

# Assignment 2: Convolutional Neural Networks with Pytorch

For this assignment, we're going to use one of most popular deep learning frameworks: PyTorch. And build our way through Convolutional Neural Networks.

## 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.

## Why?

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

## PyTorch versions

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

**If you are running on datahub, you shouldn't face any problem.**

You can also find the detailed PyTorch [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.

## 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-100 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, ResNet10 Implementation: we will implement ResNet10 from scratch given the architecture details
6. Part VI, CIFAR-100 open-ended challenge: please implement your own network to get as high accuracy as possible on CIFAR-100. 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

## Part I. Preparation

First, we load the CIFAR-100 dataset. This might take a couple minutes the first time you do it, but the files should stay cached after that.

```
In [21]: import os
os.environ["CUDA_DEVICE_ORDER"]="PCI_BUS_ID"   # see issue #152
os.environ["CUDA_VISIBLE_DEVICES"]="0"
```

```
In [13]: # Add official website of pytorch

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
```

```
In [14]: NUM_TRAIN = 49000
batch_size= 64

# 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.

#=====#
# You should try changing the transform for the training data to include      #
# data augmentation such as RandomCrop and HorizontalFlip                      #
# when running the final part of the notebook where you have to achieve      #
# as high accuracy as possible on CIFAR-100.                                  #
# Of course you will have to re-run this block for the effect to take place  #
#=====#
train_transform = transform = T.Compose([
    T.ToTensor(),
    T.Normalize((0.5071, 0.4867, 0.4408), (0.2675, 0.2565, 0.2761))
])

# 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-100
# training set into train and val sets by passing a Sampler object to the
# DataLoader telling how it should sample from the underlying Dataset.
cifar100_train = dset.CIFAR100('./datasets/cifar100', train=True, download=True,
                               transform=train_transform)
loader_train = DataLoader(cifar100_train, batch_size=batch_size, num_workers=2,
                           sampler=sampler.SubsetRandomSampler(range(NUM_TRAIN)))

cifar100_val = dset.CIFAR100('./datasets/cifar100', train=True, download=True,
                              transform=transform)
loader_val = DataLoader(cifar100_val, batch_size=batch_size, num_workers=2,
                        sampler=sampler.SubsetRandomSampler(range(NUM_TRAIN, 50000)))

cifar100_test = dset.CIFAR100('./datasets/cifar100', train=False, download=True,
                               transform=transform)
loader_test = DataLoader(cifar100_test, batch_size=batch_size, num_workers=2)

```

Files already downloaded and verified

Files already downloaded and verified

Files already downloaded and verified

You have an option to **use GPU by setting the flag to True below** (recommended). It is not necessary to use GPU for this assignment.

Note that if your computer does not have CUDA enabled, `torch.cuda.is_available()` will return False and this notebook will fallback to CPU mode. **You can run on GPU on datahub.**

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

```

In [15]: USE_GPU = True
num_class = 100
dtype = torch.float32 # we will be using float throughout this tutorial

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

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

```

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

```
using device: cuda
```

## Part II. Barebones PyTorch (10% of Grade)

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-100 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.

### PyTorch Tensors: Flatten Function

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

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

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

This is the right way to represent the data when we are doing something like a 2D convolution, that needs spatial understanding of where the intermediate features are relative to each other. When we use fully connected affine layers to process the image, however, we want each datapoint to be represented by a single vector -- it's no longer useful to segregate the different channels, rows, and columns of the data. So, we use a "flatten" operation to collapse the  $C \times H \times W$  values per representation into a single long vector. The flatten function below first reads in the N, C, H, and W values from a given batch of data, and then returns a "view" of that data. "View" is analogous to numpy's "reshape" method: it reshapes x's dimensions to be  $N \times ??$ , where ?? is allowed to be anything (in this case, it will be  $C \times H \times W$ , but we don't need to specify that explicitly).

```
In [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]])
```

## 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.

```
In [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  $d1 * \dots * dM = D$ . 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, num_class), dtype=dtype)
    scores = two_layer_fc(x, [w1, w2])
    print(scores.size()) # you should see [64, 100]

two_layer_fc_test()

torch.Size([64, 100])
```

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

```
In [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)*****

x = F.relu(F.conv2d(x, conv_w1, conv_b1, padding=2))
x = F.relu(F.conv2d(x, conv_w2, conv_b2, padding=1))
x = flatten(x)
x = x.mm(fc_w) + fc_b
scores = x

# *****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, 100).

```
In [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, num_class))
    fc_b = torch.zeros(num_class)

    scores = three_layer_convnet(x, [conv_w1, conv_b1, conv_w2, conv_b2, fc_w, fc_b])
    print(scores.size()) # you should see [64, 100]
three_layer_convnet_test()

torch.Size([64, 100])
```

## 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>

```
In [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))

```

```

Out[8]: tensor([[ -1.4538,  0.9761,  0.0736,  0.0133,  0.4468],
               [ 0.7883, -0.5139, -0.2843, -1.0967, -0.1476],
               [-0.5588, -0.1139,  0.0527,  0.1278,  0.8968]], device='cuda:0',
               requires_grad=True)

```

## 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.

```

In [16]: 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: 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))
return acc

```

## 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.

```

In [17]: 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: The accuracy of the model
"""
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 + 1) % print_every == 0:
        print('Iteration %d, loss = %.4f' % (t + 1, loss.item()))
        check_accuracy_part2(loader_val, model_fn, params)
        print()
return check_accuracy_part2(loader_val, model_fn, params)

```

## 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 100-dimensional vector that represents the probability distribution over 100 classes.

You don't need to tune any hyperparameters but you should see accuracies above 15% after training for one epoch.

