641 / 4900) loss: 1.808135
(Iteration 3651 / 4900) loss: 1.694002
(Iteration 3661 / 4900) loss: 1.762519
(Iteration 3671 / 4900) loss: 1.897809
(Iteration 3681 / 4900) loss: 1.743049
(Iteration 3691 / 4900) loss: 1.718383
(Iteration 3701 / 4900) loss: 1.702980
(Iteration 3711 / 4900) loss: 1.749403
(Iteration 3721 / 4900) loss: 1.566922
(Iteration 3731 / 4900) loss: 1.689535
(Iteration 3741 / 4900) loss: 1.743281
(Iteration 3751 / 4900) loss: 1.584100
(Iteration 3761 / 4900) loss: 1.629881
(Iteration 3771 / 4900) loss: 1.819011
(Iteration 3781 / 4900) loss: 1.788917
(Iteration 3791 / 4900) loss: 1.703337
(Iteration 3801 / 4900) loss: 1.858178
(Iteration 3811 / 4900) loss: 1.780955
(Iteration 3821 / 4900) loss: 1.790777
(Iteration 3831 / 4900) loss: 1.579376
(Iteration 3841 / 4900) loss: 1.663560
(Iteration 3851 / 4900) loss: 1.699659
(Iteration 3861 / 4900) loss: 1.717622
(Iteration 3871 / 4900) loss: 1.677709
(Iteration 3881 / 4900) loss: 1.518268
(Iteration 3891 / 4900) loss: 1.758975
(Iteration 3901 / 4900) loss: 1.650993
(Iteration 3911 / 4900) loss: 1.495586

(Epoch 8 / 10) train acc: 0.535000; val_acc: 0.474000
(Iteration 3921 / 4900) loss: 1.739283
(Iteration 3931 / 4900) loss: 1.734251
(Iteration 3941 / 4900) loss: 1.682810
(Iteration 3951 / 4900) loss: 1.578223
(Iteration 3961 / 4900) loss: 1.865077
(Iteration 3971 / 4900) loss: 1.734190
(Iteration 3981 / 4900) loss: 1.633722
(Iteration 3991 / 4900) loss: 1.617658
(Iteration 4001 / 4900) loss: 1.793436
(Iteration 4011 / 4900) loss: 1.697813
(Iteration 4021 / 4900) loss: 1.631996
(Iteration 4031 / 4900) loss: 1.742932
(Iteration 4041 / 4900) loss: 1.767911
(Iteration 4051 / 4900) loss: 1.662956
(Iteration 4061 / 4900) loss: 1.798554
(Iteration 4071 / 4900) loss: 1.620395
(Iteration 4081 / 4900) loss: 1.632987
(Iteration 4091 / 4900) loss: 1.728377
(Iteration 4101 / 4900) loss: 1.592749
(Iteration 4111 / 4900) loss: 1.587567
(Iteration 4121 / 4900) loss: 1.695208
(Iteration 4131 / 4900) loss: 1.685817
(Iteration 4141 / 4900) loss: 1.731919
(Iteration 4151 / 4900) loss: 1.568822
(Iteration 4161 / 4900) loss: 1.827337
(Iteration 4171 / 4900) loss: 1.824277
(Iteration 4181 / 4900) loss: 1.740495
(Iteration 4191 / 4900) loss: 1.652143
(Iteration 4201 / 4900) loss: 1.590637
(Iteration 4211 / 4900) loss: 1.850203
(Iteration 4221 / 4900) loss: 1.569353
(Iteration 4231 / 4900) loss: 1.610918
(Iteration 4241 / 4900) loss: 1.669749
(Iteration 4251 / 4900) loss: 1.960138
(Iteration 4261 / 4900) loss: 1.593132
(Iteration 4271 / 4900) loss: 1.495227
(Iteration 4281 / 4900) loss: 1.506143
(Iteration 4291 / 4900) loss: 1.584029
(Iteration 4301 / 4900) loss: 1.640150
(Iteration 4311 / 4900) loss: 1.425570
(Iteration 4321 / 4900) loss: 1.733288
(Iteration 4331 / 4900) loss: 1.627631
(Iteration 4341 / 4900) loss: 1.696685
(Iteration 4351 / 4900) loss: 1.627196
(Iteration 4361 / 4900) loss: 1.509448
(Iteration 4371 / 4900) loss: 1.639318
(Iteration 4381 / 4900) loss: 1.695450

(Iteration 4391 / 4900) loss: 1.611222
(Iteration 4401 / 4900) loss: 1.586930
(Epoch 9 / 10) train acc: 0.509000; val_acc: 0.473000
(Iteration 4411 / 4900) loss: 1.498096
(Iteration 4421 / 4900) loss: 1.693841
(Iteration 4431 / 4900) loss: 1.602059
(Iteration 4441 / 4900) loss: 1.654045
(Iteration 4451 / 4900) loss: 1.622450
(Iteration 4461 / 4900) loss: 1.709002
(Iteration 4471 / 4900) loss: 1.573683
(Iteration 4481 / 4900) loss: 1.714457
(Iteration 4491 / 4900) loss: 1.574970
(Iteration 4501 / 4900) loss: 1.609705
(Iteration 4511 / 4900) loss: 1.719282
(Iteration 4521 / 4900) loss: 1.725841
(Iteration 4531 / 4900) loss: 1.755162
(Iteration 4541 / 4900) loss: 1.533384
(Iteration 4551 / 4900) loss: 1.586245
(Iteration 4561 / 4900) loss: 1.505944
(Iteration 4571 / 4900) loss: 1.666541
(Iteration 4581 / 4900) loss: 1.652965
(Iteration 4591 / 4900) loss: 1.504666
(Iteration 4601 / 4900) loss: 1.520801
(Iteration 4611 / 4900) loss: 1.762950
(Iteration 4621 / 4900) loss: 1.333247
(Iteration 4631 / 4900) loss: 1.570440
(Iteration 4641 / 4900) loss: 1.616155
(Iteration 4651 / 4900) loss: 1.716586
(Iteration 4661 / 4900) loss: 1.605718
(Iteration 4671 / 4900) loss: 1.709310
(Iteration 4681 / 4900) loss: 1.569977
(Iteration 4691 / 4900) loss: 1.596912
(Iteration 4701 / 4900) loss: 1.539862
(Iteration 4711 / 4900) loss: 1.471177
(Iteration 4721 / 4900) loss: 1.631972
(Iteration 4731 / 4900) loss: 1.534986
(Iteration 4741 / 4900) loss: 1.582402
(Iteration 4751 / 4900) loss: 1.567581
(Iteration 4761 / 4900) loss: 1.599855
(Iteration 4771 / 4900) loss: 1.546536
(Iteration 4781 / 4900) loss: 1.621544
(Iteration 4791 / 4900) loss: 1.432462
(Iteration 4801 / 4900) loss: 1.774172
(Iteration 4811 / 4900) loss: 1.594071
(Iteration 4821 / 4900) loss: 1.711673
(Iteration 4831 / 4900) loss: 1.662067
(Iteration 4841 / 4900) loss: 1.570997
(Iteration 4851 / 4900) loss: 1.623583

(Iteration 4861 / 4900) loss: 1.581816
(Iteration 4871 / 4900) loss: 1.599145
(Iteration 4881 / 4900) loss: 1.606949
(Iteration 4891 / 4900) loss: 1.634506
(Epoch 10 / 10) train acc: 0.513000; val_acc: 0.485000
[3072, 75, 10]
(Iteration 1 / 4900) loss: 15.411905
(Epoch 0 / 10) train acc: 0.116000; val_acc: 0.118000
(Iteration 11 / 4900) loss: 13.823285
(Iteration 21 / 4900) loss: 13.542849
(Iteration 31 / 4900) loss: 13.447668
(Iteration 41 / 4900) loss: 13.172847
(Iteration 51 / 4900) loss: 12.850146
(Iteration 61 / 4900) loss: 12.932720
(Iteration 71 / 4900) loss: 12.745747
(Iteration 81 / 4900) loss: 12.506118
(Iteration 91 / 4900) loss: 12.483908
(Iteration 101 / 4900) loss: 12.241004
(Iteration 111 / 4900) loss: 12.072383
(Iteration 121 / 4900) loss: 12.074197
(Iteration 131 / 4900) loss: 12.005641
(Iteration 141 / 4900) loss: 11.662590
(Iteration 151 / 4900) loss: 11.729775
(Iteration 161 / 4900) loss: 11.444464
(Iteration 171 / 4900) loss: 11.617797
(Iteration 181 / 4900) loss: 11.433036
(Iteration 191 / 4900) loss: 11.260824
(Iteration 201 / 4900) loss: 11.246753
(Iteration 211 / 4900) loss: 11.133510
(Iteration 221 / 4900) loss: 10.950841
(Iteration 231 / 4900) loss: 10.866785
(Iteration 241 / 4900) loss: 10.784354
(Iteration 251 / 4900) loss: 10.789494
(Iteration 261 / 4900) loss: 10.728761
(Iteration 271 / 4900) loss: 10.363377
(Iteration 281 / 4900) loss: 10.426041
(Iteration 291 / 4900) loss: 10.297578
(Iteration 301 / 4900) loss: 10.116158
(Iteration 311 / 4900) loss: 10.274701
(Iteration 321 / 4900) loss: 9.879666
(Iteration 331 / 4900) loss: 10.000430
(Iteration 341 / 4900) loss: 9.922706
(Iteration 351 / 4900) loss: 9.770449
(Iteration 361 / 4900) loss: 9.941370
(Iteration 371 / 4900) loss: 9.688670
(Iteration 381 / 4900) loss: 9.394746
(Iteration 391 / 4900) loss: 9.386707
(Iteration 401 / 4900) loss: 9.528063

(Iteration 411 / 4900) loss: 9.311938
(Iteration 421 / 4900) loss: 9.416844
(Iteration 431 / 4900) loss: 9.198236
(Iteration 441 / 4900) loss: 9.136867
(Iteration 451 / 4900) loss: 8.958826
(Iteration 461 / 4900) loss: 9.012747
(Iteration 471 / 4900) loss: 8.797889
(Iteration 481 / 4900) loss: 8.769560
(Epoch 1 / 10) train acc: 0.421000; val_acc: 0.411000
(Iteration 491 / 4900) loss: 8.542584
(Iteration 501 / 4900) loss: 8.467985
(Iteration 511 / 4900) loss: 8.623202
(Iteration 521 / 4900) loss: 8.528508
(Iteration 531 / 4900) loss: 8.504655
(Iteration 541 / 4900) loss: 8.440586
(Iteration 551 / 4900) loss: 8.415498
(Iteration 561 / 4900) loss: 8.140516
(Iteration 571 / 4900) loss: 8.133130
(Iteration 581 / 4900) loss: 8.218553
(Iteration 591 / 4900) loss: 8.096985
(Iteration 601 / 4900) loss: 8.069734
(Iteration 611 / 4900) loss: 7.893843
(Iteration 621 / 4900) loss: 7.753526
(Iteration 631 / 4900) loss: 7.744377
(Iteration 641 / 4900) loss: 7.603955
(Iteration 651 / 4900) loss: 7.681747
(Iteration 661 / 4900) loss: 7.573497
(Iteration 671 / 4900) loss: 7.587098
(Iteration 681 / 4900) loss: 7.421574
(Iteration 691 / 4900) loss: 7.296259
(Iteration 701 / 4900) loss: 7.464697
(Iteration 711 / 4900) loss: 7.252852
(Iteration 721 / 4900) loss: 7.339172
(Iteration 731 / 4900) loss: 7.178400
(Iteration 741 / 4900) loss: 7.073019
(Iteration 751 / 4900) loss: 7.241789
(Iteration 761 / 4900) loss: 6.903308
(Iteration 771 / 4900) loss: 6.958415
(Iteration 781 / 4900) loss: 6.959964
(Iteration 791 / 4900) loss: 6.802635
(Iteration 801 / 4900) loss: 6.559536
(Iteration 811 / 4900) loss: 6.751965
(Iteration 821 / 4900) loss: 6.841544
(Iteration 831 / 4900) loss: 6.754257
(Iteration 841 / 4900) loss: 6.562527
(Iteration 851 / 4900) loss: 6.590150
(Iteration 861 / 4900) loss: 6.503537
(Iteration 871 / 4900) loss: 6.396270

(Iteration 881 / 4900) loss: 6.493344
(Iteration 891 / 4900) loss: 6.289538
(Iteration 901 / 4900) loss: 6.249118
(Iteration 911 / 4900) loss: 6.277627
(Iteration 921 / 4900) loss: 6.177942
(Iteration 931 / 4900) loss: 6.142603
(Iteration 941 / 4900) loss: 5.995083
(Iteration 951 / 4900) loss: 5.918794
(Iteration 961 / 4900) loss: 6.064577
(Iteration 971 / 4900) loss: 5.894160
(Epoch 2 / 10) train acc: 0.449000; val_acc: 0.450000
(Iteration 981 / 4900) loss: 5.938206
(Iteration 991 / 4900) loss: 5.844028
(Iteration 1001 / 4900) loss: 5.995651
(Iteration 1011 / 4900) loss: 5.711307
(Iteration 1021 / 4900) loss: 5.625174
(Iteration 1031 / 4900) loss: 5.496967
(Iteration 1041 / 4900) loss: 5.501075
(Iteration 1051 / 4900) loss: 5.679336
(Iteration 1061 / 4900) loss: 5.472335
(Iteration 1071 / 4900) loss: 5.435863
(Iteration 1081 / 4900) loss: 5.649263
(Iteration 1091 / 4900) loss: 5.472952
(Iteration 1101 / 4900) loss: 5.259068
(Iteration 1111 / 4900) loss: 5.324463
(Iteration 1121 / 4900) loss: 5.241094
(Iteration 1131 / 4900) loss: 5.204853
(Iteration 1141 / 4900) loss: 5.216389
(Iteration 1151 / 4900) loss: 5.283409
(Iteration 1161 / 4900) loss: 5.175273
(Iteration 1171 / 4900) loss: 5.071160
(Iteration 1181 / 4900) loss: 5.046246
(Iteration 1191 / 4900) loss: 4.970761
(Iteration 1201 / 4900) loss: 5.243748
(Iteration 1211 / 4900) loss: 5.085856
(Iteration 1221 / 4900) loss: 5.068268
(Iteration 1231 / 4900) loss: 4.811699
(Iteration 1241 / 4900) loss: 4.998644
(Iteration 1251 / 4900) loss: 4.929046
(Iteration 1261 / 4900) loss: 4.980539
(Iteration 1271 / 4900) loss: 4.868526
(Iteration 1281 / 4900) loss: 4.859050
(Iteration 1291 / 4900) loss: 4.859291
(Iteration 1301 / 4900) loss: 4.690207
(Iteration 1311 / 4900) loss: 4.369216
(Iteration 1321 / 4900) loss: 4.649294
(Iteration 1331 / 4900) loss: 4.628230
(Iteration 1341 / 4900) loss: 4.598076

(Iteration 1351 / 4900) loss: 4.559622
(Iteration 1361 / 4900) loss: 4.673052
(Iteration 1371 / 4900) loss: 4.508673
(Iteration 1381 / 4900) loss: 4.371158
(Iteration 1391 / 4900) loss: 4.387171
(Iteration 1401 / 4900) loss: 4.388240
(Iteration 1411 / 4900) loss: 4.188988
(Iteration 1421 / 4900) loss: 4.262877
(Iteration 1431 / 4900) loss: 4.567482
(Iteration 1441 / 4900) loss: 4.370824
(Iteration 1451 / 4900) loss: 4.356736
(Iteration 1461 / 4900) loss: 4.148722
(Epoch 3 / 10) train acc: 0.483000; val_acc: 0.459000
(Iteration 1471 / 4900) loss: 4.097899
(Iteration 1481 / 4900) loss: 4.202987
(Iteration 1491 / 4900) loss: 4.150164
(Iteration 1501 / 4900) loss: 4.212076
(Iteration 1511 / 4900) loss: 4.213498
(Iteration 1521 / 4900) loss: 4.022372
(Iteration 1531 / 4900) loss: 4.081263
(Iteration 1541 / 4900) loss: 4.027801
(Iteration 1551 / 4900) loss: 3.936915
(Iteration 1561 / 4900) loss: 3.883623
(Iteration 1571 / 4900) loss: 3.887338
(Iteration 1581 / 4900) loss: 3.840766
(Iteration 1591 / 4900) loss: 3.781112
(Iteration 1601 / 4900) loss: 3.888201
(Iteration 1611 / 4900) loss: 3.814892
(Iteration 1621 / 4900) loss: 3.709475
(Iteration 1631 / 4900) loss: 3.869601
(Iteration 1641 / 4900) loss: 3.711695
(Iteration 1651 / 4900) loss: 3.852577
(Iteration 1661 / 4900) loss: 3.648840
(Iteration 1671 / 4900) loss: 3.802220
(Iteration 1681 / 4900) loss: 3.754788
(Iteration 1691 / 4900) loss: 3.739387
(Iteration 1701 / 4900) loss: 3.765209
(Iteration 1711 / 4900) loss: 3.585665
(Iteration 1721 / 4900) loss: 3.550466
(Iteration 1731 / 4900) loss: 3.444384
(Iteration 1741 / 4900) loss: 3.688976
(Iteration 1751 / 4900) loss: 3.776184
(Iteration 1761 / 4900) loss: 3.431694
(Iteration 1771 / 4900) loss: 3.573048
(Iteration 1781 / 4900) loss: 3.486366
(Iteration 1791 / 4900) loss: 3.388739
(Iteration 1801 / 4900) loss: 3.466451
(Iteration 1811 / 4900) loss: 3.377596

(Iteration 1821 / 4900) loss: 3.494737
(Iteration 1831 / 4900) loss: 3.356660
(Iteration 1841 / 4900) loss: 3.333428
(Iteration 1851 / 4900) loss: 3.352641
(Iteration 1861 / 4900) loss: 3.436673
(Iteration 1871 / 4900) loss: 3.414054
(Iteration 1881 / 4900) loss: 3.501688
(Iteration 1891 / 4900) loss: 3.277677
(Iteration 1901 / 4900) loss: 3.290037
(Iteration 1911 / 4900) loss: 3.326131
(Iteration 1921 / 4900) loss: 3.312574
(Iteration 1931 / 4900) loss: 3.088551
(Iteration 1941 / 4900) loss: 3.266179
(Iteration 1951 / 4900) loss: 3.384234
(Epoch 4 / 10) train acc: 0.501000; val_acc: 0.464000
(Iteration 1961 / 4900) loss: 3.459043
(Iteration 1971 / 4900) loss: 3.179969
(Iteration 1981 / 4900) loss: 3.107273
(Iteration 1991 / 4900) loss: 3.043423
(Iteration 2001 / 4900) loss: 2.976852
(Iteration 2011 / 4900) loss: 2.936474
(Iteration 2021 / 4900) loss: 2.977105
(Iteration 2031 / 4900) loss: 3.051216
(Iteration 2041 / 4900) loss: 2.871116
(Iteration 2051 / 4900) loss: 2.997437
(Iteration 2061 / 4900) loss: 3.143728
(Iteration 2071 / 4900) loss: 3.004803
(Iteration 2081 / 4900) loss: 3.010440
(Iteration 2091 / 4900) loss: 2.922936
(Iteration 2101 / 4900) loss: 2.940728
(Iteration 2111 / 4900) loss: 2.818236
(Iteration 2121 / 4900) loss: 3.057994
(Iteration 2131 / 4900) loss: 2.763880
(Iteration 2141 / 4900) loss: 3.118986
(Iteration 2151 / 4900) loss: 2.819378
(Iteration 2161 / 4900) loss: 2.801153
(Iteration 2171 / 4900) loss: 2.870282
(Iteration 2181 / 4900) loss: 2.690951
(Iteration 2191 / 4900) loss: 2.634393
(Iteration 2201 / 4900) loss: 2.823241
(Iteration 2211 / 4900) loss: 2.704861
(Iteration 2221 / 4900) loss: 2.786015
(Iteration 2231 / 4900) loss: 2.935901
(Iteration 2241 / 4900) loss: 2.856445
(Iteration 2251 / 4900) loss: 2.907901
(Iteration 2261 / 4900) loss: 2.565630
(Iteration 2271 / 4900) loss: 2.827611
(Iteration 2281 / 4900) loss: 2.851751

(Iteration 2291 / 4900) loss: 2.618146
(Iteration 2301 / 4900) loss: 2.675830
(Iteration 2311 / 4900) loss: 2.761522
(Iteration 2321 / 4900) loss: 2.521166
(Iteration 2331 / 4900) loss: 2.673851
(Iteration 2341 / 4900) loss: 2.624789
(Iteration 2351 / 4900) loss: 2.585308
(Iteration 2361 / 4900) loss: 2.690116
(Iteration 2371 / 4900) loss: 2.579623
(Iteration 2381 / 4900) loss: 2.614082
(Iteration 2391 / 4900) loss: 2.731081
(Iteration 2401 / 4900) loss: 2.624162
(Iteration 2411 / 4900) loss: 2.380180
(Iteration 2421 / 4900) loss: 2.565143
(Iteration 2431 / 4900) loss: 2.675891
(Iteration 2441 / 4900) loss: 2.420533
(Epoch 5 / 10) train acc: 0.469000; val_acc: 0.472000
(Iteration 2451 / 4900) loss: 2.690454
(Iteration 2461 / 4900) loss: 2.391141
(Iteration 2471 / 4900) loss: 2.356608
(Iteration 2481 / 4900) loss: 2.463313
(Iteration 2491 / 4900) loss: 2.579916
(Iteration 2501 / 4900) loss: 2.501184
(Iteration 2511 / 4900) loss: 2.380698
(Iteration 2521 / 4900) loss: 2.497429
(Iteration 2531 / 4900) loss: 2.463244
(Iteration 2541 / 4900) loss: 2.517669
(Iteration 2551 / 4900) loss: 2.398385
(Iteration 2561 / 4900) loss: 2.388333
(Iteration 2571 / 4900) loss: 2.382099
(Iteration 2581 / 4900) loss: 2.157723
(Iteration 2591 / 4900) loss: 2.340986
(Iteration 2601 / 4900) loss: 2.343534
(Iteration 2611 / 4900) loss: 2.590582
(Iteration 2621 / 4900) loss: 2.327598
(Iteration 2631 / 4900) loss: 2.406373
(Iteration 2641 / 4900) loss: 2.193580
(Iteration 2651 / 4900) loss: 2.328167
(Iteration 2661 / 4900) loss: 2.204311
(Iteration 2671 / 4900) loss: 2.452048
(Iteration 2681 / 4900) loss: 2.242722
(Iteration 2691 / 4900) loss: 2.427503
(Iteration 2701 / 4900) loss: 2.213045
(Iteration 2711 / 4900) loss: 2.357302
(Iteration 2721 / 4900) loss: 2.433159
(Iteration 2731 / 4900) loss: 2.469051
(Iteration 2741 / 4900) loss: 2.158850
(Iteration 2751 / 4900) loss: 2.157786

(Iteration 2761 / 4900) loss: 2.328191
(Iteration 2771 / 4900) loss: 2.130382
(Iteration 2781 / 4900) loss: 2.334143
(Iteration 2791 / 4900) loss: 2.280603
(Iteration 2801 / 4900) loss: 2.272896
(Iteration 2811 / 4900) loss: 2.139038
(Iteration 2821 / 4900) loss: 2.250149
(Iteration 2831 / 4900) loss: 2.446232
(Iteration 2841 / 4900) loss: 2.139894
(Iteration 2851 / 4900) loss: 2.199743
(Iteration 2861 / 4900) loss: 2.099746
(Iteration 2871 / 4900) loss: 2.030193
(Iteration 2881 / 4900) loss: 2.356836
(Iteration 2891 / 4900) loss: 2.242350
(Iteration 2901 / 4900) loss: 2.005509
(Iteration 2911 / 4900) loss: 2.063781
(Iteration 2921 / 4900) loss: 2.112202
(Iteration 2931 / 4900) loss: 2.067082
(Epoch 6 / 10) train acc: 0.490000; val_acc: 0.486000
(Iteration 2941 / 4900) loss: 2.147396
(Iteration 2951 / 4900) loss: 2.260730
(Iteration 2961 / 4900) loss: 2.322408
(Iteration 2971 / 4900) loss: 2.055329
(Iteration 2981 / 4900) loss: 2.113732
(Iteration 2991 / 4900) loss: 2.078769
(Iteration 3001 / 4900) loss: 2.076381
(Iteration 3011 / 4900) loss: 2.204471
(Iteration 3021 / 4900) loss: 2.057119
(Iteration 3031 / 4900) loss: 2.123089
(Iteration 3041 / 4900) loss: 2.104896
(Iteration 3051 / 4900) loss: 1.991226
(Iteration 3061 / 4900) loss: 1.975701
(Iteration 3071 / 4900) loss: 1.930349
(Iteration 3081 / 4900) loss: 1.937695
(Iteration 3091 / 4900) loss: 2.126250
(Iteration 3101 / 4900) loss: 2.031332
(Iteration 3111 / 4900) loss: 2.138309
(Iteration 3121 / 4900) loss: 2.094693
(Iteration 3131 / 4900) loss: 1.860991
(Iteration 3141 / 4900) loss: 2.142686
(Iteration 3151 / 4900) loss: 2.092222
(Iteration 3161 / 4900) loss: 2.018931
(Iteration 3171 / 4900) loss: 1.957695
(Iteration 3181 / 4900) loss: 1.964852
(Iteration 3191 / 4900) loss: 2.088814
(Iteration 3201 / 4900) loss: 1.924259
(Iteration 3211 / 4900) loss: 1.870023
(Iteration 3221 / 4900) loss: 1.810704

(Iteration 3231 / 4900) loss: 1.967996
(Iteration 3241 / 4900) loss: 2.005076
(Iteration 3251 / 4900) loss: 1.821501
(Iteration 3261 / 4900) loss: 1.940167
(Iteration 3271 / 4900) loss: 2.099844
(Iteration 3281 / 4900) loss: 1.951703
(Iteration 3291 / 4900) loss: 1.943462
(Iteration 3301 / 4900) loss: 1.839084
(Iteration 3311 / 4900) loss: 1.841205
(Iteration 3321 / 4900) loss: 2.066767
(Iteration 3331 / 4900) loss: 1.920561
(Iteration 3341 / 4900) loss: 1.997770
(Iteration 3351 / 4900) loss: 1.756867
(Iteration 3361 / 4900) loss: 2.114333
(Iteration 3371 / 4900) loss: 1.957617
(Iteration 3381 / 4900) loss: 1.845779
(Iteration 3391 / 4900) loss: 1.932920
(Iteration 3401 / 4900) loss: 1.713751
(Iteration 3411 / 4900) loss: 1.926255
(Iteration 3421 / 4900) loss: 2.146776
(Epoch 7 / 10) train acc: 0.540000; val_acc: 0.503000
(Iteration 3431 / 4900) loss: 1.935413
(Iteration 3441 / 4900) loss: 1.970535
(Iteration 3451 / 4900) loss: 1.823747
(Iteration 3461 / 4900) loss: 1.946521
(Iteration 3471 / 4900) loss: 1.732191
(Iteration 3481 / 4900) loss: 1.827014
(Iteration 3491 / 4900) loss: 1.849129
(Iteration 3501 / 4900) loss: 1.901268
(Iteration 3511 / 4900) loss: 1.868119
(Iteration 3521 / 4900) loss: 1.813957
(Iteration 3531 / 4900) loss: 1.827925
(Iteration 3541 / 4900) loss: 1.678180
(Iteration 3551 / 4900) loss: 1.956036
(Iteration 3561 / 4900) loss: 1.935665
(Iteration 3571 / 4900) loss: 1.811212
(Iteration 3581 / 4900) loss: 1.846843
(Iteration 3591 / 4900) loss: 2.041065
(Iteration 3601 / 4900) loss: 1.752126
(Iteration 3611 / 4900) loss: 1.750585
(Iteration 3621 / 4900) loss: 1.765751
(Iteration 3631 / 4900) loss: 1.855585
(Iteration 3641 / 4900) loss: 1.815596
(Iteration 3651 / 4900) loss: 1.841152
(Iteration 3661 / 4900) loss: 1.671215
(Iteration 3671 / 4900) loss: 1.860435
(Iteration 3681 / 4900) loss: 1.751276
(Iteration 3691 / 4900) loss: 1.841803

(Iteration 3701 / 4900) loss: 1.932285
(Iteration 3711 / 4900) loss: 1.896882
(Iteration 3721 / 4900) loss: 1.678466
(Iteration 3731 / 4900) loss: 1.733332
(Iteration 3741 / 4900) loss: 1.738808
(Iteration 3751 / 4900) loss: 1.874552
(Iteration 3761 / 4900) loss: 1.708116
(Iteration 3771 / 4900) loss: 1.482612
(Iteration 3781 / 4900) loss: 1.687947
(Iteration 3791 / 4900) loss: 1.831090
(Iteration 3801 / 4900) loss: 1.695215
(Iteration 3811 / 4900) loss: 1.677359
(Iteration 3821 / 4900) loss: 1.875571
(Iteration 3831 / 4900) loss: 1.679121
(Iteration 3841 / 4900) loss: 1.839907
(Iteration 3851 / 4900) loss: 1.680710
(Iteration 3861 / 4900) loss: 1.736416
(Iteration 3871 / 4900) loss: 1.815990
(Iteration 3881 / 4900) loss: 1.685324
(Iteration 3891 / 4900) loss: 1.766178
(Iteration 3901 / 4900) loss: 1.605881
(Iteration 3911 / 4900) loss: 1.583431
(Epoch 8 / 10) train acc: 0.535000; val_acc: 0.475000
(Iteration 3921 / 4900) loss: 1.743489
(Iteration 3931 / 4900) loss: 1.734753
(Iteration 3941 / 4900) loss: 1.810547
(Iteration 3951 / 4900) loss: 1.635343
(Iteration 3961 / 4900) loss: 1.648624
(Iteration 3971 / 4900) loss: 1.777953
(Iteration 3981 / 4900) loss: 1.661880
(Iteration 3991 / 4900) loss: 1.702688
(Iteration 4001 / 4900) loss: 1.561372
(Iteration 4011 / 4900) loss: 1.703150
(Iteration 4021 / 4900) loss: 1.718848
(Iteration 4031 / 4900) loss: 1.608839
(Iteration 4041 / 4900) loss: 1.580262
(Iteration 4051 / 4900) loss: 1.687884
(Iteration 4061 / 4900) loss: 1.765798
(Iteration 4071 / 4900) loss: 1.634781
(Iteration 4081 / 4900) loss: 1.722229
(Iteration 4091 / 4900) loss: 1.828494
(Iteration 4101 / 4900) loss: 1.616314
(Iteration 4111 / 4900) loss: 1.673942
(Iteration 4121 / 4900) loss: 1.764314
(Iteration 4131 / 4900) loss: 1.704997
(Iteration 4141 / 4900) loss: 1.771653
(Iteration 4151 / 4900) loss: 1.931344
(Iteration 4161 / 4900) loss: 1.848987

(Iteration 4171 / 4900) loss: 1.609662
(Iteration 4181 / 4900) loss: 1.783420
(Iteration 4191 / 4900) loss: 1.719323
(Iteration 4201 / 4900) loss: 1.818183
(Iteration 4211 / 4900) loss: 1.695946
(Iteration 4221 / 4900) loss: 1.615818
(Iteration 4231 / 4900) loss: 1.678978
(Iteration 4241 / 4900) loss: 1.704928
(Iteration 4251 / 4900) loss: 1.686824
(Iteration 4261 / 4900) loss: 1.713536
(Iteration 4271 / 4900) loss: 1.854956
(Iteration 4281 / 4900) loss: 1.792072
(Iteration 4291 / 4900) loss: 1.645052
(Iteration 4301 / 4900) loss: 1.787017
(Iteration 4311 / 4900) loss: 1.697385
(Iteration 4321 / 4900) loss: 1.775241
(Iteration 4331 / 4900) loss: 1.775835
(Iteration 4341 / 4900) loss: 1.775533
(Iteration 4351 / 4900) loss: 1.587769
(Iteration 4361 / 4900) loss: 1.602738
(Iteration 4371 / 4900) loss: 1.781761
(Iteration 4381 / 4900) loss: 1.518057
(Iteration 4391 / 4900) loss: 1.731836
(Iteration 4401 / 4900) loss: 1.530800
(Epoch 9 / 10) train acc: 0.578000; val_acc: 0.475000
(Iteration 4411 / 4900) loss: 1.675745
(Iteration 4421 / 4900) loss: 1.836844
(Iteration 4431 / 4900) loss: 1.488795
(Iteration 4441 / 4900) loss: 1.884461
(Iteration 4451 / 4900) loss: 1.590660
(Iteration 4461 / 4900) loss: 1.734832
(Iteration 4471 / 4900) loss: 1.660236
(Iteration 4481 / 4900) loss: 1.382602
(Iteration 4491 / 4900) loss: 1.598552
(Iteration 4501 / 4900) loss: 1.595762
(Iteration 4511 / 4900) loss: 1.650344
(Iteration 4521 / 4900) loss: 1.574396
(Iteration 4531 / 4900) loss: 1.584248
(Iteration 4541 / 4900) loss: 1.650981
(Iteration 4551 / 4900) loss: 1.493766
(Iteration 4561 / 4900) loss: 1.876589
(Iteration 4571 / 4900) loss: 1.846185
(Iteration 4581 / 4900) loss: 1.826528
(Iteration 4591 / 4900) loss: 1.626806
(Iteration 4601 / 4900) loss: 1.618162
(Iteration 4611 / 4900) loss: 1.776087
(Iteration 4621 / 4900) loss: 1.703740
(Iteration 4631 / 4900) loss: 1.832554

```
(Iteration 4641 / 4900) loss: 1.534299
(Iteration 4651 / 4900) loss: 1.654393
(Iteration 4661 / 4900) loss: 1.432543
(Iteration 4671 / 4900) loss: 1.627118
(Iteration 4681 / 4900) loss: 1.608307
(Iteration 4691 / 4900) loss: 1.556269
(Iteration 4701 / 4900) loss: 1.569141
(Iteration 4711 / 4900) loss: 1.625043
(Iteration 4721 / 4900) loss: 1.606670
(Iteration 4731 / 4900) loss: 1.763936
(Iteration 4741 / 4900) loss: 1.728297
(Iteration 4751 / 4900) loss: 1.546293
(Iteration 4761 / 4900) loss: 1.556788
(Iteration 4771 / 4900) loss: 1.656464
(Iteration 4781 / 4900) loss: 1.809289
(Iteration 4791 / 4900) loss: 1.586060
(Iteration 4801 / 4900) loss: 1.752674
(Iteration 4811 / 4900) loss: 1.611331
(Iteration 4821 / 4900) loss: 1.721972
(Iteration 4831 / 4900) loss: 1.641282
(Iteration 4841 / 4900) loss: 1.483835
(Iteration 4851 / 4900) loss: 1.725171
(Iteration 4861 / 4900) loss: 1.513358
(Iteration 4871 / 4900) loss: 1.522193
(Iteration 4881 / 4900) loss: 1.652130
(Iteration 4891 / 4900) loss: 1.775364
(Epoch 10 / 10) train acc: 0.551000; val_acc: 0.512000
[3072, 100, 10]
(Iteration 1 / 4900) loss: 19.969652
(Epoch 0 / 10) train acc: 0.119000; val_acc: 0.129000
(Iteration 11 / 4900) loss: 17.910930
(Iteration 21 / 4900) loss: 17.599519
(Iteration 31 / 4900) loss: 17.210232
(Iteration 41 / 4900) loss: 16.893238
(Iteration 51 / 4900) loss: 16.643332
(Iteration 61 / 4900) loss: 16.535568
(Iteration 71 / 4900) loss: 16.286427
(Iteration 81 / 4900) loss: 16.154316
(Iteration 91 / 4900) loss: 15.705068
(Iteration 101 / 4900) loss: 15.904124
(Iteration 111 / 4900) loss: 15.788704
(Iteration 121 / 4900) loss: 15.232352
(Iteration 131 / 4900) loss: 15.304335
(Iteration 141 / 4900) loss: 15.280494
(Iteration 151 / 4900) loss: 15.097123
(Iteration 161 / 4900) loss: 14.961664
(Iteration 171 / 4900) loss: 14.882705
(Iteration 181 / 4900) loss: 14.666763
```