```
In [11]: hidden_layer_size = 4000
         learning_rate = 1e-2

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

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

```
Iteration 100, loss = 3.9504  
Checking accuracy on the val set  
Got 89 / 1000 correct (8.90%)
```

```
Iteration 200, loss = 3.9258  
Checking accuracy on the val set  
Got 111 / 1000 correct (11.10%)
```

```
Iteration 300, loss = 4.1650  
Checking accuracy on the val set  
Got 126 / 1000 correct (12.60%)
```

```
Iteration 400, loss = 3.9018  
Checking accuracy on the val set  
Got 131 / 1000 correct (13.10%)
```

```
Iteration 500, loss = 3.4050  
Checking accuracy on the val set  
Got 143 / 1000 correct (14.30%)
```

```
Iteration 600, loss = 3.4194  
Checking accuracy on the val set  
Got 142 / 1000 correct (14.20%)
```

```
Iteration 700, loss = 3.8442  
Checking accuracy on the val set  
Got 149 / 1000 correct (14.90%)
```

```
Checking accuracy on the val set  
Got 175 / 1000 correct (17.50%)
```

```
Out[11]: 0.175
```

## BareBones PyTorch: Training a ConvNet

In the below cell 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 100 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 **12% after one epoch**.

```
In [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)*****

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

# *****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 100, loss = 4.4661
Checking accuracy on the val set
Got 62 / 1000 correct (6.20%)

Iteration 200, loss = 4.2148
Checking accuracy on the val set
Got 84 / 1000 correct (8.40%)

Iteration 300, loss = 4.1647
Checking accuracy on the val set
Got 104 / 1000 correct (10.40%)

Iteration 400, loss = 4.0284
Checking accuracy on the val set
Got 115 / 1000 correct (11.50%)

Iteration 500, loss = 3.5134
Checking accuracy on the val set
Got 131 / 1000 correct (13.10%)

Iteration 600, loss = 3.6170
Checking accuracy on the val set
Got 128 / 1000 correct (12.80%)

Iteration 700, loss = 3.5105
Checking accuracy on the val set
Got 134 / 1000 correct (13.40%)

Checking accuracy on the val set
Got 138 / 1000 correct (13.80%)
```

Out[12]: 0.138

## Part III. PyTorch Module API (10% of Grade)

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.

## Module API: Two-Layer Network

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

```
In [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, num_class)
    scores = model(x)
    print(scores.size()) # you should see [64, 100]
test_TwoLayerFC()

torch.Size([64, 100])

```

## 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.

```

In [14]: class ThreeLayerConvNet(nn.Module):
    def __init__(self, in_channel, channel_1, channel_2, num_classes):
        super().__init__()

        self.conv_1 = nn.Conv2d(in_channel, channel_1, (5,5), padding=2)
        nn.init.kaiming_normal_(self.conv_1.weight)
        self.conv_2 = nn.Conv2d(channel_1, channel_2, (3,3), padding=1)
        nn.init.kaiming_normal_(self.conv_2.weight)

        self.fc1 = nn.Linear(65536, num_classes)
        self.relu = nn.ReLU(inplace=True)

    def forward(self, x):
        scores = None
        #####
        # TODO: Implement the forward function for a 3-layer ConvNet. you      #
        # should use the layers you defined in __init__ and specify the        #
        # connectivity of those layers in forward()                            #
        #####
        # *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

        x = self.conv_1(x)
        x = self.relu(x)
        x = self.conv_2(x)
        x = self.relu(x)
        x = flatten(x)
        scores = self.fc1(x)

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

def test_ThreeLayerConvNet():
    x = torch.zeros((64, 3, 32, 32), dtype=dtype) # minibatch size 64, image size [3, 32, 32]
    model = ThreeLayerConvNet(in_channel=3, channel_1=32, channel_2=64, num_classes=num_class)
    scores = model(x)

```

```
print(scores.size()) # you should see [64, 100]
test_ThreeLayerConvNet()

torch.Size([64, 100])
```

## 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.

```
In [18]: 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))
        return acc
```

## 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.

```
In [19]: 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: The accuracy of the model
        """
        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 + 1) % print_every == 0:
                print('Epoch %d, Iteration %d, loss = %.4f' % (e, t + 1, loss.item()))
                check_accuracy_part34(loader_val, model)
                print()
        return check_accuracy_part34(loader_val, model)
```

## 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 8% after training for one epoch.

```
In [17]: hidden_layer_size = 4000
         learning_rate = 1e-3
         model = TwoLayerFC(3 * 32 * 32, hidden_layer_size, num_class)
         optimizer = optim.SGD(model.parameters(), lr=learning_rate)

         train_part34(model, optimizer)
```

```
Epoch 0, Iteration 100, loss = 4.5533  
Checking accuracy on validation set  
Got 33 / 1000 correct (3.30)
```

```
Epoch 0, Iteration 200, loss = 4.5609  
Checking accuracy on validation set  
Got 44 / 1000 correct (4.40)
```

```
Epoch 0, Iteration 300, loss = 4.6454  
Checking accuracy on validation set  
Got 54 / 1000 correct (5.40)
```

```
Epoch 0, Iteration 400, loss = 4.2244  
Checking accuracy on validation set  
Got 68 / 1000 correct (6.80)
```

```
Epoch 0, Iteration 500, loss = 4.1602  
Checking accuracy on validation set  
Got 70 / 1000 correct (7.00)
```

```
Epoch 0, Iteration 600, loss = 3.9680  
Checking accuracy on validation set  
Got 82 / 1000 correct (8.20)
```

```
Epoch 0, Iteration 700, loss = 4.1322  
Checking accuracy on validation set  
Got 89 / 1000 correct (8.90)
```

```
Checking accuracy on validation set  
Got 91 / 1000 correct (9.10)
```

```
Out[17]: 0.091
```

## 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 **accuracy above 14% after training for one epoch**.

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



```

In [18]: learning_rate = 1e-3
channel_1 = 32
channel_2 = 64

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

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

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

train_part34(model, optimizer, epochs=1)

```

```
Epoch 0, Iteration 100, loss = 4.3022  
Checking accuracy on validation set  
Got 92 / 1000 correct (9.20)
```

```
Epoch 0, Iteration 200, loss = 3.9539  
Checking accuracy on validation set  
Got 113 / 1000 correct (11.30)
```

```
Epoch 0, Iteration 300, loss = 3.8934  
Checking accuracy on validation set  
Got 133 / 1000 correct (13.30)
```

```
Epoch 0, Iteration 400, loss = 3.7666  
Checking accuracy on validation set  
Got 145 / 1000 correct (14.50)
```

```
Epoch 0, Iteration 500, loss = 3.7812  
Checking accuracy on validation set  
Got 146 / 1000 correct (14.60)
```

```
Epoch 0, Iteration 600, loss = 3.8107  
Checking accuracy on validation set  
Got 152 / 1000 correct (15.20)
```

```
Epoch 0, Iteration 700, loss = 3.3768  
Checking accuracy on validation set  
Got 170 / 1000 correct (17.00)
```

```
Checking accuracy on validation set  
Got 169 / 1000 correct (16.90)
```

```
Out[18]: 0.169
```

## Part IV. PyTorch Sequential API (10% of Grade)

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.

## 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 17% accuracy after one epoch of training.

```
In [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, num_class),
)

# 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)
```

```
Epoch 0, Iteration 100, loss = 3.9349  
Checking accuracy on validation set  
Got 117 / 1000 correct (11.70)
```

```
Epoch 0, Iteration 200, loss = 3.8074  
Checking accuracy on validation set  
Got 139 / 1000 correct (13.90)
```

```
Epoch 0, Iteration 300, loss = 3.8191  
Checking accuracy on validation set  
Got 141 / 1000 correct (14.10)
```

```
Epoch 0, Iteration 400, loss = 3.8954  
Checking accuracy on validation set  
Got 153 / 1000 correct (15.30)
```

```
Epoch 0, Iteration 500, loss = 3.5454  
Checking accuracy on validation set  
Got 150 / 1000 correct (15.00)
```

```
Epoch 0, Iteration 600, loss = 3.7902  
Checking accuracy on validation set  
Got 155 / 1000 correct (15.50)
```

```
Epoch 0, Iteration 700, loss = 3.4674  
Checking accuracy on validation set  
Got 183 / 1000 correct (18.30)
```

```
Checking accuracy on validation set  
Got 191 / 1000 correct (19.10)
```

```
Out[19]: 0.191
```

## 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 100 classes

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

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

Again, you don't need to tune any hyperparameters but you should see **accuracy above 14% after one epoch** of training.

```
In [20]: channel_1 = 32
channel_2 = 16
learning_rate = 1e-3

model = None
optimizer = None

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

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

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, epochs=1)
```

```
Epoch 0, Iteration 100, loss = 4.4582  
Checking accuracy on validation set  
Got 56 / 1000 correct (5.60)
```

```
Epoch 0, Iteration 200, loss = 4.0210  
Checking accuracy on validation set  
Got 90 / 1000 correct (9.00)
```

```
Epoch 0, Iteration 300, loss = 3.7806  
Checking accuracy on validation set  
Got 125 / 1000 correct (12.50)
```

```
Epoch 0, Iteration 400, loss = 3.7087  
Checking accuracy on validation set  
Got 137 / 1000 correct (13.70)
```

```
Epoch 0, Iteration 500, loss = 3.7134  
Checking accuracy on validation set  
Got 166 / 1000 correct (16.60)
```

```
Epoch 0, Iteration 600, loss = 3.5467  
Checking accuracy on validation set  
Got 158 / 1000 correct (15.80)
```

```
Epoch 0, Iteration 700, loss = 3.6695  
Checking accuracy on validation set  
Got 167 / 1000 correct (16.70)
```

```
Checking accuracy on validation set  
Got 171 / 1000 correct (17.10)
```

```
Out[20]: 0.171
```

## Part V. Resnet10 Implementation (35% of Grade)

In this section, you will use the tools introduced above to implement the Resnet architecture. The Resnet architecture was introduced in: <https://arxiv.org/pdf/1512.03385.pdf> and it has become one of the most popular architectures used for computer vision. The key feature of the resnet architecture is the presence of skip connections which allow for better gradient flow even for very deep networks. Therefore, unlike vanilla CNNs introduced above, we can effectively build Resnets models having more than 100 layers. However, for the purposes of this exercise we will be using a smaller Resnet-10 architecture shown in the diagram below:

layer name	output size	layer
conv1	16 x 16	7 x 7, 64, stride 2
conv2_x	8 x 8	3 x 3, maxpool, stride 2
		3 x 3, 64
		3 x 3, 64
conv3_x	8 x 8	3 x 3, 128
		3 x 3, 128
conv4_x	8 x 8	3 x 3, 256
		3 x 3, 256
conv5_x	4 x 4	3 x 3, 512
		3 x 3, 512
	1 x 1	average pool, 100-c Softmax