(Iteration 191 / 4900) loss: 14.418517
(Iteration 201 / 4900) loss: 14.334117
(Iteration 211 / 4900) loss: 14.301233
(Iteration 221 / 4900) loss: 14.258858
(Iteration 231 / 4900) loss: 14.083495
(Iteration 241 / 4900) loss: 13.835254
(Iteration 251 / 4900) loss: 13.981914
(Iteration 261 / 4900) loss: 13.580762
(Iteration 271 / 4900) loss: 13.392962
(Iteration 281 / 4900) loss: 13.309649
(Iteration 291 / 4900) loss: 13.386592
(Iteration 301 / 4900) loss: 13.048754
(Iteration 311 / 4900) loss: 13.021427
(Iteration 321 / 4900) loss: 13.024285
(Iteration 331 / 4900) loss: 12.736955
(Iteration 341 / 4900) loss: 12.559551
(Iteration 351 / 4900) loss: 12.551816
(Iteration 361 / 4900) loss: 12.367017
(Iteration 371 / 4900) loss: 12.448107
(Iteration 381 / 4900) loss: 12.326769
(Iteration 391 / 4900) loss: 12.068099
(Iteration 401 / 4900) loss: 12.167992
(Iteration 411 / 4900) loss: 11.755004
(Iteration 421 / 4900) loss: 11.867051
(Iteration 431 / 4900) loss: 11.497921
(Iteration 441 / 4900) loss: 11.467332
(Iteration 451 / 4900) loss: 11.404977
(Iteration 461 / 4900) loss: 11.164981
(Iteration 471 / 4900) loss: 11.331510
(Iteration 481 / 4900) loss: 11.184672
(Epoch 1 / 10) train acc: 0.414000; val_acc: 0.418000
(Iteration 491 / 4900) loss: 11.142430
(Iteration 501 / 4900) loss: 10.985490
(Iteration 511 / 4900) loss: 10.996054
(Iteration 521 / 4900) loss: 10.810692
(Iteration 531 / 4900) loss: 10.648419
(Iteration 541 / 4900) loss: 10.566614
(Iteration 551 / 4900) loss: 10.341601
(Iteration 561 / 4900) loss: 10.453595
(Iteration 571 / 4900) loss: 10.488936
(Iteration 581 / 4900) loss: 10.236080
(Iteration 591 / 4900) loss: 10.140152
(Iteration 601 / 4900) loss: 10.171465
(Iteration 611 / 4900) loss: 9.863047
(Iteration 621 / 4900) loss: 9.714131
(Iteration 631 / 4900) loss: 9.737817
(Iteration 641 / 4900) loss: 9.562521
(Iteration 651 / 4900) loss: 9.845058

(Iteration 661 / 4900) loss: 9.343305
(Iteration 671 / 4900) loss: 9.527677
(Iteration 681 / 4900) loss: 9.469457
(Iteration 691 / 4900) loss: 9.323889
(Iteration 701 / 4900) loss: 9.202865
(Iteration 711 / 4900) loss: 9.294707
(Iteration 721 / 4900) loss: 9.332792
(Iteration 731 / 4900) loss: 8.766095
(Iteration 741 / 4900) loss: 8.826768
(Iteration 751 / 4900) loss: 8.896346
(Iteration 761 / 4900) loss: 8.816113
(Iteration 771 / 4900) loss: 8.766097
(Iteration 781 / 4900) loss: 8.815383
(Iteration 791 / 4900) loss: 8.570272
(Iteration 801 / 4900) loss: 8.589501
(Iteration 811 / 4900) loss: 8.445095
(Iteration 821 / 4900) loss: 8.341177
(Iteration 831 / 4900) loss: 8.299542
(Iteration 841 / 4900) loss: 8.109231
(Iteration 851 / 4900) loss: 8.278753
(Iteration 861 / 4900) loss: 7.922106
(Iteration 871 / 4900) loss: 7.828173
(Iteration 881 / 4900) loss: 7.816324
(Iteration 891 / 4900) loss: 7.992747
(Iteration 901 / 4900) loss: 7.782163
(Iteration 911 / 4900) loss: 7.946956
(Iteration 921 / 4900) loss: 7.608738
(Iteration 931 / 4900) loss: 7.771130
(Iteration 941 / 4900) loss: 7.631451
(Iteration 951 / 4900) loss: 7.748477
(Iteration 961 / 4900) loss: 7.729550
(Iteration 971 / 4900) loss: 7.556277
(Epoch 2 / 10) train acc: 0.446000; val_acc: 0.449000
(Iteration 981 / 4900) loss: 7.540035
(Iteration 991 / 4900) loss: 7.501450
(Iteration 1001 / 4900) loss: 7.191758
(Iteration 1011 / 4900) loss: 7.205421
(Iteration 1021 / 4900) loss: 7.123696
(Iteration 1031 / 4900) loss: 7.164633
(Iteration 1041 / 4900) loss: 6.896850
(Iteration 1051 / 4900) loss: 7.020629
(Iteration 1061 / 4900) loss: 6.816375
(Iteration 1071 / 4900) loss: 6.785112
(Iteration 1081 / 4900) loss: 6.783520
(Iteration 1091 / 4900) loss: 6.633254
(Iteration 1101 / 4900) loss: 6.571052
(Iteration 1111 / 4900) loss: 6.617481
(Iteration 1121 / 4900) loss: 6.584521

(Iteration 1131 / 4900) loss: 6.354380
(Iteration 1141 / 4900) loss: 6.540647
(Iteration 1151 / 4900) loss: 6.481034
(Iteration 1161 / 4900) loss: 6.275787
(Iteration 1171 / 4900) loss: 6.389690
(Iteration 1181 / 4900) loss: 6.399293
(Iteration 1191 / 4900) loss: 6.153651
(Iteration 1201 / 4900) loss: 6.151776
(Iteration 1211 / 4900) loss: 6.314670
(Iteration 1221 / 4900) loss: 6.128239
(Iteration 1231 / 4900) loss: 5.882951
(Iteration 1241 / 4900) loss: 6.059318
(Iteration 1251 / 4900) loss: 5.900347
(Iteration 1261 / 4900) loss: 5.915342
(Iteration 1271 / 4900) loss: 5.848702
(Iteration 1281 / 4900) loss: 5.881053
(Iteration 1291 / 4900) loss: 5.880453
(Iteration 1301 / 4900) loss: 5.935147
(Iteration 1311 / 4900) loss: 5.635606
(Iteration 1321 / 4900) loss: 5.619482
(Iteration 1331 / 4900) loss: 5.557492
(Iteration 1341 / 4900) loss: 5.860804
(Iteration 1351 / 4900) loss: 5.342924
(Iteration 1361 / 4900) loss: 5.465712
(Iteration 1371 / 4900) loss: 5.468179
(Iteration 1381 / 4900) loss: 5.457202
(Iteration 1391 / 4900) loss: 5.225572
(Iteration 1401 / 4900) loss: 5.182870
(Iteration 1411 / 4900) loss: 5.489920
(Iteration 1421 / 4900) loss: 5.041038
(Iteration 1431 / 4900) loss: 5.328541
(Iteration 1441 / 4900) loss: 5.093867
(Iteration 1451 / 4900) loss: 4.881827
(Iteration 1461 / 4900) loss: 4.949661
(Epoch 3 / 10) train acc: 0.476000; val_acc: 0.454000
(Iteration 1471 / 4900) loss: 5.260479
(Iteration 1481 / 4900) loss: 4.959245
(Iteration 1491 / 4900) loss: 4.953393
(Iteration 1501 / 4900) loss: 5.113067
(Iteration 1511 / 4900) loss: 4.794341
(Iteration 1521 / 4900) loss: 5.095970
(Iteration 1531 / 4900) loss: 4.930543
(Iteration 1541 / 4900) loss: 4.713758
(Iteration 1551 / 4900) loss: 4.856503
(Iteration 1561 / 4900) loss: 4.806204
(Iteration 1571 / 4900) loss: 4.692843
(Iteration 1581 / 4900) loss: 4.675027
(Iteration 1591 / 4900) loss: 4.739186

(Iteration 1601 / 4900) loss: 4.773986
(Iteration 1611 / 4900) loss: 4.605448
(Iteration 1621 / 4900) loss: 4.597430
(Iteration 1631 / 4900) loss: 4.455373
(Iteration 1641 / 4900) loss: 4.485009
(Iteration 1651 / 4900) loss: 4.708115
(Iteration 1661 / 4900) loss: 4.616059
(Iteration 1671 / 4900) loss: 4.468410
(Iteration 1681 / 4900) loss: 4.480976
(Iteration 1691 / 4900) loss: 4.371126
(Iteration 1701 / 4900) loss: 4.259136
(Iteration 1711 / 4900) loss: 4.424915
(Iteration 1721 / 4900) loss: 4.216385
(Iteration 1731 / 4900) loss: 4.208714
(Iteration 1741 / 4900) loss: 4.208455
(Iteration 1751 / 4900) loss: 4.083254
(Iteration 1761 / 4900) loss: 4.184089
(Iteration 1771 / 4900) loss: 4.120522
(Iteration 1781 / 4900) loss: 4.234023
(Iteration 1791 / 4900) loss: 4.126803
(Iteration 1801 / 4900) loss: 3.949510
(Iteration 1811 / 4900) loss: 4.213698
(Iteration 1821 / 4900) loss: 3.983670
(Iteration 1831 / 4900) loss: 3.972976
(Iteration 1841 / 4900) loss: 3.889866
(Iteration 1851 / 4900) loss: 3.980828
(Iteration 1861 / 4900) loss: 3.877266
(Iteration 1871 / 4900) loss: 3.927041
(Iteration 1881 / 4900) loss: 3.954859
(Iteration 1891 / 4900) loss: 3.840177
(Iteration 1901 / 4900) loss: 3.703063
(Iteration 1911 / 4900) loss: 3.784274
(Iteration 1921 / 4900) loss: 3.661125
(Iteration 1931 / 4900) loss: 3.665429
(Iteration 1941 / 4900) loss: 3.861913
(Iteration 1951 / 4900) loss: 3.762587
(Epoch 4 / 10) train acc: 0.473000; val_acc: 0.460000
(Iteration 1961 / 4900) loss: 3.674423
(Iteration 1971 / 4900) loss: 3.426795
(Iteration 1981 / 4900) loss: 3.541926
(Iteration 1991 / 4900) loss: 3.722667
(Iteration 2001 / 4900) loss: 3.703114
(Iteration 2011 / 4900) loss: 3.604524
(Iteration 2021 / 4900) loss: 3.478735
(Iteration 2031 / 4900) loss: 3.472980
(Iteration 2041 / 4900) loss: 3.562988
(Iteration 2051 / 4900) loss: 3.447603
(Iteration 2061 / 4900) loss: 3.543475

(Iteration 2071 / 4900) loss: 3.497663
(Iteration 2081 / 4900) loss: 3.504628
(Iteration 2091 / 4900) loss: 3.401848
(Iteration 2101 / 4900) loss: 3.375801
(Iteration 2111 / 4900) loss: 3.389166
(Iteration 2121 / 4900) loss: 3.361076
(Iteration 2131 / 4900) loss: 3.295218
(Iteration 2141 / 4900) loss: 3.379735
(Iteration 2151 / 4900) loss: 3.388879
(Iteration 2161 / 4900) loss: 3.344373
(Iteration 2171 / 4900) loss: 3.250633
(Iteration 2181 / 4900) loss: 3.239239
(Iteration 2191 / 4900) loss: 3.210144
(Iteration 2201 / 4900) loss: 3.413237
(Iteration 2211 / 4900) loss: 3.041461
(Iteration 2221 / 4900) loss: 3.154345
(Iteration 2231 / 4900) loss: 2.994706
(Iteration 2241 / 4900) loss: 3.097533
(Iteration 2251 / 4900) loss: 3.100350
(Iteration 2261 / 4900) loss: 3.105808
(Iteration 2271 / 4900) loss: 3.238787
(Iteration 2281 / 4900) loss: 3.083206
(Iteration 2291 / 4900) loss: 3.026341
(Iteration 2301 / 4900) loss: 3.029498
(Iteration 2311 / 4900) loss: 2.801566
(Iteration 2321 / 4900) loss: 3.048572
(Iteration 2331 / 4900) loss: 2.959843
(Iteration 2341 / 4900) loss: 3.059553
(Iteration 2351 / 4900) loss: 2.887118
(Iteration 2361 / 4900) loss: 2.891981
(Iteration 2371 / 4900) loss: 2.965660
(Iteration 2381 / 4900) loss: 2.950990
(Iteration 2391 / 4900) loss: 2.967055
(Iteration 2401 / 4900) loss: 2.921026
(Iteration 2411 / 4900) loss: 2.947966
(Iteration 2421 / 4900) loss: 2.848919
(Iteration 2431 / 4900) loss: 2.965048
(Iteration 2441 / 4900) loss: 2.846254
(Epoch 5 / 10) train acc: 0.516000; val_acc: 0.486000
(Iteration 2451 / 4900) loss: 2.833870
(Iteration 2461 / 4900) loss: 2.864035
(Iteration 2471 / 4900) loss: 2.787332
(Iteration 2481 / 4900) loss: 2.899850
(Iteration 2491 / 4900) loss: 2.763212
(Iteration 2501 / 4900) loss: 2.663921
(Iteration 2511 / 4900) loss: 2.802801
(Iteration 2521 / 4900) loss: 2.616525
(Iteration 2531 / 4900) loss: 2.715648

(Iteration 2541 / 4900) loss: 2.732247
(Iteration 2551 / 4900) loss: 2.665328
(Iteration 2561 / 4900) loss: 2.611848
(Iteration 2571 / 4900) loss: 2.628225
(Iteration 2581 / 4900) loss: 2.684238
(Iteration 2591 / 4900) loss: 2.739176
(Iteration 2601 / 4900) loss: 2.673656
(Iteration 2611 / 4900) loss: 2.733609
(Iteration 2621 / 4900) loss: 2.611630
(Iteration 2631 / 4900) loss: 2.572282
(Iteration 2641 / 4900) loss: 2.715918
(Iteration 2651 / 4900) loss: 2.739398
(Iteration 2661 / 4900) loss: 2.487793
(Iteration 2671 / 4900) loss: 2.591922
(Iteration 2681 / 4900) loss: 2.632135
(Iteration 2691 / 4900) loss: 2.452727
(Iteration 2701 / 4900) loss: 2.692677
(Iteration 2711 / 4900) loss: 2.371725
(Iteration 2721 / 4900) loss: 2.538296
(Iteration 2731 / 4900) loss: 2.334296
(Iteration 2741 / 4900) loss: 2.704140
(Iteration 2751 / 4900) loss: 2.492912
(Iteration 2761 / 4900) loss: 2.631581
(Iteration 2771 / 4900) loss: 2.369792
(Iteration 2781 / 4900) loss: 2.531437
(Iteration 2791 / 4900) loss: 2.498092
(Iteration 2801 / 4900) loss: 2.387276
(Iteration 2811 / 4900) loss: 2.682209
(Iteration 2821 / 4900) loss: 2.507492
(Iteration 2831 / 4900) loss: 2.364775
(Iteration 2841 / 4900) loss: 2.186009
(Iteration 2851 / 4900) loss: 2.286835
(Iteration 2861 / 4900) loss: 2.371368
(Iteration 2871 / 4900) loss: 2.438725
(Iteration 2881 / 4900) loss: 2.565858
(Iteration 2891 / 4900) loss: 2.541166
(Iteration 2901 / 4900) loss: 2.367862
(Iteration 2911 / 4900) loss: 2.172605
(Iteration 2921 / 4900) loss: 2.446317
(Iteration 2931 / 4900) loss: 2.322109
(Epoch 6 / 10) train acc: 0.514000; val_acc: 0.489000
(Iteration 2941 / 4900) loss: 2.395294
(Iteration 2951 / 4900) loss: 2.419484
(Iteration 2961 / 4900) loss: 2.221272
(Iteration 2971 / 4900) loss: 2.302176
(Iteration 2981 / 4900) loss: 2.291177
(Iteration 2991 / 4900) loss: 2.121861
(Iteration 3001 / 4900) loss: 2.200309

(Iteration 3011 / 4900) loss: 2.093930
(Iteration 3021 / 4900) loss: 2.427210
(Iteration 3031 / 4900) loss: 2.353273
(Iteration 3041 / 4900) loss: 2.344799
(Iteration 3051 / 4900) loss: 2.385714
(Iteration 3061 / 4900) loss: 2.232494
(Iteration 3071 / 4900) loss: 2.192399
(Iteration 3081 / 4900) loss: 2.242355
(Iteration 3091 / 4900) loss: 2.052682
(Iteration 3101 / 4900) loss: 2.045771
(Iteration 3111 / 4900) loss: 2.044016
(Iteration 3121 / 4900) loss: 2.082524
(Iteration 3131 / 4900) loss: 2.049038
(Iteration 3141 / 4900) loss: 2.254731
(Iteration 3151 / 4900) loss: 2.152143
(Iteration 3161 / 4900) loss: 2.256391
(Iteration 3171 / 4900) loss: 2.170240
(Iteration 3181 / 4900) loss: 2.161758
(Iteration 3191 / 4900) loss: 2.101652
(Iteration 3201 / 4900) loss: 2.131502
(Iteration 3211 / 4900) loss: 2.217932
(Iteration 3221 / 4900) loss: 2.156305
(Iteration 3231 / 4900) loss: 2.164889
(Iteration 3241 / 4900) loss: 2.172170
(Iteration 3251 / 4900) loss: 2.144654
(Iteration 3261 / 4900) loss: 2.050790
(Iteration 3271 / 4900) loss: 2.004149
(Iteration 3281 / 4900) loss: 2.186140
(Iteration 3291 / 4900) loss: 2.124236
(Iteration 3301 / 4900) loss: 1.946378
(Iteration 3311 / 4900) loss: 2.095886
(Iteration 3321 / 4900) loss: 2.102549
(Iteration 3331 / 4900) loss: 2.114751
(Iteration 3341 / 4900) loss: 1.956407
(Iteration 3351 / 4900) loss: 2.069781
(Iteration 3361 / 4900) loss: 1.976552
(Iteration 3371 / 4900) loss: 2.088758
(Iteration 3381 / 4900) loss: 1.999645
(Iteration 3391 / 4900) loss: 2.226962
(Iteration 3401 / 4900) loss: 2.153707
(Iteration 3411 / 4900) loss: 2.065382
(Iteration 3421 / 4900) loss: 2.026098
(Epoch 7 / 10) train acc: 0.529000; val_acc: 0.514000
(Iteration 3431 / 4900) loss: 1.924352
(Iteration 3441 / 4900) loss: 2.068299
(Iteration 3451 / 4900) loss: 2.152286
(Iteration 3461 / 4900) loss: 2.065492
(Iteration 3471 / 4900) loss: 1.921102

(Iteration 3481 / 4900) loss: 1.821509
(Iteration 3491 / 4900) loss: 2.194940
(Iteration 3501 / 4900) loss: 1.966174
(Iteration 3511 / 4900) loss: 1.724470
(Iteration 3521 / 4900) loss: 1.998987
(Iteration 3531 / 4900) loss: 1.888796
(Iteration 3541 / 4900) loss: 2.064935
(Iteration 3551 / 4900) loss: 1.805351
(Iteration 3561 / 4900) loss: 1.866970
(Iteration 3571 / 4900) loss: 1.934097
(Iteration 3581 / 4900) loss: 1.799281
(Iteration 3591 / 4900) loss: 1.872493
(Iteration 3601 / 4900) loss: 2.149437
(Iteration 3611 / 4900) loss: 1.914572
(Iteration 3621 / 4900) loss: 1.796573
(Iteration 3631 / 4900) loss: 1.855173
(Iteration 3641 / 4900) loss: 1.998832
(Iteration 3651 / 4900) loss: 1.819336
(Iteration 3661 / 4900) loss: 1.903333
(Iteration 3671 / 4900) loss: 1.823593
(Iteration 3681 / 4900) loss: 2.066985
(Iteration 3691 / 4900) loss: 2.074769
(Iteration 3701 / 4900) loss: 1.825404
(Iteration 3711 / 4900) loss: 1.892390
(Iteration 3721 / 4900) loss: 1.866557
(Iteration 3731 / 4900) loss: 2.049274
(Iteration 3741 / 4900) loss: 1.883973
(Iteration 3751 / 4900) loss: 1.829363
(Iteration 3761 / 4900) loss: 1.827123
(Iteration 3771 / 4900) loss: 1.719791
(Iteration 3781 / 4900) loss: 1.811372
(Iteration 3791 / 4900) loss: 1.860304
(Iteration 3801 / 4900) loss: 1.882898
(Iteration 3811 / 4900) loss: 1.962877
(Iteration 3821 / 4900) loss: 1.865228
(Iteration 3831 / 4900) loss: 1.916250
(Iteration 3841 / 4900) loss: 1.787199
(Iteration 3851 / 4900) loss: 1.867711
(Iteration 3861 / 4900) loss: 1.861366
(Iteration 3871 / 4900) loss: 1.673367
(Iteration 3881 / 4900) loss: 1.739753
(Iteration 3891 / 4900) loss: 1.892322
(Iteration 3901 / 4900) loss: 1.663844
(Iteration 3911 / 4900) loss: 1.715593
(Epoch 8 / 10) train acc: 0.551000; val_acc: 0.495000
(Iteration 3921 / 4900) loss: 1.707365
(Iteration 3931 / 4900) loss: 1.784963
(Iteration 3941 / 4900) loss: 1.895019

(Iteration 3951 / 4900) loss: 1.599772
(Iteration 3961 / 4900) loss: 1.763460
(Iteration 3971 / 4900) loss: 1.857975
(Iteration 3981 / 4900) loss: 1.917691
(Iteration 3991 / 4900) loss: 1.726436
(Iteration 4001 / 4900) loss: 1.862269
(Iteration 4011 / 4900) loss: 1.784882
(Iteration 4021 / 4900) loss: 1.505198
(Iteration 4031 / 4900) loss: 1.703950
(Iteration 4041 / 4900) loss: 1.999842
(Iteration 4051 / 4900) loss: 1.769459
(Iteration 4061 / 4900) loss: 1.927022
(Iteration 4071 / 4900) loss: 1.769735
(Iteration 4081 / 4900) loss: 1.822924
(Iteration 4091 / 4900) loss: 1.795063
(Iteration 4101 / 4900) loss: 1.704855
(Iteration 4111 / 4900) loss: 1.904038
(Iteration 4121 / 4900) loss: 1.696172
(Iteration 4131 / 4900) loss: 1.800506
(Iteration 4141 / 4900) loss: 1.811055
(Iteration 4151 / 4900) loss: 1.611631
(Iteration 4161 / 4900) loss: 1.703751
(Iteration 4171 / 4900) loss: 1.640666
(Iteration 4181 / 4900) loss: 1.681930
(Iteration 4191 / 4900) loss: 1.736151
(Iteration 4201 / 4900) loss: 1.695749
(Iteration 4211 / 4900) loss: 1.657195
(Iteration 4221 / 4900) loss: 1.687532
(Iteration 4231 / 4900) loss: 1.858841
(Iteration 4241 / 4900) loss: 1.834316
(Iteration 4251 / 4900) loss: 1.659708
(Iteration 4261 / 4900) loss: 1.710122
(Iteration 4271 / 4900) loss: 1.578381
(Iteration 4281 / 4900) loss: 1.948141
(Iteration 4291 / 4900) loss: 1.696576
(Iteration 4301 / 4900) loss: 1.632947
(Iteration 4311 / 4900) loss: 1.758933
(Iteration 4321 / 4900) loss: 1.830270
(Iteration 4331 / 4900) loss: 1.823758
(Iteration 4341 / 4900) loss: 1.697238
(Iteration 4351 / 4900) loss: 1.655538
(Iteration 4361 / 4900) loss: 1.802328
(Iteration 4371 / 4900) loss: 1.655496
(Iteration 4381 / 4900) loss: 1.705712
(Iteration 4391 / 4900) loss: 1.708969
(Iteration 4401 / 4900) loss: 1.762811
(Epoch 9 / 10) train acc: 0.510000; val_acc: 0.508000
(Iteration 4411 / 4900) loss: 1.602339

(Iteration 4421 / 4900) loss: 1.722949
(Iteration 4431 / 4900) loss: 1.713641
(Iteration 4441 / 4900) loss: 1.736121
(Iteration 4451 / 4900) loss: 1.657753
(Iteration 4461 / 4900) loss: 1.733440
(Iteration 4471 / 4900) loss: 1.686093
(Iteration 4481 / 4900) loss: 1.560429
(Iteration 4491 / 4900) loss: 1.593723
(Iteration 4501 / 4900) loss: 1.619667
(Iteration 4511 / 4900) loss: 1.863108
(Iteration 4521 / 4900) loss: 1.718570
(Iteration 4531 / 4900) loss: 1.466632
(Iteration 4541 / 4900) loss: 1.762400
(Iteration 4551 / 4900) loss: 1.814818
(Iteration 4561 / 4900) loss: 1.607734
(Iteration 4571 / 4900) loss: 1.778536
(Iteration 4581 / 4900) loss: 1.613957
(Iteration 4591 / 4900) loss: 1.720630
(Iteration 4601 / 4900) loss: 1.726105
(Iteration 4611 / 4900) loss: 1.678162
(Iteration 4621 / 4900) loss: 1.644455
(Iteration 4631 / 4900) loss: 1.775320
(Iteration 4641 / 4900) loss: 1.849526
(Iteration 4651 / 4900) loss: 1.714540
(Iteration 4661 / 4900) loss: 1.538750
(Iteration 4671 / 4900) loss: 1.718718
(Iteration 4681 / 4900) loss: 1.648744
(Iteration 4691 / 4900) loss: 1.713342
(Iteration 4701 / 4900) loss: 1.620783
(Iteration 4711 / 4900) loss: 1.529502
(Iteration 4721 / 4900) loss: 1.602889
(Iteration 4731 / 4900) loss: 1.553425
(Iteration 4741 / 4900) loss: 1.506195
(Iteration 4751 / 4900) loss: 1.587967
(Iteration 4761 / 4900) loss: 1.550274
(Iteration 4771 / 4900) loss: 1.512113
(Iteration 4781 / 4900) loss: 1.814492
(Iteration 4791 / 4900) loss: 1.656026
(Iteration 4801 / 4900) loss: 1.636440
(Iteration 4811 / 4900) loss: 1.841821
(Iteration 4821 / 4900) loss: 1.658743
(Iteration 4831 / 4900) loss: 1.506660
(Iteration 4841 / 4900) loss: 1.677667
(Iteration 4851 / 4900) loss: 1.569695
(Iteration 4861 / 4900) loss: 1.619733
(Iteration 4871 / 4900) loss: 1.655337
(Iteration 4881 / 4900) loss: 1.689747
(Iteration 4891 / 4900) loss: 1.686663

```
(Epoch 10 / 10) train acc: 0.552000; val_acc: 0.486000  
0.514
```

4 Test Your Model!

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

```
[ ]: 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:  0.514  
Test set accuracy:  0.501
```

BatchNormalization

November 4, 2022

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

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'cs231n/assignments/assignment1/'
FOLDERNAME = 'enpm809K/assignments/assignment2/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

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

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

Mounted at /content/drive

/content/drive/My Drive/enpm809K/assignments/assignment2/cs231n/datasets

/content/drive/My Drive/enpm809K/assignments/assignment2

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, proposed by [1] in 2015.

To understand the goal of batch normalization, it is important to first recognize that machine learning methods tend to perform better with input data consisting of uncorrelated features with zero mean and unit variance. When training a neural network, we can preprocess the data before feeding it to the network to explicitly decorrelate its features. This will ensure that the first layer of the network sees data that follows a nice distribution. However, even if we preprocess the input

data, the activations at deeper layers of the network will likely no longer be decorrelated and will no longer have zero mean or unit variance, since they are output from earlier layers in the network. Even worse, during the training process the distribution of features at each layer of the network will shift as the weights of each layer are updated.

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

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

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

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

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

%load_ext autoreload
%autoreload 2

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

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

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

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

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

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

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

2 Batch Normalization: Forward Pass

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

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

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

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

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

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

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

gamma = np.asarray([1.0, 2.0, 3.0])
```

```

beta = np.asarray([11.0, 12.0, 13.0])

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

```

Before batch normalization:

```

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

```

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

```

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

```

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

```

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

```

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

```

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

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

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

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

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

```

```
print_mean_std(a_norm,axis=0)
```

After batch normalization (test-time):

```
means: [-0.03927354 -0.04349152 -0.10452688]
```

```
stds: [1.01531428 1.01238373 0.97819988]
```

3 Batch Normalization: Backward Pass

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

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

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

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

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

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

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

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

```
dx error: 1.702926968594948e-09
```

```
dgamma error: 7.420414216247087e-13
```

```
dbeta error: 2.8795057655839487e-12
```


4 Batch Normalization: Alternative Backward Pass

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

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

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

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

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

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

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

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

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

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

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

```
[7]: np.random.seed(231)
N, D = 100, 500
x = 5 * np.random.randn(N, D) + 12
gamma = np.random.randn(D)
```

```

beta = np.random.randn(D)
dout = np.random.randn(N, D)

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

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

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

```

```

dx difference:  9.989010819624349e-13
dgamma difference:  0.0
dbeta difference:  0.0
speedup: 1.79x

```

5 Fully Connected Networks with Batch Normalization

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

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

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

```

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

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

```

```

reg=reg, weight_scale=5e-2, dtype=np.float64,
normalization='batchnorm')

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

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

```

```

Running check with reg = 0
Initial loss: 2.2611955101340957
W1 relative error: 1.10e-04
W2 relative error: 5.65e-06
W3 relative error: 4.23e-10
b1 relative error: 4.44e-08
b2 relative error: 5.55e-09
b3 relative error: 1.77e-10
beta1 relative error: 1.77e-08
beta2 relative error: 1.71e-09
gamma1 relative error: 6.96e-09
gamma2 relative error: 3.35e-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: 6.94e-10
b2 relative error: 4.44e-03
b3 relative error: 3.81e-10
beta1 relative error: 6.32e-09
beta2 relative error: 4.23e-09
gamma1 relative error: 6.27e-09
gamma2 relative error: 5.18e-09

```

6 Batch Normalization for Deep Networks

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

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

```

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

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

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

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

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

```

Solver with batch norm:

```

(Iteration 1 / 200) loss: 2.340975
(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: 2.039345
(Epoch 2 / 10) train acc: 0.394000; val_acc: 0.286000
(Iteration 41 / 200) loss: 2.045770
(Epoch 3 / 10) train acc: 0.479000; val_acc: 0.321000
(Iteration 61 / 200) loss: 1.772057
(Epoch 4 / 10) train acc: 0.528000; val_acc: 0.317000

```

```

(Iteration 81 / 200) loss: 1.238056
(Epoch 5 / 10) train acc: 0.619000; val_acc: 0.330000
(Iteration 101 / 200) loss: 1.366125
(Epoch 6 / 10) train acc: 0.650000; val_acc: 0.312000
(Iteration 121 / 200) loss: 1.129776
(Epoch 7 / 10) train acc: 0.671000; val_acc: 0.304000
(Iteration 141 / 200) loss: 1.215410
(Epoch 8 / 10) train acc: 0.730000; val_acc: 0.320000
(Iteration 161 / 200) loss: 0.808850
(Epoch 9 / 10) train acc: 0.778000; val_acc: 0.334000
(Iteration 181 / 200) loss: 0.992587
(Epoch 10 / 10) train acc: 0.779000; val_acc: 0.310000

```

Solver without batch norm:

```

(Iteration 1 / 200) loss: 2.302332
(Epoch 0 / 10) train acc: 0.129000; val_acc: 0.131000
(Epoch 1 / 10) train acc: 0.283000; val_acc: 0.250000
(Iteration 21 / 200) loss: 2.041970
(Epoch 2 / 10) train acc: 0.316000; val_acc: 0.277000
(Iteration 41 / 200) loss: 1.900473
(Epoch 3 / 10) train acc: 0.373000; val_acc: 0.282000
(Iteration 61 / 200) loss: 1.713156
(Epoch 4 / 10) train acc: 0.390000; val_acc: 0.310000
(Iteration 81 / 200) loss: 1.662209
(Epoch 5 / 10) train acc: 0.434000; val_acc: 0.300000
(Iteration 101 / 200) loss: 1.696059
(Epoch 6 / 10) train acc: 0.535000; val_acc: 0.345000
(Iteration 121 / 200) loss: 1.557987
(Epoch 7 / 10) train acc: 0.530000; val_acc: 0.304000
(Iteration 141 / 200) loss: 1.432189
(Epoch 8 / 10) train acc: 0.628000; val_acc: 0.339000
(Iteration 161 / 200) loss: 1.033931
(Epoch 9 / 10) train acc: 0.661000; val_acc: 0.340000
(Iteration 181 / 200) loss: 0.901034
(Epoch 10 / 10) train acc: 0.726000; val_acc: 0.318000

```

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

```

[10]: def plot_training_history(title, label, baseline, bn_solvers, plot_fn,
    ↪bl_marker='.', bn_marker='.', labels=None):
    """utility function for plotting training history"""
    plt.title(title)
    plt.xlabel(label)
    bn_plots = [plot_fn(bn_solver) for bn_solver in bn_solvers]
    bl_plot = plot_fn(baseline)
    num_bn = len(bn_plots)
    for i in range(num_bn):

```

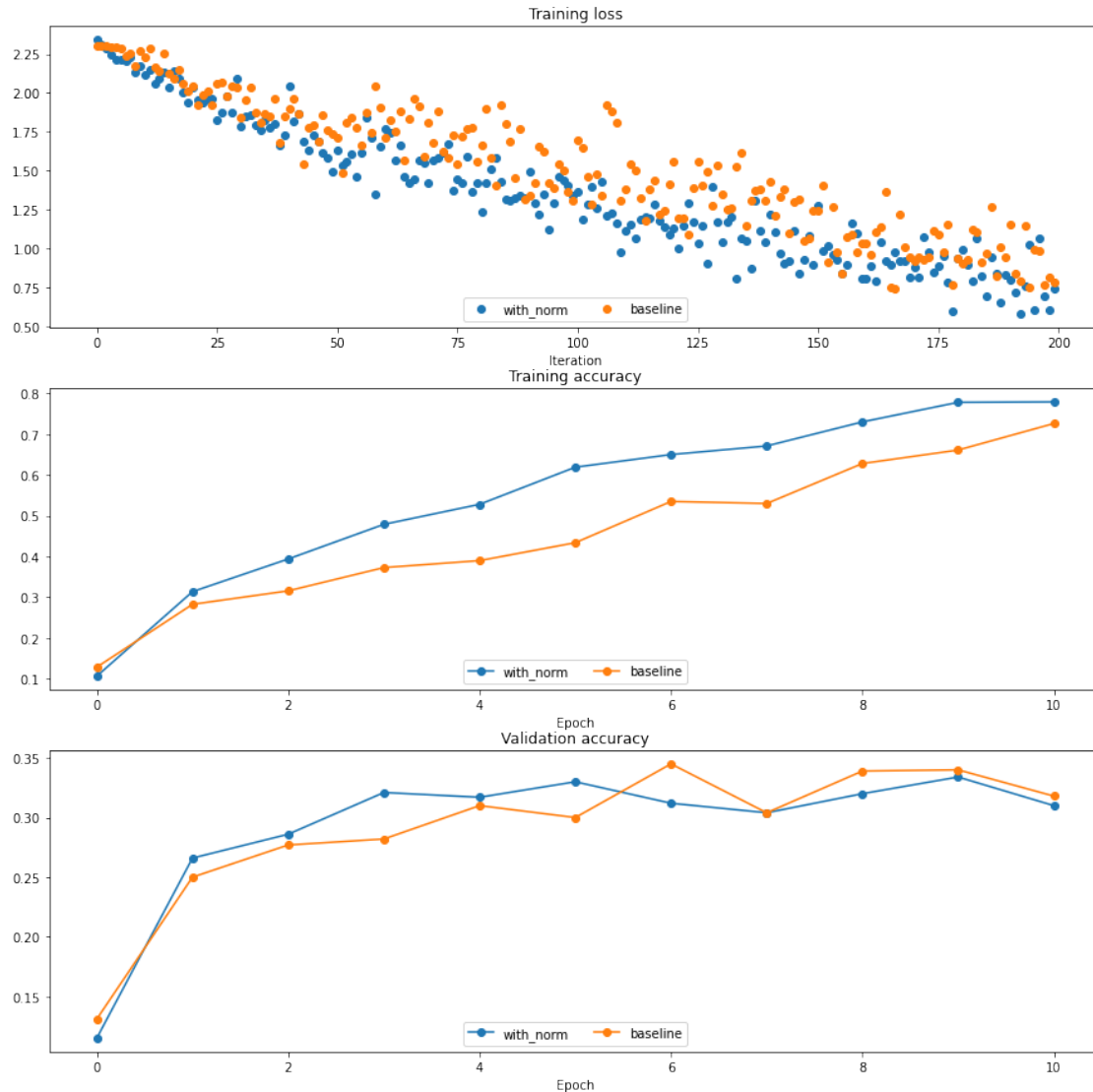
```

        label='with_norm'
        if labels is not None:
            label += str(labels[i])
        plt.plot(bn_plots[i], bn_marker, label=label)
    label='baseline'
    if labels is not None:
        label += str(labels[0])
    plt.plot(bl_plot, bl_marker, label=label)
    plt.legend(loc='lower center', ncol=num_bn+1)