In the architecture above, the downsampling is performed in conv5\_1. We recommend using the adam optimizer for training Resnet. You should see about 45% accuracy in 10 epochs. The template below is based on the Module API but you are allowed to use other Pytorch APIs if you prefer.

```
In [21]: #####
# TODO: Implement the forward function for the Resnet specified #
# above. HINT: You might need to create a helper class to      #
```

```

# define a Resnet block and then use that block here to create          #
# the resnet layers i.e. conv2_x, conv3_x, conv4_x and conv5_x          #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

class ResidualBlock(nn.Module):

    def __init__(self, in_channels, out_channels, hidden_size = 64, skip_connection = True, kernel_size = 3, stride=1):
        super(ResidualBlock, self).__init__()

        self.skip_connection = skip_connection

        self.res_block = nn.Sequential(
            nn.Sequential(
                nn.Conv2d(in_channels = in_channels, out_channels=hidden_size, kernel_size = kernel_size, stride=1, padding=1),
                nn.ReLU(inplace=True)
            ),
            nn.Sequential(
                nn.Conv2d(in_channels = hidden_size, out_channels=out_channels, kernel_size = kernel_size, stride=1, padding=1),
                nn.ReLU(inplace=True)
            )
        )

    def forward(self, x):
        if self.skip_connection:
            return x + self.res_block(x)
        else:
            return self.res_block(x)

class ResNet(nn.Module):

    def __init__(self, in_channel, num_classes=10):
        super(ResNet, self).__init__()

        self.conv1 = nn.Sequential(
            nn.Conv2d(in_channels=in_channel, out_channels=64, kernel_size=7, stride=2, padding=3),
            nn.ReLU(inplace=True)
        )

        self.conv2_x = nn.Sequential(
            nn.MaxPool2d(kernel_size=3, stride=2, padding=1),

```



```

        ResidualBlock(in_channels=64, out_channels=64, hidden_size=64),
    )

    self.conv3_x = nn.Sequential(
        nn.Conv2d(in_channels=64, out_channels=128, kernel_size=1, stride=1, padding=0),
        nn.ReLU(inplace=True),
        ResidualBlock(in_channels=128, out_channels=128, hidden_size=128)
    )

    self.conv4_x = nn.Sequential(
        nn.Conv2d(in_channels=128, out_channels=256, kernel_size=1, stride=1, padding=0),
        nn.ReLU(inplace=True),
        ResidualBlock(in_channels=256, out_channels=256, hidden_size=256)
    )

    self.conv5_x = nn.Sequential(
        nn.Conv2d(in_channels=256, out_channels=512, kernel_size=1, stride=2, padding=0),
        nn.ReLU(inplace=True),
        ResidualBlock(in_channels=512, out_channels=512, hidden_size=512)
    )

    self.avgpool = nn.AdaptiveAvgPool2d((1, 1))
    self.fc = nn.Linear(512, num_classes)

    def forward(self, x):
        x = self.conv1(x)
        x = self.conv2_x(x)
        x = self.conv3_x(x)
        x = self.conv4_x(x)
        x = self.conv5_x(x)
        x = self.avgpool(x)
        x = x.view(x.size(0), -1)
        x = self.fc(x)
        return x

model = ResNet(in_channel=3, num_classes=num_class)
optimizer = optim.SGD(model.parameters(), lr=learning_rate, momentum=0.9, nesterov=True)

#####

```

```
#                                END OF YOUR CODE                                #
#####
```

```
In [22]: learning_rate = 1e-3

model = None
optimizer = None

#####
# TODO: Instantiate and train Resnet-10.                                     #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

model = ResNet(in_channel=3, num_classes=num_class)
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                                #
#####

print_every = 700
train_part34(model, optimizer, epochs=10)
print_every = 100
```

Epoch 0, Iteration 700, loss = 4.6075  
Checking accuracy on validation set  
Got 10 / 1000 correct (1.00)

Epoch 1, Iteration 700, loss = 4.5996  
Checking accuracy on validation set  
Got 19 / 1000 correct (1.90)

Epoch 2, Iteration 700, loss = 4.5070  
Checking accuracy on validation set  
Got 23 / 1000 correct (2.30)

Epoch 3, Iteration 700, loss = 4.3649  
Checking accuracy on validation set  
Got 43 / 1000 correct (4.30)

Epoch 4, Iteration 700, loss = 4.1418  
Checking accuracy on validation set  
Got 57 / 1000 correct (5.70)

Epoch 5, Iteration 700, loss = 3.8435  
Checking accuracy on validation set  
Got 76 / 1000 correct (7.60)

Epoch 6, Iteration 700, loss = 3.8109  
Checking accuracy on validation set  
Got 84 / 1000 correct (8.40)

Epoch 7, Iteration 700, loss = 3.4089  
Checking accuracy on validation set  
Got 108 / 1000 correct (10.80)

Epoch 8, Iteration 700, loss = 3.6177  
Checking accuracy on validation set  
Got 113 / 1000 correct (11.30)

Epoch 9, Iteration 700, loss = 3.5271  
Checking accuracy on validation set  
Got 141 / 1000 correct (14.10)

Checking accuracy on validation set  
Got 143 / 1000 correct (14.30)

## BatchNorm

Now you will also introduce the Batch-Normalization layer within the Resnet architecture implemented above. Please add a batch normalization layer after each conv in your network before applying the activation function (i.e. the order should be conv->BatchNorm->Relu). Please read the section 3.4 from the Resnet paper (<https://arxiv.org/pdf/1512.03385.pdf>).