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

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

```



7 Batch Normalization and Initialization

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

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

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

# Try training a very deep net with batchnorm.
```

```

hidden_dims = [50, 50, 50, 50, 50, 50, 50]
num_train = 1000
small_data = {
    'X_train': data['X_train'][:num_train],
    'y_train': data['y_train'][:num_train],
    'X_val': data['X_val'],
    'y_val': data['y_val'],
}

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

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

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

```

```

Running weight scale 1 / 20
Running weight scale 2 / 20
Running weight scale 3 / 20
Running weight scale 4 / 20
Running weight scale 5 / 20
Running weight scale 6 / 20
Running weight scale 7 / 20
Running weight scale 8 / 20

```



```

Running weight scale 9 / 20
Running weight scale 10 / 20
Running weight scale 11 / 20
Running weight scale 12 / 20
Running weight scale 13 / 20
Running weight scale 14 / 20
Running weight scale 15 / 20
Running weight scale 16 / 20

/content/drive/My Drive/enpm809K/assignments/assignment2/cs231n/layers.py:168:
RuntimeWarning: divide by zero encountered in log
  loss = -np.sum(np.log(sy) / np.sum(P, axis=1))

Running weight scale 17 / 20
Running weight scale 18 / 20
Running weight scale 19 / 20
Running weight scale 20 / 20

```

```

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

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

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

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

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

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

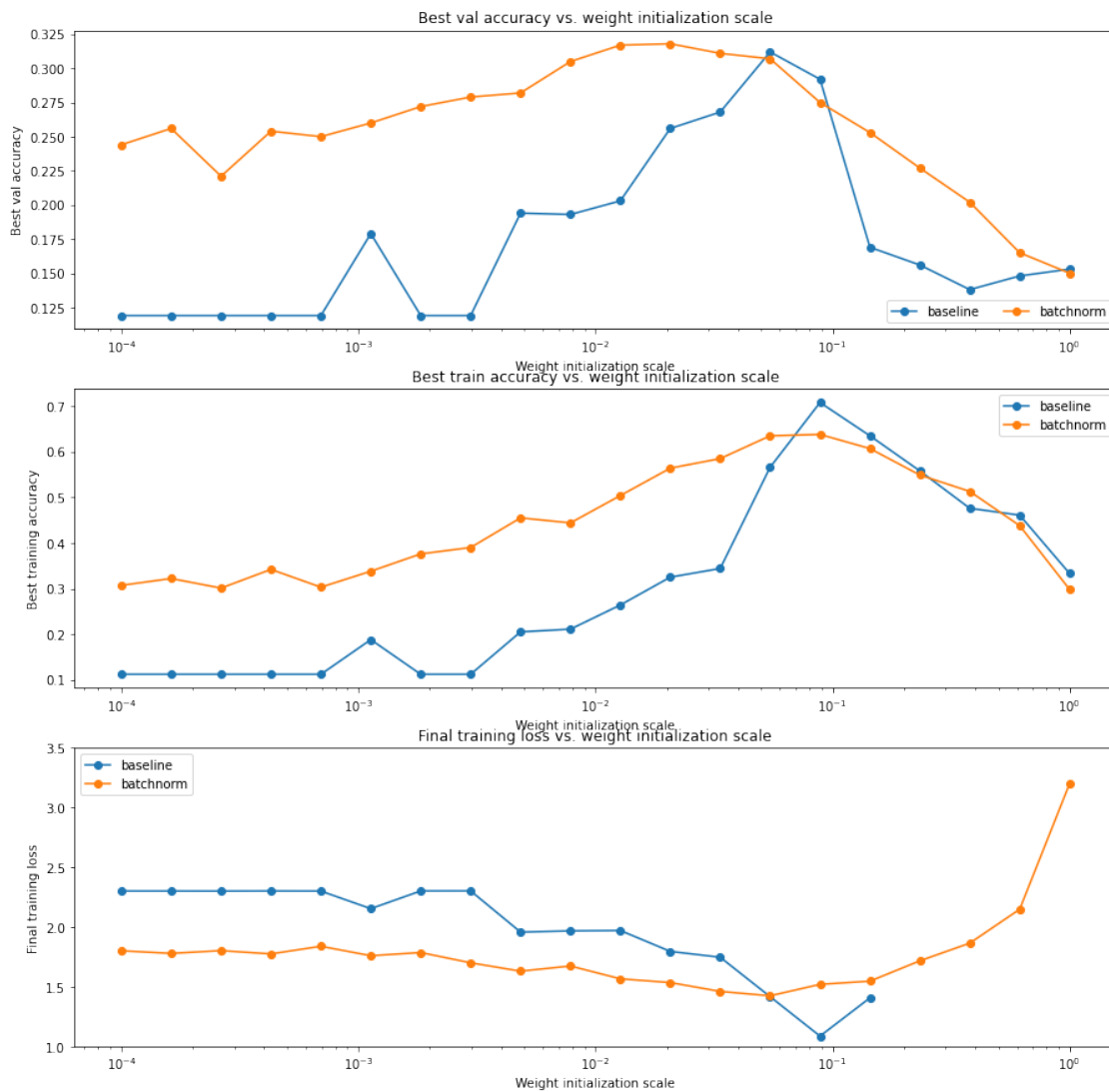
```

```

plt.subplot(3, 1, 3)
plt.title('Final training loss vs. weight initialization scale')
plt.xlabel('Weight initialization scale')
plt.ylabel('Final training loss')
plt.semilogx(weight_scales, final_train_loss, '-o', label='baseline')
plt.semilogx(weight_scales, bn_final_train_loss, '-o', label='batchnorm')
plt.legend()
plt.gca().set_ylim(1.0, 3.5)

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

```



7.1 Inline Question 1:

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

7.2 Answer:

if without normalization, the baseline only works when the weight initialization around the $1e-1$, it seems like very sensitive to the weight initialization, on the other hand when I applied the batchnormalization, the performance is more stable when we change the weight initialization scale.

The reason is because 1. we do not exactly know the distribution of the data, it might be very dense and hard to find a boundary of classification, after normalization will help solve this problem. 2. The centered data can prevent gradient vanish or explode, or all elements in the gradient have the same sign with x .

8 Batch Normalization and Batch Size

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

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

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

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

    print('No normalization: batch size = ',solver_bsize)
    model = FullyConnectedNet(hidden_dims, weight_scale=weight_scale,
    ↪normalization=None)
    solver = Solver(model, small_data,
                    num_epochs=n_epochs, batch_size=solver_bsize,
```

```

        update_rule='adam',
        optim_config={
            'learning_rate': lr,
        },
        verbose=False)
    solver.train()

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

    return bn_solvers, solver, batch_sizes

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

```

```

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

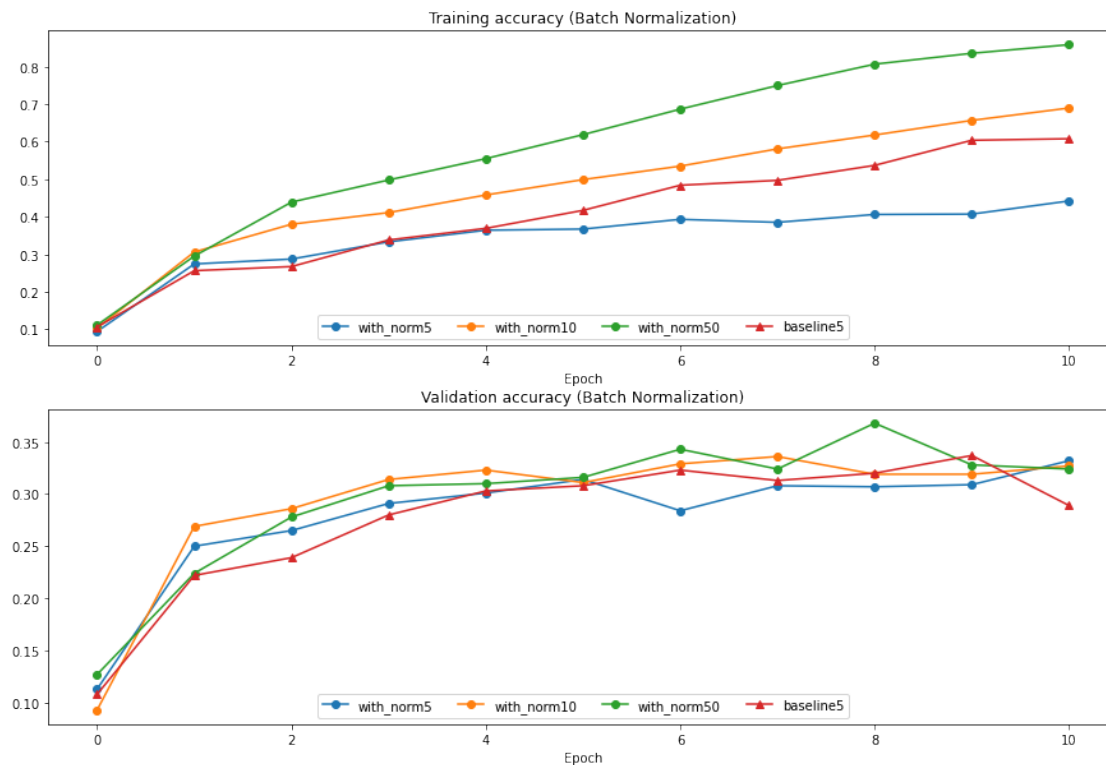
```

```

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

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

```



8.1 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?

8.2 Answer:

When the batch normalization value is larger, the validate accuracy is higher. This might be because if we have a larger data set our training beta and gamma will have a better result for normalizing the data even the testing or validate data.

9 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>)

9.1 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.

9.2 Answer:

Batch normalization:(3) The batch normalization is design to get the average of all feature, (3) find a mean value from every data and applied it to all data, it just like the idea of batch normalization.

Layer normalization:(2) The layer normalization is design to get the average of the input data.

10 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.

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

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

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

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

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

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

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

Before layer normalization:

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

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

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

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

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

```
[16]: # Gradient check batchnorm backward pass.
```

```
np.random.seed(231)
N, D = 4, 5
x = 5 * np.random.randn(N, D) + 12
```

```

gamma = np.random.randn(D)
beta = np.random.randn(D)
dout = np.random.randn(N, D)

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

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

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

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

```

```

dx error: 1.433616168873336e-09
dgamma error: 4.519489546032799e-12
dbeta error: 2.276445013433725e-12

```

11 Layer Normalization and Batch Size

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

```

[17]: ln_solvers_bsize, solver_bsize, batch_sizes = \
    ↪run_batchsize_experiments('layernorm')

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

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

```



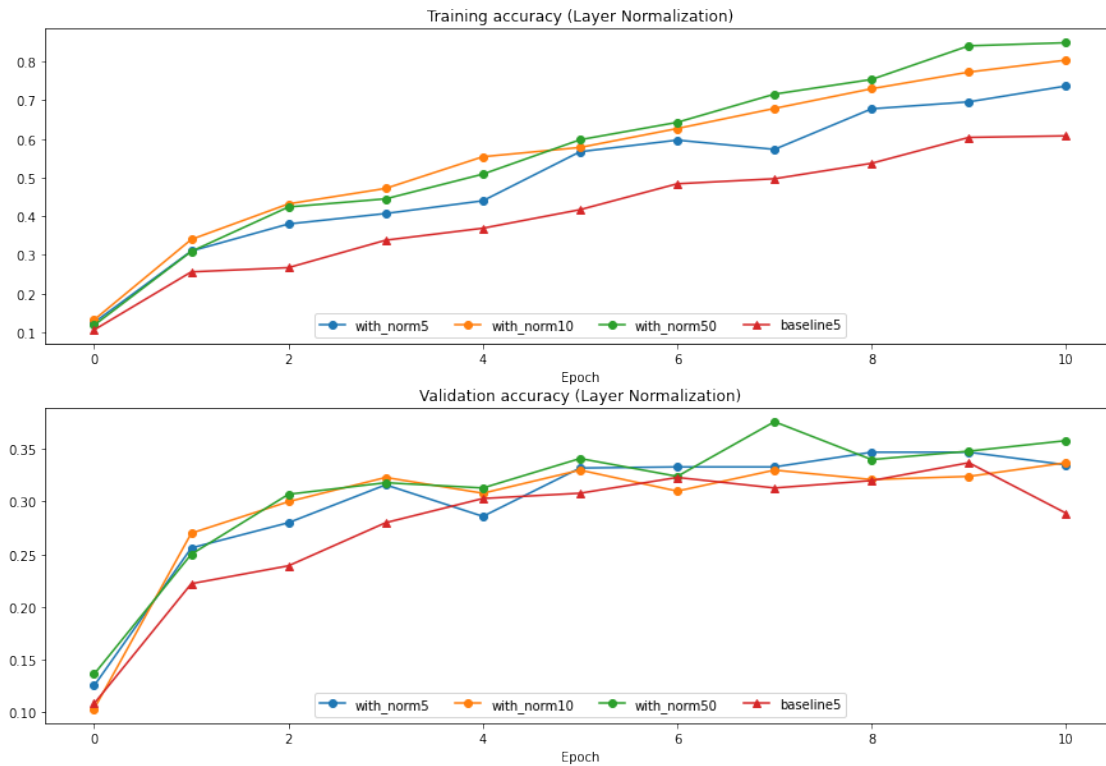
```
plt.show()
```

No normalization: batch size = 5

Normalization: batch size = 5

Normalization: batch size = 10

Normalization: batch size = 50



11.1 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

11.2 Answer:

- (2) It is not representable for the average when your data is too small, for layer normalization, it takes the average over the feature so when the features is in a small dimesions it will not work well.

- (3) When the regularization term is too high it will underfitting so that the layer normalization will not work well. Actually, everything will not work well when the regularization term is too high.

Dropout

November 4, 2022

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

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'cs231n/assignments/assignment1/'
FOLDERNAME = 'enpm809K/assignments/assignment2/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

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

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

Mounted at /content/drive

/content/drive/My Drive/enpm809K/assignments/assignment2/cs231n/datasets

/content/drive/My Drive/enpm809K/assignments/assignment2

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>)

```
[2]: # Setup cell.
import time
import numpy as np
```

```

import matplotlib.pyplot as plt
from cs231n.classifiers.fc_net import *
from cs231n.data_utils import get_CIFAR10_data
from cs231n.gradient_check import eval_numerical_gradient, \
    eval_numerical_gradient_array
from cs231n.solver import Solver

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

%load_ext autoreload
%autoreload 2

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

```

```

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

```

```

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

```

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.

```

[4]: 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:  0.0
Mean of test-time output:  10.000207878477502
Fraction of train-time output set to zero:  1.0
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:  0.0
Mean of test-time output:  10.000207878477502
Fraction of train-time output set to zero:  1.0
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:  0.0
Mean of test-time output:  10.000207878477502
Fraction of train-time output set to zero:  1.0
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.

```

[5]: 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: 0.0

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:

When we apply dropout in our system, we only use the p (dropout ratio) * number of neurons in the training, however, in the test time, we do not apply mask in the inputs therefore every neurons will be used, this caused the number of neurons while training and testing are different, to make it become the same, we need to divide p in the training step.

4 Fully Connected Networks 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_keep_ratio` 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.

```
[6]: 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_keep_ratio in [1, 0.75, 0.5]:
    print('Running check with dropout = ', dropout_keep_ratio)
    model = FullyConnectedNet(
        [H1, H2],
        input_dim=D,
        num_classes=C,
        weight_scale=5e-2,
        dtype=np.float64,
        dropout_keep_ratio=dropout_keep_ratio,
        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_keep_ratio=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: 2.42e-07
W2 relative error: 2.06e-04
W3 relative error: 7.88e-07
b1 relative error: 3.57e-09
b2 relative error: 2.09e-09
b3 relative error: 1.89e-10

```

```

Running check with dropout = 0.75
Initial loss: 2.3016482157750753
W1 relative error: 1.00e+00
W2 relative error: 1.00e+00
W3 relative error: 1.81e-07
b1 relative error: 1.00e+00
b2 relative error: 1.00e+00
b3 relative error: 1.28e-10

```

```

Running check with dropout = 0.5
Initial loss: 2.30485979047391
W1 relative error: 1.00e+00
W2 relative error: 1.00e+00
W3 relative error: 6.89e-08
b1 relative error: 1.00e+00
b2 relative error: 1.00e+00
b3 relative error: 1.41e-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.

```

[7]: # 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_keep_ratio in dropout_choices:
    model = FullyConnectedNet(
        [500],
        dropout_keep_ratio=dropout_keep_ratio
    )
    print(dropout_keep_ratio)

    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_keep_ratio] = solver
    print()

```

1

```

(Iteration 1 / 125) loss: 7.856644
(Epoch 0 / 25) train acc: 0.260000; val_acc: 0.184000
(Epoch 1 / 25) train acc: 0.416000; val_acc: 0.258000
(Epoch 2 / 25) train acc: 0.482000; val_acc: 0.276000
(Epoch 3 / 25) train acc: 0.532000; val_acc: 0.277000
(Epoch 4 / 25) train acc: 0.600000; val_acc: 0.271000
(Epoch 5 / 25) train acc: 0.708000; val_acc: 0.299000
(Epoch 6 / 25) train acc: 0.722000; val_acc: 0.282000
(Epoch 7 / 25) train acc: 0.832000; val_acc: 0.255000
(Epoch 8 / 25) train acc: 0.880000; val_acc: 0.268000
(Epoch 9 / 25) train acc: 0.902000; val_acc: 0.277000
(Epoch 10 / 25) train acc: 0.898000; val_acc: 0.261000
(Epoch 11 / 25) train acc: 0.924000; val_acc: 0.263000
(Epoch 12 / 25) train acc: 0.960000; val_acc: 0.300000
(Epoch 13 / 25) train acc: 0.972000; val_acc: 0.314000
(Epoch 14 / 25) train acc: 0.972000; val_acc: 0.310000
(Epoch 15 / 25) train acc: 0.974000; val_acc: 0.314000
(Epoch 16 / 25) train acc: 0.994000; val_acc: 0.304000

```



```

(Epoch 17 / 25) train acc: 0.970000; val_acc: 0.307000
(Epoch 18 / 25) train acc: 0.992000; val_acc: 0.311000
(Epoch 19 / 25) train acc: 0.992000; val_acc: 0.311000
(Epoch 20 / 25) train acc: 0.990000; val_acc: 0.288000
(Iteration 101 / 125) loss: 0.001514
(Epoch 21 / 25) train acc: 0.996000; val_acc: 0.290000
(Epoch 22 / 25) train acc: 0.998000; val_acc: 0.306000
(Epoch 23 / 25) train acc: 0.996000; val_acc: 0.307000
(Epoch 24 / 25) train acc: 0.998000; val_acc: 0.309000
(Epoch 25 / 25) train acc: 0.998000; val_acc: 0.303000

```

0.25

```

(Iteration 1 / 125) loss: 10.430470
(Epoch 0 / 25) train acc: 0.116000; val_acc: 0.117000
(Epoch 1 / 25) train acc: 0.162000; val_acc: 0.155000
(Epoch 2 / 25) train acc: 0.254000; val_acc: 0.201000
(Epoch 3 / 25) train acc: 0.312000; val_acc: 0.207000
(Epoch 4 / 25) train acc: 0.376000; val_acc: 0.226000
(Epoch 5 / 25) train acc: 0.414000; val_acc: 0.235000
(Epoch 6 / 25) train acc: 0.494000; val_acc: 0.228000
(Epoch 7 / 25) train acc: 0.524000; val_acc: 0.239000
(Epoch 8 / 25) train acc: 0.606000; val_acc: 0.246000
(Epoch 9 / 25) train acc: 0.650000; val_acc: 0.252000
(Epoch 10 / 25) train acc: 0.664000; val_acc: 0.241000
(Epoch 11 / 25) train acc: 0.702000; val_acc: 0.245000
(Epoch 12 / 25) train acc: 0.738000; val_acc: 0.255000
(Epoch 13 / 25) train acc: 0.754000; val_acc: 0.249000
(Epoch 14 / 25) train acc: 0.788000; val_acc: 0.245000
(Epoch 15 / 25) train acc: 0.828000; val_acc: 0.254000
(Epoch 16 / 25) train acc: 0.842000; val_acc: 0.246000
(Epoch 17 / 25) train acc: 0.852000; val_acc: 0.254000
(Epoch 18 / 25) train acc: 0.876000; val_acc: 0.257000
(Epoch 19 / 25) train acc: 0.914000; val_acc: 0.252000
(Epoch 20 / 25) train acc: 0.930000; val_acc: 0.252000
(Iteration 101 / 125) loss: 0.269993
(Epoch 21 / 25) train acc: 0.946000; val_acc: 0.259000
(Epoch 22 / 25) train acc: 0.952000; val_acc: 0.261000
(Epoch 23 / 25) train acc: 0.960000; val_acc: 0.266000
(Epoch 24 / 25) train acc: 0.970000; val_acc: 0.262000
(Epoch 25 / 25) train acc: 0.978000; val_acc: 0.279000

```

```

[8]: # Plot train and validation accuracies of the two models.
train_accs = []
val_accs = []
for dropout_keep_ratio in dropout_choices:
    solver = solvers[dropout_keep_ratio]

```

```

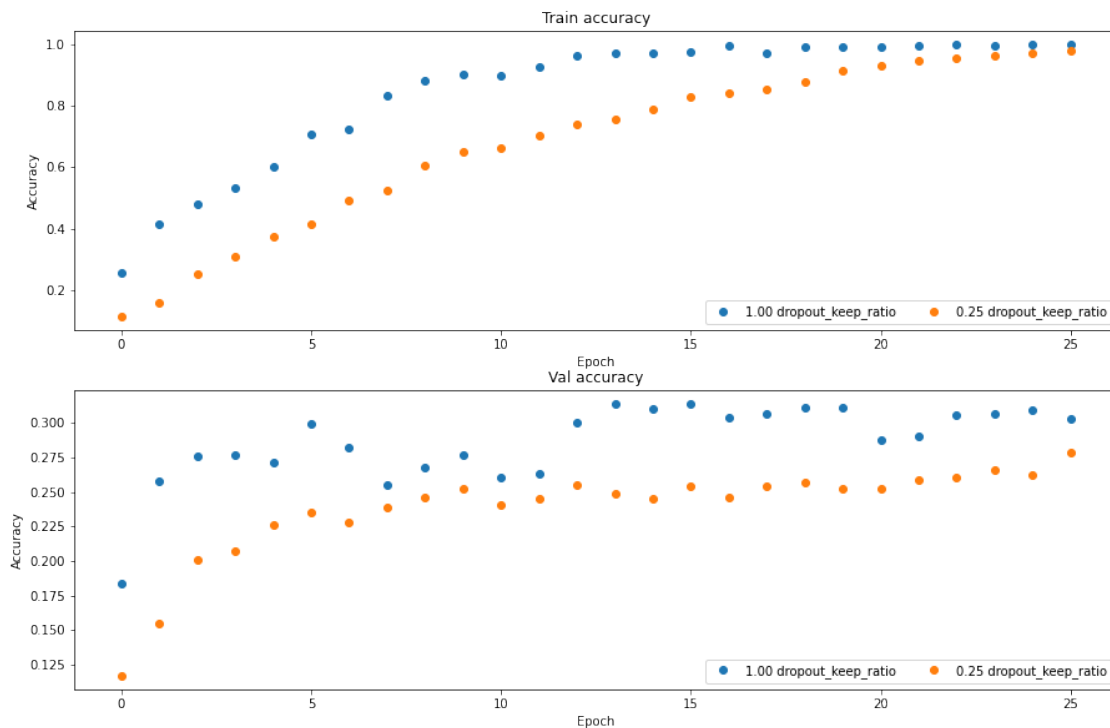
train_accs.append(solver.train_acc_history[-1])
val_accs.append(solver.val_acc_history[-1])

plt.subplot(3, 1, 1)
for dropout_keep_ratio in dropout_choices:
    plt.plot(
        solvers[dropout_keep_ratio].train_acc_history, 'o', label='%0.2f_
        ↳dropout_keep_ratio' % dropout_keep_ratio)
plt.title('Train accuracy')
plt.xlabel('Epoch')
plt.ylabel('Accuracy')
plt.legend(ncol=2, loc='lower right')

plt.subplot(3, 1, 2)
for dropout_keep_ratio in dropout_choices:
    plt.plot(
        solvers[dropout_keep_ratio].val_acc_history, 'o', label='%0.2f_
        ↳dropout_keep_ratio' % dropout_keep_ratio)
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:

In the result, we can see that when we do not apply dropout in our architecture, the training set become 100% accuracy but the validation set is not that accurate, which means overfitting. After adding dropout, we can prevent overfitting, just like the regularize does.

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:

Since you have already decrease the size of the hidden layers, it have already have some sense of preventing overfitting. I would recommend to keep the drop ratio or increase, it will better for preventing overfitting.

[8]:

Convolutional Networks

November 4, 2022

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

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'cs231n/assignments/assignment1/'
FOLDERNAME = 'enpm809K/assignments/assignment2/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

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

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

Drive already mounted at /content/drive; to attempt to forcibly remount, call `drive.mount("/content/drive", force_remount=True)`.

```
/content/drive/My Drive/enpm809K/assignments/assignment2/cs231n/datasets
/content/drive/My Drive/enpm809K/assignments/assignment2
```

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.

```
[11]: # Setup cell.
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))))
```

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

```
%reload_ext autoreload
```

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

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

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:

```
[13]: 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
```

2.1 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.

```
[14]: 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')

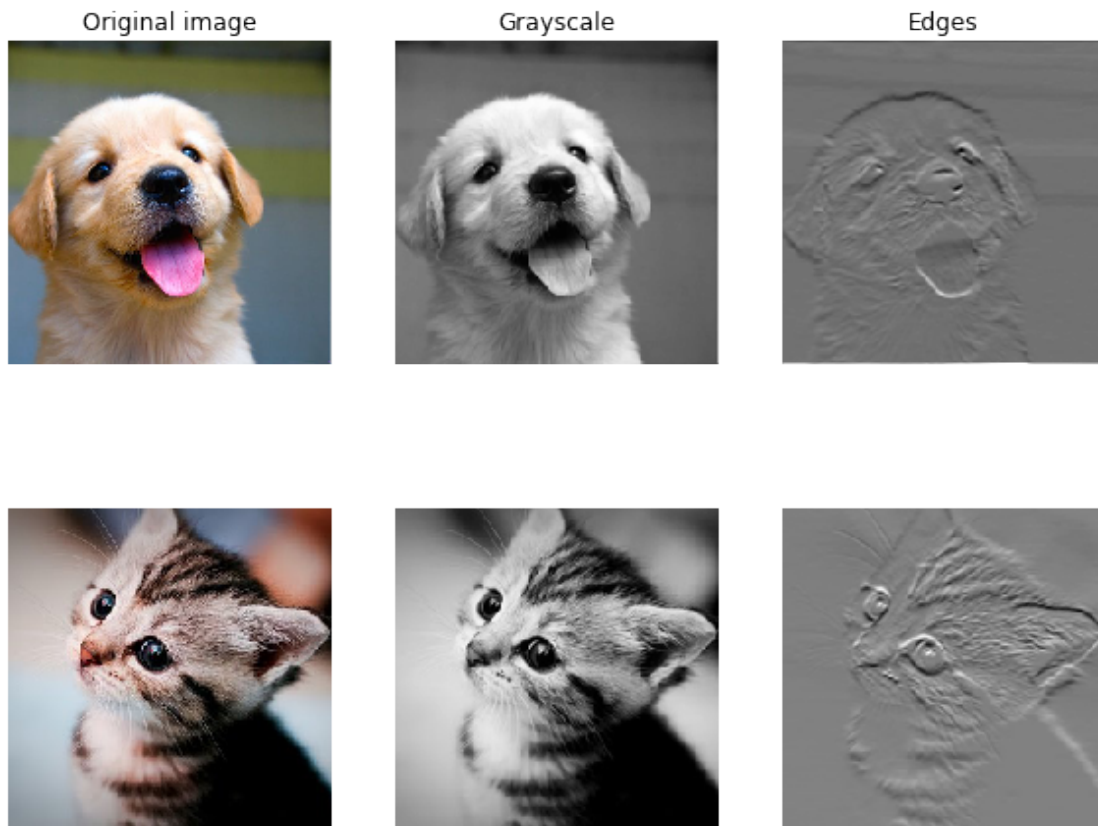
# Show the original images and the results of the conv operation
plt.subplot(2, 3, 1)
imshow_no_ax(puppy, normalize=False)
plt.title('Original image')
plt.subplot(2, 3, 2)
imshow_no_ax(out[0, 0])
plt.title('Grayscale')
plt.subplot(2, 3, 3)

```

```

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()

```



3 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.


```
[15]: 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
```

4 Max-Pooling: Naive Forward Pass

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:

```
[16]: 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

5 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:

```

[17]: 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)
# print(x)
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

6 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`.

6.0.1 Execute the below cell, save the notebook, and restart the runtime

The fast convolution implementation depends on a Cython extension; to compile it, run the cell below. Next, save the Colab notebook (File > Save) and **restart the runtime** (Runtime > Restart runtime). You can then re-execute the preceeding cells from top to bottom and skip the cell below as you only need to run it once for the compilation step.

```
[18]: # Remember to restart the runtime after executing this cell!
%cd /content/drive/My\ Drive/$FOLDERNAME/cs231n/
!python setup.py build_ext --inplace
%cd /content/drive/My\ Drive/$FOLDERNAME/

/content/drive/My Drive/enpm809K/assignments/assignment2/cs231n
Compiling im2col_cython.pyx because it depends on /usr/local/lib/python3.7/dist-
packages/Cython/Includes/libc/string.pxd.
[1/1] Cythonizing im2col_cython.pyx
/usr/local/lib/python3.7/dist-packages/Cython/Compiler/Main.py:369:
FutureWarning: Cython directive 'language_level' not set, using 2 for now (Py2).
This will change in a later release! File: /content/drive/My
Drive/enpm809K/assignments/assignment2/cs231n/im2col_cython.pyx
    tree = Parsing.p_module(s, pxd, full_module_name)
running build_ext
building 'im2col_cython' extension
x86_64-linux-gnu-gcc -pthread -Wno-unused-result -Wsign-compare -DNDEBUG -g
-fwrapv -O2 -Wall -g -fstack-protector-strong -Wformat -Werror=format-security
-g -fwrapv -O2 -g -fstack-protector-strong -Wformat -Werror=format-security
-Wdate-time -D_FORTIFY_SOURCE=2 -fPIC -I/usr/local/lib/python3.7/dist-
packages/numpy/core/include -I/usr/include/python3.7m -c im2col_cython.c -o
build/temp.linux-x86_64-3.7/im2col_cython.o
In file included from /usr/local/lib/python3.7/dist-
packages/numpy/core/include/numpy/ndarraytypes.h:1969:0,
         from /usr/local/lib/python3.7/dist-
packages/numpy/core/include/numpy/ndarrayobject.h:12,
         from /usr/local/lib/python3.7/dist-
packages/numpy/core/include/numpy/arrayobject.h:4,
         from im2col_cython.c:768:
/usr/local/lib/python3.7/dist-
packages/numpy/core/include/numpy/npymath/npymath.h:17:2:
warning: #warning "Using deprecated NumPy API, disable it with
" "#define NPY_NO_DEPRECATED_API NPY_1_7_API_VERSION" [-Wcpp]
#warning "Using deprecated NumPy API, disable it with " \
~~~~~
x86_64-linux-gnu-gcc -pthread -shared -Wl,-O1 -Wl,-Bsymbolic-functions
```

```
-Wl,-Bsymbolic-functions -g -fwrapv -O2 -Wl,-Bsymbolic-functions -g -fwrapv -O2
-g -fstack-protector-strong -Wformat -Werror=format-security -Wdate-time
-D_FORTIFY_SOURCE=2 build/temp.linux-x86_64-3.7/im2col_cython.o -o
/content/drive/My Drive/enpm809K/assignments/assignment2/cs231n/im2col_cython.cp
ython-37m-x86_64-linux-gnu.so
/content/drive/My Drive/enpm809K/assignments/assignment2
```

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:

```
[19]: # 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: 5.555731s

Fast: 0.013207s

Speedup: 420.659374x

Difference: 4.926407851494105e-11

Testing conv_backward_fast:

Naive: 8.114237s

Fast: 0.013322s

Speedup: 609.080266x

dx difference: 1.949764775345631e-11

dw difference: 3.681156828004736e-13

db difference: 3.481354613192702e-14

```

[20]: # 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: 0.435685s
fast: 0.008776s
speedup: 49.642634x
difference: 0.0

```

Testing pool_backward_fast:

```

Naive: 0.483390s
fast: 0.014009s
speedup: 34.505080x
dx difference: 0.0

```

7 Convolutional "Sandwich" Layers

In the previous assignment, 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 their usage.

```

[21]: 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

```
[22]: 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

8 Three-Layer Convolutional Network

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:

8.1 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.

```
[23]: 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): 2.508255638232932
```

8.2 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 $e-2$.

```
[24]: 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: 6.284982e-04
W2 max relative error: 5.934202e-02
W3 max relative error: 2.375651e-04
b1 max relative error: 3.397321e-06
b2 max relative error: 2.517459e-03
b3 max relative error: 1.422510e-09

```

8.3 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.