Feel free to re-use the Resnet class that you have implemented above by introducing a boolean flag for batch normalization.

After trying out batch-norm, please discuss the performance comparison between Resnet with BatchNorm and without BatchNorm and possible reasons for why one performs better than the other.

```
In [28]: learning_rate = 1e-3

model = None
optimizer = None

#####
# TODO: InstantiateResnet with BatchNorm                                     #
#####
# *****START OF YOUR CODE (DO NOT DELETE/MODIFY THIS LINE)*****

class ResidualBlock(nn.Module):

    def __init__(self, in_channels, out_channels, hidden_size = 64, skip_connection = True, kernel_size = 3, stride

        super(ResidualBlock, self).__init__()

        self.skip_connection = skip_connection

        self.res_block = nn.Sequential(
            nn.Sequential(
                nn.Conv2d(in_channels = in_channels, out_channels=hidden_size, kernel_size = kernel_size, stride
                nn.BatchNorm2d(num_features=hidden_size),
                nn.ReLU(inplace=True)
            ),
            nn.Sequential(
```

```

        nn.Conv2d(in_channels = hidden_size, out_channels=out_channels, kernel_size = kernel_size, stri
        nn.BatchNorm2d(num_features=out_channels),
        nn.ReLU(inplace=True)
    )
)

def forward(self, x):
    if self.skip_connection:
        return x + self.res_block(x)
    else:
        return self.res_block(x)

class ResNetBN(nn.Module):
    def __init__(self, in_channel, num_classes=10):
        super(ResNetBN, self).__init__()

        self.conv1 = nn.Sequential(
            nn.Conv2d(in_channels=in_channel, out_channels=64, kernel_size=7, stride=2, padding=3),
            nn.BatchNorm2d(num_features=64),
            nn.ReLU(inplace=True)
        )

        self.conv2_x = nn.Sequential(
            nn.MaxPool2d(kernel_size=3, stride=2, padding=1),
            ResidualBlock(in_channels=64, out_channels=64, hidden_size=64),
        )

        self.conv3_x = nn.Sequential(
            nn.Conv2d(in_channels=64, out_channels=128, kernel_size=1, stride=1, padding=0),
            nn.BatchNorm2d(num_features=128),
            nn.ReLU(inplace=True),
            ResidualBlock(in_channels=128, out_channels=128, hidden_size=128)
        )

        self.conv4_x = nn.Sequential(
            nn.Conv2d(in_channels=128, out_channels=256, kernel_size=1, stride=1, padding=0),
            nn.BatchNorm2d(num_features=256),
            nn.ReLU(inplace=True),
            ResidualBlock(in_channels=256, out_channels=256, hidden_size=256)
        )

```

```

self.conv5_x = nn.Sequential(
    nn.Conv2d(in_channels=256, out_channels=512, kernel_size=1, stride=2, padding=0),
    nn.BatchNorm2d(num_features=512),
    nn.ReLU(inplace=True),
    ResidualBlock(in_channels=512, out_channels=512, hidden_size=512)
)

self.avgpool = nn.AdaptiveAvgPool2d((1, 1))
self.fc = nn.Linear(512, num_classes)

def forward(self, x):
    x = self.conv1(x)
    x = self.conv2_x(x)
    x = self.conv3_x(x)
    x = self.conv4_x(x)
    x = self.conv5_x(x)
    x = self.avgpool(x)
    x = x.view(x.size(0), -1)
    x = self.fc(x)
    return x

model = ResNetBN(in_channel=3, num_classes=num_class)
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
#####

print_every = 700
train_part34(model, optimizer, epochs=10)
print_every = 100

```

Epoch 0, Iteration 700, loss = 3.5591  
Checking accuracy on validation set  
Got 179 / 1000 correct (17.90)

Epoch 1, Iteration 700, loss = 3.0492  
Checking accuracy on validation set  
Got 256 / 1000 correct (25.60)

Epoch 2, Iteration 700, loss = 2.8418  
Checking accuracy on validation set  
Got 326 / 1000 correct (32.60)

Epoch 3, Iteration 700, loss = 2.2903  
Checking accuracy on validation set  
Got 342 / 1000 correct (34.20)

Epoch 4, Iteration 700, loss = 2.3426  
Checking accuracy on validation set  
Got 383 / 1000 correct (38.30)

Epoch 5, Iteration 700, loss = 2.1310  
Checking accuracy on validation set  
Got 409 / 1000 correct (40.90)

Epoch 6, Iteration 700, loss = 1.7745  
Checking accuracy on validation set  
Got 430 / 1000 correct (43.00)

Epoch 7, Iteration 700, loss = 1.8238  
Checking accuracy on validation set  
Got 400 / 1000 correct (40.00)

Epoch 8, Iteration 700, loss = 1.3763  
Checking accuracy on validation set  
Got 435 / 1000 correct (43.50)

Epoch 9, Iteration 700, loss = 1.3580  
Checking accuracy on validation set  
Got 448 / 1000 correct (44.80)

Checking accuracy on validation set  
Got 453 / 1000 correct (45.30)

## Discussion on BatchNorm

### TODO: Write your answer here

I see a significant jump in accuracy from ~14% to ~43% after using Batchnorm. This is because the goal of BN is to normalize the features (the output of each layer after going through the activation) to zero-mean state with standard deviation 1. This, supposedly, allows each layer to learn on a more stable distribution of inputs, and would thus accelerate the training of the network. Allows the values of a particular layer to become stable around a mean and variance with respect to change in input. Therefore, if input brings unstable changes, the derivative may be too big, or too small, resulting in an unstable learning model. Batch Normalization also acts as a form of regularization that helps to minimize over-fitting.

## Batch Size

In this exercise, we will study the effect of batch size on performance of ResNet (with BatchNorm).

Specifically, you should try batch sizes of 32, 64 and 128 and describe the effect of varying batch size. You should print the validation accuracy of using each batch size in different rows.

After trying out different batch size, please discuss the effect of different batch sizes and possible reasons for that (either they are showing some trend or not).

```
In [29]: print_every = 9999
batch_sizes = [32, 64, 128]
learning_rate = 1e-3
model = None
optimizer = None

#####
# TODO: Try Resnet with different batch sizes. Hint: You will need to      #
# create a new dataloader with appropriate batch size for each experiment.  #
# You will also need to store the final accuracy for each experiment        #
```



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

model = ResNetBN(in_channel=3, num_classes=num_class)
optimizer = optim.SGD(model.parameters(), lr=learning_rate, momentum=0.9, nesterov=True)

for batch_size in batch_sizes:
    print("Batch size: ", batch_size)
    loader_train = DataLoader(cifar100_train, batch_size=batch_size, num_workers=2,
                              sampler=sampler.SubsetRandomSampler(range(NUM_TRAIN)))
    loader_val = DataLoader(cifar100_val, batch_size=batch_size, num_workers=2,
                            sampler=sampler.SubsetRandomSampler(range(NUM_TRAIN, 50000)))

    print_every = 700
    train_part34(model, optimizer, epochs=10)
    print("-----")

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

Batch size: 32

Epoch 0, Iteration 700, loss = 3.4582

Checking accuracy on validation set

Got 158 / 1000 correct (15.80)

Epoch 0, Iteration 1400, loss = 3.1590

Checking accuracy on validation set

Got 205 / 1000 correct (20.50)

Epoch 1, Iteration 700, loss = 3.1147

Checking accuracy on validation set

Got 270 / 1000 correct (27.00)

Epoch 1, Iteration 1400, loss = 2.9196

Checking accuracy on validation set

Got 306 / 1000 correct (30.60)

Epoch 2, Iteration 700, loss = 2.6960

Checking accuracy on validation set

Got 345 / 1000 correct (34.50)

Epoch 2, Iteration 1400, loss = 2.5932

Checking accuracy on validation set

Got 358 / 1000 correct (35.80)

Epoch 3, Iteration 700, loss = 2.3798

Checking accuracy on validation set

Got 379 / 1000 correct (37.90)

Epoch 3, Iteration 1400, loss = 2.2778

Checking accuracy on validation set

Got 395 / 1000 correct (39.50)

Epoch 4, Iteration 700, loss = 2.5746

Checking accuracy on validation set

Got 426 / 1000 correct (42.60)

Epoch 4, Iteration 1400, loss = 1.8106

Checking accuracy on validation set

Got 414 / 1000 correct (41.40)

Epoch 5, Iteration 700, loss = 2.1971  
Checking accuracy on validation set  
Got 439 / 1000 correct (43.90)

Epoch 5, Iteration 1400, loss = 1.6585  
Checking accuracy on validation set  
Got 442 / 1000 correct (44.20)

Epoch 6, Iteration 700, loss = 1.8987  
Checking accuracy on validation set  
Got 459 / 1000 correct (45.90)

Epoch 6, Iteration 1400, loss = 1.6806  
Checking accuracy on validation set  
Got 475 / 1000 correct (47.50)

Epoch 7, Iteration 700, loss = 1.7205  
Checking accuracy on validation set  
Got 451 / 1000 correct (45.10)

Epoch 7, Iteration 1400, loss = 1.5817  
Checking accuracy on validation set  
Got 473 / 1000 correct (47.30)

Epoch 8, Iteration 700, loss = 1.0182  
Checking accuracy on validation set  
Got 459 / 1000 correct (45.90)

Epoch 8, Iteration 1400, loss = 1.3458  
Checking accuracy on validation set  
Got 458 / 1000 correct (45.80)

Epoch 9, Iteration 700, loss = 1.2405  
Checking accuracy on validation set  
Got 466 / 1000 correct (46.60)

Epoch 9, Iteration 1400, loss = 0.7795  
Checking accuracy on validation set  
Got 484 / 1000 correct (48.40)

Checking accuracy on validation set  
Got 453 / 1000 correct (45.30)

-----  
Batch size: 64

Epoch 0, Iteration 700, loss = 0.5807

Checking accuracy on validation set  
Got 492 / 1000 correct (49.20)

Epoch 1, Iteration 700, loss = 0.3393

Checking accuracy on validation set  
Got 462 / 1000 correct (46.20)

Epoch 2, Iteration 700, loss = 0.2340

Checking accuracy on validation set  
Got 493 / 1000 correct (49.30)

Epoch 3, Iteration 700, loss = 0.1611

Checking accuracy on validation set  
Got 470 / 1000 correct (47.00)

Epoch 4, Iteration 700, loss = 0.0731

Checking accuracy on validation set  
Got 474 / 1000 correct (47.40)

Epoch 5, Iteration 700, loss = 0.0619

Checking accuracy on validation set  
Got 477 / 1000 correct (47.70)

Epoch 6, Iteration 700, loss = 0.0360

Checking accuracy on validation set  
Got 473 / 1000 correct (47.30)

Epoch 7, Iteration 700, loss = 0.0343

Checking accuracy on validation set  
Got 483 / 1000 correct (48.30)

Epoch 8, Iteration 700, loss = 0.0317

Checking accuracy on validation set  
Got 479 / 1000 correct (47.90)

```
Epoch 9, Iteration 700, loss = 0.0197  
Checking accuracy on validation set  
Got 480 / 1000 correct (48.00)
```

```
Checking accuracy on validation set  
Got 479 / 1000 correct (47.90)
```

```
-----  
Batch size: 128  
Checking accuracy on validation set  
Got 477 / 1000 correct (47.70)  
-----
```

## Discuss effect of Batch Size

I see slight increase in accuracy when batch size increased from 32 to 64 and the accuracy almost remains same for batch size 64 and 128. I reckon there are 2 opposing factors playing out here:

1. Normalisation: When the batch size is large can lead to more effective normalization. When batch size is small, the mean and standard deviation computed over the batch can be noisy, which can lead to less effective normalization.
2. Batch size: Using a largert batch size reduces regularization as batchnorm operation adds a noise to activation of each layer. This effect reduced when we use larger batch size but leads to a better generalisation when the batch size is smaller. This may be due to higher noise added to activations at smaller batch size.

Hence, when I increase batch size from 32 to 64, the activation are more effectively normalized across the batch. However, regularisation effect decreases when batch size is increased from 64 to 128.

## Part VI. CIFAR-100 open-ended challenge (25% of Grade)

In this section, you can experiment with whatever ConvNet architecture you'd like on CIFAR-100 **except Resnet** because we already tried it.

Now it's your job to experiment with architectures, hyperparameters, loss functions, and optimizers to train a model that achieves **at least 48%** accuracy on the CIFAR-100 **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.

- 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>

## Things you might try:

- **Filter size:** Above we used 5x5; would smaller filters be more efficient?
- **Adam Optimizer:** Above we used SGD optimizer, would an Adam optimizer do better?
- **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? You can also try out LayerNorm and GroupNorm.
- **Network architecture:** 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.

## 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.

## Want more improvements?

There are many other features you can implement to try and improve your performance.

- 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
  - [DenseNets](#) where inputs into previous layers are concatenated together.

## Have fun and may the gradients be with you!

```
In [36]: #####
# TODO:
# Experiment with any architectures, optimizers, and hyperparameters.
# Achieve AT LEAST 48% 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.
#####

import torch.nn.functional as F
```

```
model = None
optimizer = None

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

class InceptionModule(nn.Module):

    def __init__(self, in_channels, out_channels, filters = (3, 5, 7)):
        super(InceptionModule, self).__init__()
        """
        Module to apply multi-scale convolutions in parallel and finally concatenating
        Args:
            in_channels: number of input channels
            out_channels: number of output channels
            filters: tuple of filters to apply
        Returns:
            output of the module
        """

        self.conv1 = nn.Sequential(
            nn.Conv2d(in_channels=in_channels, out_channels=out_channels, kernel_size=filters[0], stride=1, padding=1),
            nn.BatchNorm2d(num_features=out_channels),
            nn.ReLU(inplace=True)
        )

        self.conv2 = nn.Sequential(
            nn.Conv2d(in_channels=in_channels, out_channels=out_channels, kernel_size=filters[1], stride=1, padding=1),
            nn.BatchNorm2d(num_features=out_channels),
            nn.ReLU(inplace=True)
        )

        self.conv3 = nn.Sequential(
            nn.Conv2d(in_channels=in_channels, out_channels=out_channels, kernel_size=filters[2], stride=1, padding=1),
            nn.BatchNorm2d(num_features=out_channels),
            nn.ReLU(inplace=True)
        )

        self.maxpool = nn.Sequential(
            nn.MaxPool2d(kernel_size=3, stride=1, padding=1),
            nn.Conv2d(in_channels=in_channels, out_channels=out_channels, kernel_size=1, stride=1, padding=0),
```



```

        nn.BatchNorm2d(num_features=out_channels),
        nn.ReLU(inplace=True)
    )

    self.conv4 = nn.Sequential(
        nn.Conv2d(in_channels=4*out_channels, out_channels=out_channels, kernel_size=1, stride=1, padding=0),
        nn.BatchNorm2d(num_features=out_channels),
        nn.ReLU(inplace=True)
    )

    def forward(self, x):
        x1 = self.conv1(x)
        x2 = self.conv2(x)
        x3 = self.conv3(x)
        x4 = self.maxpool(x)
        return self.conv4(torch.cat([x1, x2, x3, x4], dim=1))

class ConvNetBlock(nn.Module):

    def __init__(self, in_channels, out_channels, hidden_size = 64, skip_connection = True, kernel_size = 3, stride
        super(ConvNetBlock, self).__init__()

        self.skip_connection = skip_connection

        self.res_block = nn.Sequential(
            nn.Sequential(
                nn.Conv2d(in_channels = in_channels, out_channels=hidden_size, kernel_size = kernel_size, stride
                nn.BatchNorm2d(num_features=hidden_size),
                nn.ReLU(inplace=True)
            ),
            nn.Sequential(
                nn.Conv2d(in_channels = hidden_size, out_channels=out_channels, kernel_size = kernel_size, stride
                nn.BatchNorm2d(num_features=out_channels),
                nn.ReLU(inplace=True)
            )
        )

    def forward(self, x):
        return self.res_block(x)

```

```
class CustomNetwork(nn.Module):
    def __init__(self, in_channel, num_classes=10):
        super(CustomNetwork, self).__init__()

        self.conv1 = nn.Sequential(
            nn.Conv2d(in_channels=in_channel, out_channels=64, kernel_size=7, stride=2, padding=3),
            nn.BatchNorm2d(num_features=64),
            nn.ReLU(inplace=True)
        )

        # Module to apply multiple filters to the same input and concatenate the results
        self.conv2_x = nn.Sequential(
            nn.MaxPool2d(kernel_size=3, stride=2, padding=1),
            ConvNetBlock(in_channels=64, out_channels=64, hidden_size=64),
        )

        self.inception1 = InceptionModule(in_channels=64, out_channels=64)

        self.conv3_x = nn.Sequential(
            nn.Conv2d(in_channels=64, out_channels=128, kernel_size=1, stride=1, padding=0),
            nn.BatchNorm2d(num_features=128),
            nn.ReLU(inplace=True),
            ConvNetBlock(in_channels=128, out_channels=128, hidden_size=128)
        )

        self.inception2 = InceptionModule(in_channels=128, out_channels=128)

        self.conv4_x = nn.Sequential(
            nn.Conv2d(in_channels=128, out_channels=256, kernel_size=1, stride=1, padding=0),
            nn.BatchNorm2d(num_features=256),
            nn.ReLU(inplace=True),
            ConvNetBlock(in_channels=256, out_channels=256, hidden_size=256)
        )

        #self.inception3 = InceptionModule(in_channels=256, out_channels=256)

        self.conv5_x = nn.Sequential(
            nn.Conv2d(in_channels=256, out_channels=512, kernel_size=1, stride=2, padding=0),
            nn.BatchNorm2d(num_features=512),
            nn.ReLU(inplace=True),
            ConvNetBlock(in_channels=512, out_channels=512, hidden_size=512)
        )
```

```

        #self.inception4 = InceptionModule(in_channels=512, out_channels=512)

        self.avgpool = nn.AdaptiveAvgPool2d((1, 1))
        self.fc = nn.Linear(512, num_classes)

    def forward(self, x):
        x = self.conv1(x)
        x = self.conv2_x(x)
        x = self.inception1(x) + x
        x = self.conv3_x(x)
        x = self.inception2(x) + x
        x = self.conv4_x(x)
        #x = self.inception3(x)
        x = self.conv5_x(x)
        #x = self.inception4(x)
        x = self.avgpool(x)
        x = x.view(x.size(0), -1)
        x = self.fc(x)
        return x

#model = ResNetBN(in_channel=3, num_classes=num_class)
model = CustomNetwork(in_channel=3, num_classes=num_class)
optimizer = optim.SGD(model.parameters(), lr=1e-2, momentum=0.9, nesterov=True, weight_decay=0.0001)

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

# You should get at least 48% accuracy.
print_every = 700
train_part34(model, optimizer, epochs=10)
print_every = 100

```

Epoch 0, Iteration 700, loss = 3.3624  
Checking accuracy on validation set  
Got 156 / 1000 correct (15.60)

Epoch 1, Iteration 700, loss = 2.7594  
Checking accuracy on validation set  
Got 275 / 1000 correct (27.50)

Epoch 2, Iteration 700, loss = 2.1948  
Checking accuracy on validation set  
Got 330 / 1000 correct (33.00)

Epoch 3, Iteration 700, loss = 2.4755  
Checking accuracy on validation set  
Got 387 / 1000 correct (38.70)

Epoch 4, Iteration 700, loss = 1.9746  
Checking accuracy on validation set  
Got 411 / 1000 correct (41.10)

Epoch 5, Iteration 700, loss = 1.6572  
Checking accuracy on validation set  
Got 419 / 1000 correct (41.90)

Epoch 6, Iteration 700, loss = 1.9283  
Checking accuracy on validation set  
Got 443 / 1000 correct (44.30)

Epoch 7, Iteration 700, loss = 1.4666  
Checking accuracy on validation set  
Got 461 / 1000 correct (46.10)

Epoch 8, Iteration 700, loss = 1.2876  
Checking accuracy on validation set  
Got 451 / 1000 correct (45.10)

Epoch 9, Iteration 700, loss = 1.0947  
Checking accuracy on validation set  
Got 475 / 1000 correct (47.50)

```
Checking accuracy on validation set  
Got 455 / 1000 correct (45.50)
```

## Describe what you did (10% of Grade)

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.

I tried few things

1. Added Inception module to the network: A series of (7x7), (5x5), (3x3) convolution filter and Maxpool layers are applied in parallel followed by 1d convolution. Got a best case accuracy of 45% with learning rate of 0.01 in the SGD optimizer
2. Added skip connections to the inception module: With same settings, accuracy boosted to 47%
3. Changed optimizer to Adam and got a subpar accuracy of less than 10%
4. Tuned learning rate at a log scale from 0.1 to 0.0001 and found 0.01 performs best for the network
5. Used weight decay at three different value - 0.9, 0.1 and 0.01 but didn't get the accuracy boost as desired. I think after 10 epochs, the model is still not overfitting. Finally settled with a weight decay of 0.0001 to add some regularization.

## 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.

```
In [37]: best_model = model  
         check_accuracy_part34(loader_test, best_model)
```

```
Checking accuracy on test set  
Got 4769 / 10000 correct (47.69)
```

```
Out[37]: 0.4769
```

```
In [ ]:
```