```

[25]: 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

```

```

[26]: # 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

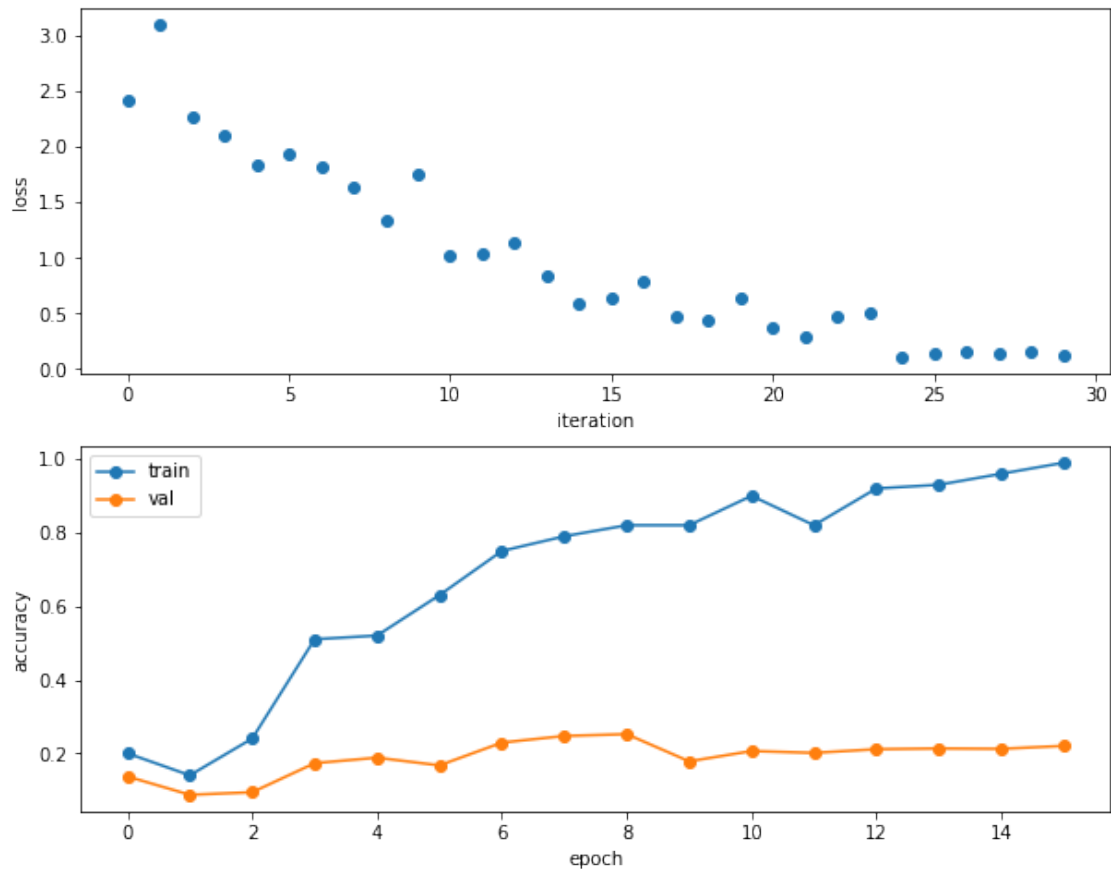
```
[27]: # Print final validation accuracy.
print(
    "Small data validation accuracy:",
    solver.check_accuracy(small_data['X_val'], small_data['y_val'])
)
```

Small data validation accuracy: 0.252

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

```
[28]: 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()
```



8.4 Train the Network

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

```
[29]: 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: 2.304740
(Epoch 0 / 1) train acc: 0.103000; val_acc: 0.107000
(Iteration 21 / 980) loss: 2.098229
(Iteration 41 / 980) loss: 1.949788
(Iteration 61 / 980) loss: 1.888398
(Iteration 81 / 980) loss: 1.877093
(Iteration 101 / 980) loss: 1.851877
(Iteration 121 / 980) loss: 1.859353
(Iteration 141 / 980) loss: 1.800181
(Iteration 161 / 980) loss: 2.143292
(Iteration 181 / 980) loss: 1.830573
(Iteration 201 / 980) loss: 2.037280
(Iteration 221 / 980) loss: 2.020304
(Iteration 241 / 980) loss: 1.823728
(Iteration 261 / 980) loss: 1.692679
(Iteration 281 / 980) loss: 1.882594
(Iteration 301 / 980) loss: 1.798261
(Iteration 321 / 980) loss: 1.851960
(Iteration 341 / 980) loss: 1.716323
(Iteration 361 / 980) loss: 1.897655
(Iteration 381 / 980) loss: 1.319744
(Iteration 401 / 980) loss: 1.738790
(Iteration 421 / 980) loss: 1.488866
(Iteration 441 / 980) loss: 1.718409
(Iteration 461 / 980) loss: 1.744440
(Iteration 481 / 980) loss: 1.605460
(Iteration 501 / 980) loss: 1.494847
(Iteration 521 / 980) loss: 1.835179
(Iteration 541 / 980) loss: 1.483923
(Iteration 561 / 980) loss: 1.676871
(Iteration 581 / 980) loss: 1.438325
(Iteration 601 / 980) loss: 1.443469
(Iteration 621 / 980) loss: 1.529369
(Iteration 641 / 980) loss: 1.763475
(Iteration 661 / 980) loss: 1.790329
(Iteration 681 / 980) loss: 1.693343
(Iteration 701 / 980) loss: 1.637078
(Iteration 721 / 980) loss: 1.644564
(Iteration 741 / 980) loss: 1.708919
(Iteration 761 / 980) loss: 1.494252
(Iteration 781 / 980) loss: 1.901751
(Iteration 801 / 980) loss: 1.898991
(Iteration 821 / 980) loss: 1.489988
(Iteration 841 / 980) loss: 1.377615
(Iteration 861 / 980) loss: 1.763751
(Iteration 881 / 980) loss: 1.540284
(Iteration 901 / 980) loss: 1.525582
(Iteration 921 / 980) loss: 1.674166

```
(Iteration 941 / 980) loss: 1.714316
(Iteration 961 / 980) loss: 1.534668
(Epoch 1 / 1) train acc: 0.504000; val_acc: 0.499000
```

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

Full data training accuracy: 0.4761836734693878

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

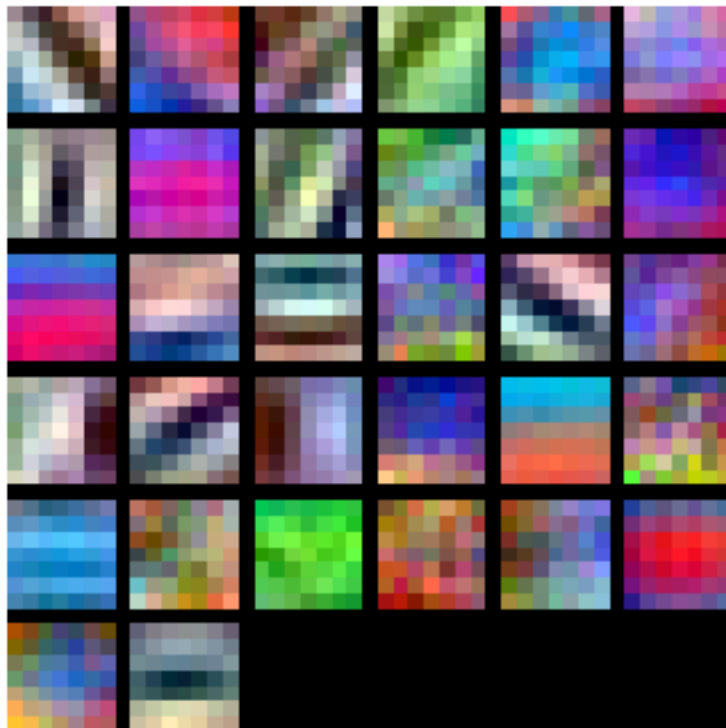
Full data validation accuracy: 0.499

8.5 Visualize Filters

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

```
[32]: 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()
```



9 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 Spatial Batch Normalization: Forward Pass

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:

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

# Check the training-time forward pass by checking means and variances
# of features both before and after spatial batch normalization.
N, C, H, W = 2, 3, 4, 5
x = 4 * np.random.randn(N, C, H, W) + 10

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

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

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

Before spatial batch normalization:

```
shape: (2, 3, 4, 5)
means: [9.33463814 8.90909116 9.11056338]
stds:  [3.61447857 3.19347686 3.5168142 ]
```

After spatial batch normalization:

```
shape: (2, 3, 4, 5)
means: [ 6.18949336e-16  5.99520433e-16 -1.22124533e-16]
stds:  [0.99999962 0.99999951 0.9999996 ]
```

After spatial batch normalization (nontrivial gamma, beta):

```
shape: (2, 3, 4, 5)
means: [6. 7. 8.]
stds:  [2.99999885 3.99999804 4.99999798]
```



```
[34]: 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]
```

11 Spatial Batch Normalization: Backward Pass

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:

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

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

```

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

#You should expect errors of magnitudes between 1e-12~1e-06
_, cache = spatial_batchnorm_forward(x, gamma, beta, bn_param)
dx, dgamma, dbeta = spatial_batchnorm_backward(dout, cache)
print('dx error: ', rel_error(dx_num, dx))
print('dgamma error: ', rel_error(da_num, dgamma))
print('dbeta error: ', rel_error(db_num, dbeta))

```

```

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

```

12 Spatial 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 handcrafted 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.

[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/>)

13 Spatial Group Normalization: Forward Pass

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:

```
[36]: 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  2.65528340e-16 -3.38618023e-16]
stds:  [0.99999963 0.99999948 0.99999973 0.99999968]
```

14 Spatial Group Normalization: Backward Pass

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:

```
[37]: 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 and 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: 7.413109384854475e-08
dgamma error: 9.468195772749234e-12
dbeta error: 3.354494437653335e-12
```

PyTorch

November 4, 2022

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

# TODO: Enter the foldername in your Drive where you have saved the unzipped
# assignment folder, e.g. 'cs231n/assignments/assignment1/'
FOLDERNAME = 'enpm809K/assignments/assignment2/'
assert FOLDERNAME is not None, "[!] Enter the foldername."

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

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

Mounted at /content/drive

/content/drive/My Drive/enpm809K/assignments/assignment2/cs231n/datasets

/content/drive/My Drive/enpm809K/assignments/assignment2

1 Introduction to PyTorch

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

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

1.1 Why do we use deep learning frameworks?

- Our code will now run on GPUs! This will allow our models to train much faster. 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).
- In 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 :)
- Finally, we want you to be exposed to the sort of deep learning code you might run into in academia or industry.

1.2 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.

1.3 How do I learn PyTorch?

One of our former instructors, Justin Johnson, 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.

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 GPU

You can manually switch to a GPU device on Colab by clicking **Runtime** -> **Change runtime type** and selecting **GPU** under **Hardware Accelerator**. You should do this before running the following cells to import packages, since the kernel gets restarted upon switching runtimes.

```
[2]: import torch
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

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: cpu

4 Part I. Preparation

Now, let's 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.

```
[3]: NUM_TRAIN = 49000

# The torchvision.transforms package provides tools for preprocessing data
# and for performing data augmentation; here we set up a transform to
# preprocess the data by subtracting the mean RGB value and dividing by the
# standard deviation of each RGB value; we've hardcoded the mean and std.
transform = T.Compose([
    T.ToTensor(),
    T.Normalize((0.4914, 0.4822, 0.4465), (0.2023, 0.1994, 0.2010))
])

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

5 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.

5.0.1 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 x C x H x 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 x H x 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).

```
[4]: 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]])

```

5.0.2 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.

```

[5]: 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, d1, ..., dM) giving a minibatch of
        input data.
    - params: A list [w1, w2] of PyTorch Tensors giving weights for the network;
        w1 has shape (D, H) and w2 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])
```

5.0.3 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!

```

[6]: 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)*****
# A convolutional layer (with bias) with channel_1 filters, each with shape
↳KW1 x KH1, and zero-padding of two
# ReLU nonlinearity
# A convolutional layer (with bias) with channel_2 filters, each with shape
↳KW2 x KH2, and zero-padding of one
# ReLU nonlinearity
# Fully-connected layer with bias, producing scores for C classes.
x = F.relu(F.conv2d(x, conv_w1, conv_b1, padding=2))
x = F.relu(F.conv2d(x, conv_w2, conv_b2, padding=1))
x2 = flatten(x)
scores = x2.mm(fc_w) + fc_b
```

```

    # pass

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

```

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

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

```

[7]: 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])
```

5.0.4 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>

```
[8]: 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))
```

```
[8]: tensor([[ -0.2564, -1.2601, -0.1910,  0.7398,  0.4345],
          [-0.5676,  0.6106,  0.7797,  1.1722, -0.1383],
          [ 1.0973, -0.4011, -0.1125, -0.8563, -0.5883]], requires_grad=True)
```

5.0.5 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.

```
[9]: 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))

```

5.0.6 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.

```

[10]: 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()

```

5.0.7 BareBones PyTorch: Train a Two-Layer Network

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

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

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

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

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

```

[11]: 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 = 3.6974
Checking accuracy on the val set
Got 172 / 1000 correct (17.20%)
```

```
Iteration 100, loss = 2.2843
Checking accuracy on the val set
Got 344 / 1000 correct (34.40%)
```

```
Iteration 200, loss = 2.0960
Checking accuracy on the val set
Got 391 / 1000 correct (39.10%)
```

```
Iteration 300, loss = 2.1643
Checking accuracy on the val set
Got 388 / 1000 correct (38.80%)
```

```
Iteration 400, loss = 2.0174
Checking accuracy on the val set
Got 405 / 1000 correct (40.50%)
```

```
Iteration 500, loss = 1.5174
Checking accuracy on the val set
Got 399 / 1000 correct (39.90%)
```

```
Iteration 600, loss = 1.4099
Checking accuracy on the val set
Got 442 / 1000 correct (44.20%)
```

```
Iteration 700, loss = 1.4928
Checking accuracy on the val set
Got 442 / 1000 correct (44.20%)
```

5.0.8 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.

```
[12]: 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)*****

# Convolutional layer (with bias) with 32 5x5 filters, with zero-padding of 2
conv_w1 = random_weight((channel_1, 3, 5, 5))
conv_b1 = random_weight((channel_1, ))

# Convolutional layer (with bias) with 16 3x3 filters, with zero-padding of 1
conv_w2 = random_weight((channel_2, channel_1, 3, 3))
conv_b2 = random_weight((channel_2, ))

# Fully-connected layer (with bias) to compute scores for 10 classes
fc_w = random_weight((channel_2 * 1024, 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 = 3.9360

Checking accuracy on the val set

Got 116 / 1000 correct (11.60%)

```
Iteration 100, loss = 1.8167
Checking accuracy on the val set
Got 342 / 1000 correct (34.20%)
```

```
Iteration 200, loss = 1.8018
Checking accuracy on the val set
Got 426 / 1000 correct (42.60%)
```

```
Iteration 300, loss = 1.6762
Checking accuracy on the val set
Got 411 / 1000 correct (41.10%)
```

```
Iteration 400, loss = 1.5482
Checking accuracy on the val set
Got 422 / 1000 correct (42.20%)
```

```
Iteration 500, loss = 1.5530
Checking accuracy on the val set
Got 439 / 1000 correct (43.90%)
```

```
Iteration 600, loss = 1.6782
Checking accuracy on the val set
Got 460 / 1000 correct (46.00%)
```

```
Iteration 700, loss = 1.4434
Checking accuracy on the val set
Got 469 / 1000 correct (46.90%)
```

6 Part III. PyTorch Module API

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

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

To use the Module API, follow the steps below:

1. Subclass `nn.Module`. Give your network class an intuitive name like `TwoLayerFC`.
2. In the constructor `__init__()`, define all the layers you need as class attributes. Layer objects like `nn.Linear` and `nn.Conv2d` are themselves `nn.Module` subclasses and contain learnable

parameters, so that you don't have to instantiate the raw tensors yourself. `nn.Module` will track these internal parameters for you. Refer to the [doc](#) to learn more about the dozens of builtin layers. **Warning:** don't forget to call the `super().__init__()` first!

3. In the `forward()` method, define the *connectivity* of your network. You should use the attributes defined in `__init__` as function calls that take tensor as input and output the "transformed" tensor. Do *not* create any new layers with learnable parameters in `forward()`! All of them must be declared upfront in `__init__`.

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

6.0.1 Module API: Two-Layer Network

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

```
[13]: 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])
```

6.0.2 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.

```
[14]: 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)*****
        # Convolutional layer with channel_1 5x5 filters with zero-padding of 2
        self.conv1 = nn.Conv2d(in_channel, channel_1, 5, padding=2)
        nn.init.kaiming_normal_(self.conv1.weight)
        # Convolutional layer with channel_1 5x5 filters with zero-padding of 2
        self.conv2 = nn.Conv2d(channel_1, channel_2, 3, padding=1)
        nn.init.kaiming_normal_(self.conv2.weight)
        # Fully-connected layer to num_classes classes
        self.fc3 = nn.Linear(channel_2 * 1024, num_classes)
        nn.init.kaiming_normal_(self.fc3.weight)
        # pass

        # *****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)*****
        x1 = F.relu(self.conv1(x))
        x2 = F.relu(self.conv2(x1))
```

```

        x3 = flatten(x2)
        scores = self.fc3(x3) # + fc_b
        # pass

        # *****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])
```

6.0.3 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.

```

[15]: 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))
```

6.0.4 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.

```
[16]: 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()
```

6.0.5 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.

```
[17]: 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 = 3.3256
Checking accuracy on validation set
Got 103 / 1000 correct (10.30)
```

```
Iteration 100, loss = 2.1397
Checking accuracy on validation set
Got 342 / 1000 correct (34.20)
```

```
Iteration 200, loss = 2.1930
Checking accuracy on validation set
Got 377 / 1000 correct (37.70)
```

```
Iteration 300, loss = 2.0020
Checking accuracy on validation set
Got 388 / 1000 correct (38.80)
```

```
Iteration 400, loss = 2.1188
Checking accuracy on validation set
Got 423 / 1000 correct (42.30)
```

```
Iteration 500, loss = 1.7618
Checking accuracy on validation set
Got 407 / 1000 correct (40.70)
```



```
Iteration 600, loss = 1.5540
Checking accuracy on validation set
Got 414 / 1000 correct (41.40)
```

```
Iteration 700, loss = 1.7637
Checking accuracy on validation set
Got 444 / 1000 correct (44.40)
```

6.0.6 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.

```
[18]: 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)*****

# define learning rate
learning_rate = 1e-2
# define ThreeLayerConvNet
model = ThreeLayerConvNet(3, channel_1, channel_2, 10)
# use SGD gradient descent updater
optimizer = optim.SGD(model.parameters(), lr=learning_rate)
#pass

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

train_part34(model, optimizer)
```

```
Iteration 0, loss = 2.7388
Checking accuracy on validation set
Got 113 / 1000 correct (11.30)
```

```
Iteration 100, loss = 1.7880
```

```
Checking accuracy on validation set
Got 384 / 1000 correct (38.40)
```

```
Iteration 200, loss = 1.4596
Checking accuracy on validation set
Got 471 / 1000 correct (47.10)
```

```
Iteration 300, loss = 1.2621
Checking accuracy on validation set
Got 485 / 1000 correct (48.50)
```

```
Iteration 400, loss = 1.3655
Checking accuracy on validation set
Got 497 / 1000 correct (49.70)
```

```
Iteration 500, loss = 1.3256
Checking accuracy on validation set
Got 516 / 1000 correct (51.60)
```

```
Iteration 600, loss = 1.4794
Checking accuracy on validation set
Got 536 / 1000 correct (53.60)
```

```
Iteration 700, loss = 1.2768
Checking accuracy on validation set
Got 536 / 1000 correct (53.60)
```

7 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.

7.0.1 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.

```
[19]: # 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 = 2.3710
Checking accuracy on validation set
Got 161 / 1000 correct (16.10)
```

```
Iteration 100, loss = 1.6606
Checking accuracy on validation set
Got 396 / 1000 correct (39.60)
```

```
Iteration 200, loss = 1.8046
Checking accuracy on validation set
Got 391 / 1000 correct (39.10)
```

```
Iteration 300, loss = 1.7647
Checking accuracy on validation set
Got 416 / 1000 correct (41.60)
```

```
Iteration 400, loss = 1.7134
Checking accuracy on validation set
Got 448 / 1000 correct (44.80)
```

```
Iteration 500, loss = 1.7621
Checking accuracy on validation set
Got 467 / 1000 correct (46.70)
```

```
Iteration 600, loss = 1.6997
Checking accuracy on validation set
Got 416 / 1000 correct (41.60)
```

```
Iteration 700, loss = 1.6850
Checking accuracy on validation set
Got 467 / 1000 correct (46.70)
```

7.0.2 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 can use the default PyTorch weight initialization.

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.

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

model = None
optimizer = None

#####
# TODO: Rewrite the 2-layer ConvNet with bias from Part III with the #
# Sequential API. #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
# write the architecture in a sequential
model = nn.Sequential(
    # Convolutional layer (with bias) with 32 5x5 filters, with zero-padding of 2
    ↪2
    nn.Conv2d(3, channel_1, kernel_size=5, padding=2),
    # Nonlinearity
    nn.ReLU(),
    # Convolutional layer (with bias) with 16 3x3 filters, with zero-padding of 1
    ↪1
```

```

    nn.Conv2d(channel_1, channel_2, kernel_size=3, padding=1),
    # Nonlinearity
    nn.ReLU(),
    Flatten(),
    # Fully-connected layer (with bias) to compute scores for 10 classes
    nn.Linear(channel_2 * 1024, 10)
)
#pass
# define optimizer
optimizer = optim.SGD(model.parameters(), lr=learning_rate,
                        momentum=0.9, nesterov=True)
# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
#####
#                                     END OF YOUR CODE                                #
#####

train_part34(model, optimizer)

```

Iteration 0, loss = 2.2918
Checking accuracy on validation set
Got 102 / 1000 correct (10.20)

Iteration 100, loss = 1.6386
Checking accuracy on validation set
Got 468 / 1000 correct (46.80)

Iteration 200, loss = 1.5816
Checking accuracy on validation set
Got 511 / 1000 correct (51.10)

Iteration 300, loss = 1.5318
Checking accuracy on validation set
Got 525 / 1000 correct (52.50)

Iteration 400, loss = 1.4170
Checking accuracy on validation set
Got 547 / 1000 correct (54.70)

Iteration 500, loss = 1.3771
Checking accuracy on validation set
Got 546 / 1000 correct (54.60)

Iteration 600, loss = 1.0575
Checking accuracy on validation set
Got 548 / 1000 correct (54.80)

Iteration 700, loss = 1.1656

Checking accuracy on validation set
Got 562 / 1000 correct (56.20)

8 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>

8.0.1 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.

8.0.2 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.

8.0.3 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](#)

8.0.4 Have fun and happy training!

[28]:

```
#####
# 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
channel_list = [3, 20, 40, 80, 160, 20, 40, 50, 10]
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
# input is 3 * 32 * 32
model = nn.Sequential(
```

```

    nn.Conv2d(channel_list[0], channel_list[1], kernel_size=5, padding=1), # 32
    ↪x 32 x 20
    nn.ReLU(), # 32
    ↪x 32 x 20
    nn.MaxPool2d(2), # 16
    ↪x 16 x 20

    nn.Conv2d(channel_list[1], channel_list[2], kernel_size=3, padding=1), # 16
    ↪x 16 x 40
    nn.ReLU(), # 16
    ↪x 16 x 40
    nn.MaxPool2d(2), # 8
    ↪x 8 x 40

    nn.Conv2d(channel_list[2], channel_list[3], kernel_size=3, padding=1), # 8
    ↪x 8 x 80
    nn.ReLU(), # 8
    ↪x 8 x 80
    nn.MaxPool2d(2), # 4
    ↪x 4 x 80

    nn.Conv2d(channel_list[3], channel_list[4], kernel_size=3, padding=2), # 4
    ↪x 4 x 160
    nn.ReLU(), # 4
    ↪x 4 x 160
    nn.MaxPool2d(2), # 2
    ↪x 2 x 160

    # nn.Conv2d(channel_list[4], channel_list[5], kernel_size=3, padding=1), #
    ↪2 x 2 x 20
    # nn.ReLU(), #
    ↪2 x 2 x 20
    # nn.MaxPool2d(2), #
    ↪1 x 1 x 20
    Flatten(),
    nn.Linear(channel_list[4] * 4, 10)
)
#pass

optimizer = optim.SGD(model.parameters(), lr=learning_rate,
                        momentum=0.9, nesterov=True)
# *****END OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****
#####
#
#                               END OF YOUR CODE
#
#####

```



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

```
Iteration 0, loss = 2.3019  
Checking accuracy on validation set  
Got 84 / 1000 correct (8.40)
```

```
Iteration 100, loss = 2.1046  
Checking accuracy on validation set  
Got 260 / 1000 correct (26.00)
```

```
Iteration 200, loss = 1.9569  
Checking accuracy on validation set  
Got 424 / 1000 correct (42.40)
```

```
Iteration 300, loss = 1.5609  
Checking accuracy on validation set  
Got 456 / 1000 correct (45.60)
```

```
Iteration 400, loss = 1.4352  
Checking accuracy on validation set  
Got 463 / 1000 correct (46.30)
```

```
Iteration 500, loss = 1.6351  
Checking accuracy on validation set  
Got 470 / 1000 correct (47.00)
```

```
Iteration 600, loss = 1.4481  
Checking accuracy on validation set  
Got 526 / 1000 correct (52.60)
```

```
Iteration 700, loss = 1.2172  
Checking accuracy on validation set  
Got 564 / 1000 correct (56.40)
```

```
Iteration 0, loss = 1.2780  
Checking accuracy on validation set  
Got 591 / 1000 correct (59.10)
```

```
Iteration 100, loss = 1.5847  
Checking accuracy on validation set  
Got 599 / 1000 correct (59.90)
```

```
Iteration 200, loss = 1.0709  
Checking accuracy on validation set  
Got 610 / 1000 correct (61.00)
```

Iteration 300, loss = 0.9493
Checking accuracy on validation set
Got 578 / 1000 correct (57.80)

Iteration 400, loss = 0.9850
Checking accuracy on validation set
Got 658 / 1000 correct (65.80)

Iteration 500, loss = 1.2627
Checking accuracy on validation set
Got 644 / 1000 correct (64.40)

Iteration 600, loss = 0.8896
Checking accuracy on validation set
Got 658 / 1000 correct (65.80)

Iteration 700, loss = 0.9878
Checking accuracy on validation set
Got 648 / 1000 correct (64.80)

Iteration 0, loss = 0.9591
Checking accuracy on validation set
Got 644 / 1000 correct (64.40)

Iteration 100, loss = 0.8385
Checking accuracy on validation set
Got 686 / 1000 correct (68.60)

Iteration 200, loss = 0.8466
Checking accuracy on validation set
Got 674 / 1000 correct (67.40)

Iteration 300, loss = 1.0611
Checking accuracy on validation set
Got 693 / 1000 correct (69.30)

Iteration 400, loss = 0.8172
Checking accuracy on validation set
Got 704 / 1000 correct (70.40)

Iteration 500, loss = 0.9829
Checking accuracy on validation set
Got 702 / 1000 correct (70.20)

Iteration 600, loss = 0.8229
Checking accuracy on validation set
Got 688 / 1000 correct (68.80)

Iteration 700, loss = 0.8284
Checking accuracy on validation set
Got 713 / 1000 correct (71.30)

Iteration 0, loss = 0.7254
Checking accuracy on validation set
Got 716 / 1000 correct (71.60)

Iteration 100, loss = 0.9415
Checking accuracy on validation set
Got 708 / 1000 correct (70.80)

Iteration 200, loss = 0.9418
Checking accuracy on validation set
Got 704 / 1000 correct (70.40)

Iteration 300, loss = 0.6893
Checking accuracy on validation set
Got 695 / 1000 correct (69.50)

Iteration 400, loss = 0.8711
Checking accuracy on validation set
Got 697 / 1000 correct (69.70)

Iteration 500, loss = 1.0158
Checking accuracy on validation set
Got 733 / 1000 correct (73.30)

Iteration 600, loss = 0.7408
Checking accuracy on validation set
Got 706 / 1000 correct (70.60)

Iteration 700, loss = 0.6487
Checking accuracy on validation set
Got 735 / 1000 correct (73.50)

Iteration 0, loss = 0.5880
Checking accuracy on validation set
Got 725 / 1000 correct (72.50)

Iteration 100, loss = 0.6468
Checking accuracy on validation set
Got 725 / 1000 correct (72.50)

Iteration 200, loss = 0.6398
Checking accuracy on validation set
Got 728 / 1000 correct (72.80)

Iteration 300, loss = 0.6191
Checking accuracy on validation set
Got 732 / 1000 correct (73.20)

Iteration 400, loss = 0.6490
Checking accuracy on validation set
Got 715 / 1000 correct (71.50)

Iteration 500, loss = 0.5973
Checking accuracy on validation set
Got 722 / 1000 correct (72.20)

Iteration 600, loss = 0.6749
Checking accuracy on validation set
Got 718 / 1000 correct (71.80)

Iteration 700, loss = 0.8764
Checking accuracy on validation set
Got 730 / 1000 correct (73.00)

Iteration 0, loss = 0.4134
Checking accuracy on validation set
Got 743 / 1000 correct (74.30)

Iteration 100, loss = 0.6500
Checking accuracy on validation set
Got 720 / 1000 correct (72.00)

Iteration 200, loss = 0.4696
Checking accuracy on validation set
Got 757 / 1000 correct (75.70)

Iteration 300, loss = 0.5112
Checking accuracy on validation set
Got 727 / 1000 correct (72.70)

Iteration 400, loss = 0.5541
Checking accuracy on validation set
Got 747 / 1000 correct (74.70)

Iteration 500, loss = 0.6739
Checking accuracy on validation set
Got 719 / 1000 correct (71.90)

Iteration 600, loss = 0.4698
Checking accuracy on validation set
Got 722 / 1000 correct (72.20)

Iteration 700, loss = 0.5921
Checking accuracy on validation set
Got 740 / 1000 correct (74.00)

Iteration 0, loss = 0.4775
Checking accuracy on validation set
Got 731 / 1000 correct (73.10)

Iteration 100, loss = 0.3251
Checking accuracy on validation set
Got 748 / 1000 correct (74.80)

Iteration 200, loss = 0.5725
Checking accuracy on validation set
Got 729 / 1000 correct (72.90)

Iteration 300, loss = 0.4928
Checking accuracy on validation set
Got 743 / 1000 correct (74.30)

Iteration 400, loss = 0.3789
Checking accuracy on validation set
Got 740 / 1000 correct (74.00)

Iteration 500, loss = 0.5245
Checking accuracy on validation set
Got 730 / 1000 correct (73.00)

Iteration 600, loss = 0.4387
Checking accuracy on validation set
Got 743 / 1000 correct (74.30)

Iteration 700, loss = 0.3640
Checking accuracy on validation set
Got 741 / 1000 correct (74.10)

Iteration 0, loss = 0.5115
Checking accuracy on validation set
Got 753 / 1000 correct (75.30)

Iteration 100, loss = 0.5292
Checking accuracy on validation set
Got 738 / 1000 correct (73.80)

Iteration 200, loss = 0.5339
Checking accuracy on validation set
Got 749 / 1000 correct (74.90)

Iteration 300, loss = 0.5919
Checking accuracy on validation set
Got 703 / 1000 correct (70.30)

Iteration 400, loss = 0.3750
Checking accuracy on validation set
Got 743 / 1000 correct (74.30)

Iteration 500, loss = 0.4302
Checking accuracy on validation set
Got 756 / 1000 correct (75.60)

Iteration 600, loss = 0.4009
Checking accuracy on validation set
Got 725 / 1000 correct (72.50)

Iteration 700, loss = 0.7203
Checking accuracy on validation set
Got 736 / 1000 correct (73.60)

Iteration 0, loss = 0.3815
Checking accuracy on validation set
Got 738 / 1000 correct (73.80)

Iteration 100, loss = 0.3449
Checking accuracy on validation set
Got 720 / 1000 correct (72.00)

Iteration 200, loss = 0.3702
Checking accuracy on validation set
Got 744 / 1000 correct (74.40)

Iteration 300, loss = 0.4365
Checking accuracy on validation set
Got 736 / 1000 correct (73.60)

Iteration 400, loss = 0.5065
Checking accuracy on validation set
Got 731 / 1000 correct (73.10)

Iteration 500, loss = 0.4587
Checking accuracy on validation set
Got 728 / 1000 correct (72.80)

Iteration 600, loss = 0.3142
Checking accuracy on validation set
Got 739 / 1000 correct (73.90)

Iteration 700, loss = 0.4847
Checking accuracy on validation set
Got 747 / 1000 correct (74.70)

Iteration 0, loss = 0.3416
Checking accuracy on validation set
Got 748 / 1000 correct (74.80)

Iteration 100, loss = 0.2779
Checking accuracy on validation set
Got 755 / 1000 correct (75.50)

Iteration 200, loss = 0.3340
Checking accuracy on validation set
Got 713 / 1000 correct (71.30)

Iteration 300, loss = 0.2803
Checking accuracy on validation set
Got 724 / 1000 correct (72.40)

Iteration 400, loss = 0.3334
Checking accuracy on validation set
Got 736 / 1000 correct (73.60)

Iteration 500, loss = 0.2965
Checking accuracy on validation set
Got 729 / 1000 correct (72.90)

Iteration 600, loss = 0.5222
Checking accuracy on validation set
Got 722 / 1000 correct (72.20)

Iteration 700, loss = 0.5789
Checking accuracy on validation set
Got 729 / 1000 correct (72.90)

8.1 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.

Answer: I mostly follow the architectures provide above [conv-relu-pool] \times N \rightarrow [affine] \times M \rightarrow [softmax or SVM]. I try to made this architerture deeper, but the property of maxpool prevent me to doing that, so I only use 4 of depth in the end. I think this architerture it has very powerful non-linearity, so the result would be better than previous two-layers fully-connected and three layers

convolution neural network.

8.2 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.

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

```
Checking accuracy on test set  
Got 7056 / 10000 correct (70.56)
```


TensorFlow

November 4, 2022

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

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

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

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

1 Introduction to TensorFlow

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

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

1.1 Why do we use deep learning frameworks?

- Our code will now run on GPUs! This will allow our models to train much faster. 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).

- In 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 :)
- Finally, we want you to be exposed to the sort of deep learning code you might run into in academia or industry.

1.2 What is TensorFlow?

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.

1.3 How do 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 Tensorflow 2.x. Most examples on the web today are still in 1.x, so be careful not to confuse the two when looking up documentation.

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 GPU

You can manually switch to a GPU device on Colab by clicking **Runtime** -> **Change runtime type** and selecting **GPU** under **Hardware Accelerator**. You should do this before running the following cells to import packages, since the kernel gets restarted upon switching runtimes.

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

%matplotlib inline

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)
```

4 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.

```
[ ]: 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)

```

```

[ ]: 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)

```

```

[ ]: # 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

```

5 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.

5.0.1 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.

5.0.2 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

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

5.0.3 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 $N \times C \times H \times W$.

```
[ ]: 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))
```

```
[ ]: 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()
```

5.0.4 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.

```
[ ]: 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  $d1 * \dots * dM = D$ . The hidden layer will have  $H_{\square}$   
    ↪units,  
    and the output layer will produce scores for C classes.  
  
    Inputs:  
    - x: A TensorFlow Tensor of shape (N, d1, ..., dM) giving a minibatch of  
        input data.  
    - params: A list [w1, w2] of TensorFlow Tensors giving weights for the  
        network, where w1 has shape (D, H) and w2 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,  $\square$   
    ↪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  
  
[ ]: 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()

```

5.0.5 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 `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.

HINT: For convolutions: 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>

```

[ ]: 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)*****

    pass

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

```

After defing 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).

```

[ ]: 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()

```

5.0.6 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

- 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
- 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
- We'll mutate the weight values stored in a TensorFlow Tensor using `tf.assign_sub` ("sub" is for subtraction): https://www.tensorflow.org/api_docs/python/tf/assign_sub

```
[ ]: 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
```

```
[ ]: 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)

```

```

[ ]: 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))

```

5.0.7 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>

```
[ ]: 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)
```

5.0.8 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.

```
[ ]: 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)
```

5.0.9 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.

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

    pass

    # *****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)
```

6 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 unwieldy 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.

6.0.1 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

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.

```
[ ]: 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()
```

6.0.2 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/Dense


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

        pass

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

    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)*****

        pass

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

```
[ ]: 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()
```

6.0.3 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.

```
[ ]: 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 trainingn
    """
    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

```

6.0.4 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.

```
[ ]: 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)
```

6.0.5 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

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

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

          pass

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

    pass

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

train_part34(model_init_fn, optimizer_init_fn)

```

7 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.

7.0.1 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.

```

[ ]: 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)

```

7.0.2 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.

```

[ ]: 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)

```

7.0.3 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.

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

    pass

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

    pass

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

train_part34(model_init_fn, optimizer_init_fn)
```

We will also train this model with the built-in training loop APIs provided by TensorFlow.

```
[ ]: 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)
```

7.1 Part IV: Functional API

7.1.1 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.

```
[ ]: 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()
```

7.1.2 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.

```
[ ]: 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)
```

8 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.

8.0.1 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?

8.0.2 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/Dropout#methods

8.0.3 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.

8.0.4 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](#)

8.0.5 Have fun and happy training!

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

        pass

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

        # END OF YOUR CODE

    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)*****

        pass

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

        # END OF YOUR CODE

        return x

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.Adam(learning_rate)

train_part34(model_init_fn, optimizer_init_fn, num_epochs=num_epochs,
    ↪is_training=True)
```

8.1 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.

Answer: