# COMS 4732: Advanced Computer Vision Homework 3

## Overview

In this assignment, we'll build a simple convolutional neural network in PyTorch and Train it to recognize natural objects using the CIFAR-10 dataset. Training a classifier on the CIFAR-10 dataset can be regarded as the hello world of image recognition. The structure of this assignment is:

```
1. Introduction
2. Setting up the Environment
3. Preparing the Data
4. Building the Network
5. Training the Model and Evaluating the Performance
6. Understand the Deep Neural Networks
```

(The full content can be better viewed in jupyter notebook).

## Submission Instructions

- Please submit the notebook (ipynb and pdf) including the output of all cells. Please note that you should export your completed notebook as a PDF (CV2_HW3_UNI.pdf) and upload it to GradeScope.
- Then, please submit the executable completed notebook (CV2_HW3_UNI.ipynb) to Cousework.
- For both files, 1) remember to replace UNI with your own uni; 2) the completed notebook should show your code, as well as the log and image.

## Introduction

Object Recognition typically aims to train a model to recognize, identify, and locate the objects in the images with a given degree of confidence. Object recognition is also called image recognition in many other kinds of literature, which is strictly related to computer vision, which we define as the art and science of making computers understand images.

Deep learning has contributed quite a lot to the development of object recognition algorithms and has achieved many surprising results on many benchmark datasets, such as ImageNet, MS COCO, LVIS. Also, different architectures are proposed to improve the accuracy and efficiency of the learned models. As a brief introduction, we first go through some basic concepts in object recognition and give the following tasks:

1. Classification.
2. Detection.
3. Segmentation.

In this assignment, we only dive into the image classification task and we start with a small model on a toy dataset. We do encourage you to explore the other tasks. Here are a few interesting references:

1. https://towardsdatascience.com/creating-your-own-object-detector-ad69dda69c85
2. https://medium.com/analytics-vidhya/iou-intersection-over-union-705a39e7acef
3. https://www.analyticsvidhya.com/blog/2018/11/implementation-faster-r-cnn-python-object-detection/
4. https://blog.paperspace.com/mask-r-cnn-in-tensorflow-2-0/
5. https://www.analyticsvidhya.com/blog/2021/05/an-introduction-to-few-shot-learning/

```
1 from google.colab import drive
2 drive.mount('/content/drive')

   Mounted at /content/drive
```

```
1 %cd /content/drive/MyDrive/CV2_HW3/Figures
2
```

## Classification

Classification is an important task in Object recognition and aims to identify what is in the image and with what level of confidence. The general pipeline of this task is straightforward. It starts with the definition of the ontology, i.e. the class of objects to detect. Then, the classification task identifies what is in the image (associated level of confidence) and picks up the class with the highest confidence. Usually, each test image used for classification has one dominant object in the image and we directly use the representation of the whole image for classification.

For the example given above, the classification algorithm will only remember that there is a dog, ignoring all other classes. Similar to classification, tagging also predicts confidence level for each class but aims to recognize multiple ones for a given image. For the example given below, it will try to return all the best classes corresponding to the image.

## Detection and segmentation

Once identified what is in the image, we want to locate the objects. There are two ways to do so: detection and segmentation.Detection outputs a rectangle, also called bounding box, where the objects are. It is a very robust technology, prone to minor errors and imprecisions.

Segmentation identifies the objects for each pixel in the image, resulting in a very precise map. However, the accuracy of segmentation depends on an extensive and often time-consuming training of the neural network.

Different from Classification, Detection and Segmentation ask the algorithm to return the location of the objects in the image. However, as no prior is given to predict the location, a set of potential locations/regions are enumerated. Then, the algorithm performs classification and localization prediction on each of the enumerated subregions. To note, the enumerated sub-regions differ from each other regarding the center location, aspect ratio, and size, while all of them are pre-defined by humans and all of the enumerated sub-regions can densely cover the full image. The location prediction is in effect predicting the offset w.r.t. each the sub-region.

For object detection/segmentation, the two-stage and the one-stage frameworks have been dominating the related research area. Recently, the end-to-end object detector has been proposed to perform prediction in a unified manner (i.e., no pre-defined sub-region will be densely enumerated).

## Setting up the Environment

To learn a deep neural network, we need to first build proper environment and we use Pytorch. For some basic overview and features offered in Colab notebooks, check out: Overview of Colaboratory Features.

For this homework, you may need you Colab for your experiments. Under your columbia account, you can open the colab and then upload this jupyter notebook. You need to use the Colab GPU for this assignment by selecting:

> **Runtime → Change runtime type → Hardware Accelerator: GPU**

You should be capable to directly run the following cells in most cases. However, in case you need, you can use the following code to install pytorch when you open the CoLab.

```
!pip install torch torchvision
!pip install Pillow==4.0.0
```

```
1  """
```

```
 2 Below are some input statements which basically imports numpy, torch, torchvision,
 3 """
 4
 5 import cv2
 6 import torch
 7 import torchvision
 8 import torchvision.transforms as transforms
 9 from torchvision.datasets import CIFAR10
10 from torch.utils.data import DataLoader
11 import torchvision.models as models
12 import matplotlib.pyplot as plt
13 import numpy as np
14 from copy import deepcopy
15 from torchvision.utils import make_grid
16
17 import torch.nn as nn
18 import torch.nn.functional as F
19 import torch.optim as optim
20 from torch.utils.data import random_split
21
22 import math
23 import os
24 import argparse
```

## Warm-Up for Pytorch

```
 1 #Pytorch is very similar to numpy. Basically, any operation in numpy can be easily
 2 # The following is a few examples.
 3
 4 # Operations
 5 y = torch.rand(2, 2)
 6 x = torch.rand(2, 2)
 7
 8 # multiplication
 9 z = x * y
10 z = torch.mul(x,y) #elementwise #must be broadcastable
11
12 #matrix multiplication
13 tensor1 = torch.randn(3, 4)
14 tensor2 = torch.randn(4, 3)
15 torch.matmul(tensor1, tensor2)
16
17 # NumPy conversion
18 x = torch.rand(2,2)
19 y = x.numpy()
20 print(type(y))
21
22 z1 = torch.from_numpy(y) #sharing the memory space with the numpy ndarray
23 z2 = torch.tensor(y) #a copy
```

```
24 print(type(z1))
25 print(type(z2))
26
27 # Pytorch attributes and functions for tensors
28 print(x.shape)
29 print(x.device)
30
31 # autograd
32
33 # requires grad equals true lets us compute gradients on the tenor
34 x = torch.tensor([2,3,5], dtype=float, requires_grad=True)
35 y =(5 * x**2).sum()
36
37 # When we finish our computation we can call .backward() and have all the gradients
38 # The gradient for this tensor will be accumulated into .grad attribute
39 y.backward()
40 #print(z.grad) # dz/dz
41 print(x.grad) # dz/dx
42
43 # autograd requires computational resources and can take time.
44 # disable autograd for model eval by writing your evaluation code in
45 # As such, with torch.no_grad() is usually used in evaluation part
```

```
<class 'numpy.ndarray'>
<class 'torch.Tensor'>
<class 'torch.Tensor'>
torch.Size([2, 2])
cpu
tensor([20., 30., 50.], dtype=torch.float64)
```

For repeatable experiments, we are recommended to set random seeds for anything using random number generation - this means numpy and random as well! It's also worth mentioning that cuDNN uses nondeterministic algorithms which can be disabled by setting torch.backends.cudnn.enabled = False.

```
1 experiment_name = 'debug'  #Provide name to model experiment
2 model_name = 'basic' # Choose between [basic, alexnet]
3 batch_size = 100  #You may not need to change this but incase you do
4
5 torch.manual_seed(42)
```

```
<torch._C.Generator at 0x7f921432bab0>
```

## Preparing the Data

We provide the data-loading functions and a few helper functions below.

To learn a deep learning algorithm, we need to prepare the data and this is where TorchVision comes into play. Usually, we have a training dataset and a test dataset. Each image in the training dataset is paired with a class label, serving as groundtruth to guide the updating of the deep neural network. In real-world test scenario, no label will be provided. However, we still have the labels in the test set to calculate the accuracy. Then, the accuracy on test data is used to quantitatively evaluate the performance, i.e., the generalization of the learned deep algorithm. After all, it will be unexpected if the model works perfectly on the training data but fails on the testing data.

During network training, we usually perform data augmentation on the training data by manipulating the value of pixels. To determine how we adjust the pixel values, we first determine the effects to be applied to the images. TorchVision offers a lot of handy transformations, such as cropping or normalization. For example, shifting the image, adjusting the light and contrast of an image, translating an RGB image into a grayscale image, image rotation, and so on. Also, torchvision can apply random sections among the effects.

To note, the effects to be added on the training data should also seriously consider the property of the dataset. For example, for the digit recognition in MNIST, can we perform 180-degree rotation on the image of digit 6 without altering the image label during network training?

```
1 def get_transform(model_name):
2
3     if model_name == 'alexnet':
4         transform = transforms.Compose([
5             transforms.Resize((227, 227)),
6             transforms.ToTensor(),
7             transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5)),
8         ])
9
10    else:
11
12        transform = transforms.Compose([
13            transforms.ToTensor(),
14            transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5)),
15        ])
16
17    return transform
18
19
20 def get_dataset(model_name, train_percent=0.9):
21     '''
22     Returns the train, val and test torch.Datasetin addition to a list of classes,
23     to the label used for it in the data
24
25
26     Reference for transforms in torchvision: https://pytorch.org/vision/stable/tran
27     @model_name: either 'basic' or 'alexnet'
28     @train_percent: percent of training data to keep for training. Rest will be val
29     '''
```

```python
30
31     transform  = get_transform(model_name)
32
33     train_data = CIFAR10(root='./data', train=True, download=first_run, transform=t
34     test_data  = CIFAR10(root='./data', train=False, download=first_run, transform=
35
36     train_size = int(train_percent * len(train_data))
37     val_size = len(train_data) - train_size
38
39
40     train_data, val_data = random_split(train_data, [train_size, val_size])
41     class_names = ('plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', '
42
43     return train_data, val_data, test_data, class_names
44
45
46 def get_dataloader(batch_size, num_workers=1, model_name='basic'):
47     '''
48     Returns the train, val and test dataloaders in addition to a list of classes, v
49     to the label used for it in the data
50
51     Reference for dataloader class: https://pytorch.org/docs/stable/data.html
52     @batch_size: batch to be used by dataloader
53     @num_workers: number of dataloader workers/instances used
54     @model_name: either 'basic' or 'alexnet'
55     '''
56
57     train_set, val_set, test_set, class_names = get_dataset(model_name)
58     trainloader = DataLoader(train_set, batch_size=batch_size, shuffle=True, num_wo
59     valloader = DataLoader(val_set, batch_size=batch_size, shuffle=False, num_worke
60     testloader = DataLoader(test_set, batch_size=batch_size, shuffle=False, num_wor
61
62
63     return trainloader, valloader, testloader, class_names
64
65 def makegrid_images(model_name='basic'):
66     '''
67     For visualization purposes
68
69     @model_name: either 'basic' or 'alexnet'
70     '''
71
72
73     _, trial_loader, _, _ = get_dataloader(32, model_name=model_name)
74     images, labels = iter(trial_loader).next()
75
76     grid = make_grid(images)
77
78     return grid
79
80 def show_img(img, mean=(0.5, 0.5, 0.5), std=(0.5, 0.5, 0.5), viz=True, norm=True):
```

```
 81        '''
 82       For visualization purposes
 83
 84       B: batch size
 85       C: channels
 86       H: height
 87       W: width
 88
 89       @img: torch.Tensor for the image of type (B, C, H, W)
 90       @mean: mean used for normalizing along 3 dimensions C, H, W in get_transform
 91       @std:  std. deviation used for normalizing along 3 dimensions C, H, WW in get_t
 92       @viz: whether or not to plt.plot or just return the unnormalized image
 93       @norm: whether or not unnormalize. Unnormalizes if true.
 94
 95       Returns:
 96       Viewable image in (H, W, C) as a numpy array
 97       '''
 98
 99
100       if norm:
101           for idx in range(img.shape[0]):
102
103               img[idx] = img[idx] * std[idx] + mean[idx]
104
105       image = np.asarray(img)
106
107       if viz:
108           if len(image.shape) == 4:
109               image = image.squeeze()
110
111           plt.imshow(np.transpose(image, (1, 2, 0)))
112           plt.show()
113
114       return np.transpose(image.squeeze(), (1, 2, 0))
```

## Building the Network

The network is essentially the key component of the deep learning design. For image classification, the network is typically a convolution neural network (i.e., CNN, ConvNet). CNN will take the full image as input and then predict the confidence level for each class, which is further used for classification. To learn a deep neural network, there are roughly two different ways, 1) training from scratch and 2) finetuning from a pre-trained model.

### Training from scratch

First, we will train a shallow convolutional neural network from scratch for practice. We first define the neural network in **GradBasicNet**, randomly initialize the value of the parameters in the network

(the reason why it is called "from scratch"), and train it on the training dataset:

1. Fill out **get_conv_layers()** with a network of **2 conv layers** each with kernel size of 5. After the first convolutional layers, add a **max pooling layer** with kernel size of 2 and stride of 2. Remember to add the non linearities immediately after the conv layers. You must choose the in-channels and out-channels for both the layers yourself: remember that the image to be input will be RGB so there is only one number that can be used for the in-channels of the first conv layer. Use the nn.Sequential API to combine all this into one layer, and return from **get_conv_layers()** method

2. Fill out **final_pool_layer()** with a max pooling layer. Typically after all the convolutional layers there is another max pooling layer. Use the same kernel size and stride as before and this time directly return the **nn.MaxPool2d**.

3. Fill out the **get_fc_layers()** method with a classifier containing 3 linear layers. Once again you are free to choose the in_channels and out_channels for these yourself. Once again use the **nn.Sequential** API and return the object from the method. Inside, alternate the Linear layers with ReLU activations. You do not need a ReLU after the final layer. Remember that the first Linear layer must take in the output of the final convolution layer so depending on your choice in (1.) there is only 1 value you can have for the in_channels of the first Linear layer. Also remember that the final Linear layer must have out_channels=10 as we are performing 10-way classification. Lastly, **remember to leave comments** on why you need to keep the relu of the first two layers while remove the relu of the last layer.

4. Finally, you should fill out the forward pass using these layers. Remember to use **self.conv_model**, **self.final_max_pool** and **self.fc_model** one after the other.

Here are a few useful reference for your implementation, which are useful for your implementation:

1. convolution layers: https://pytorch.org/docs/stable/generated/torch.nn.Conv2d.html
2. mas pooling layers: https://pytorch.org/docs/stable/generated/torch.nn.MaxPool2d.html
3. Sequential API: https://pytorch.org/docs/stable/generated/torch.nn.Sequential.html

```
 1 class GradBasicNet(nn.Module):
 2
 3     def __init__(self):
 4         super().__init__()
 5
 6         self.conv_model = self.get_conv_layers()
 7         self.final_max_pool = self.final_pool_layer()
 8         self.fc_model = self.get_fc_layers()
 9
10     def get_conv_layers(self):
11
12         #TODO: Group the convolutional layers using nn.Sequential
13
```

```
14          layers = nn.Sequential(
15              nn.Conv2d(3,16,5),
16              nn.ReLU(),
17              nn.MaxPool2d(2,stride=2),
18              nn.Conv2d(16,32,5),
19              nn.ReLU(),
20              #self.final_pool_layer()
21          )
22
23          return layers
24
25      def final_pool_layer(self):
26
27          #TODO Set this to a MaxPool layers
28
29          layer = nn.MaxPool2d(2,stride=2)
30
31          return layer
32
33      def get_fc_layers(self):
34
35          # TODO Group the linear layers using nn.Sequential
36
37          layers = nn.Sequential(
38
39              nn.Linear(800,256),
40              nn.ReLU(),
41              nn.Linear(256,64),
42              nn.ReLU(),
43              nn.Linear(64,10)
44          )
45
46          # ==================================================
47          # We do not need a ReLU after last layer because the last layer
48          # contains the logits of each class and denote the prob. with
49          # which we should select a class label. ReLU would change these values
50          # ==================================================
51          return layers
52
53      def register_grad_hook(self, grad):
54          self.grad = grad
55
56
57      def forward(self, x):
58
59          #TODO
60          x = self.conv_model(x)
61          #ignore this: relevance for gradcam section
62          h = x.register_hook(self.register_grad_hook)
63
64          x = self.final_max_pool(x)
```

```
65        x = torch.flatten(x, start_dim=1)
66        x = self.fc_model(x)
67
68        return x
69
70    def get_gradient_activations(self):
71        return self.grad
72
73    def get_final_conv_layer(self, x):
74        return self.conv_model(x)
```

## Finetuning on a pre-trained model

Secondly, we build a stronger convolutional neural network by starting from a pre-trained model **AlexNet**:

In this subsection, we might have to take advantage of a pretrained network in a process called transfer learning, where we only train a few final layers of a neural network. Here we will use the AlexNet architecture (Krizhevsky et al.)that revolutionized Deep Learning. The pretrained model (on ImageNet) is available on torchvision library and all we need to is ask PyTorch to allow updates on a few of thee final layers during training. We will put the AlexNet model also into the API we have used for our model allow, except that we will need an additionl **transition layer** function for the added **AvgPool** layer. Visualize the model by running the code cell. You will notice:

1. (features) contains most of the conv layers. We need upto (11) to include every conv layer
2. (features) (12) is the final max pooling layer
3. (avgpool) is the transition average pooling layer
4. (classifier) is the collection of linear layers

## Question

After you finish the implementation and run the experiments, pls go back and answer the questions below:

1. what is the main difference between finetuning and training from scratch.

2. Compare the performance between finetuning from a pre-trained model and training from scratch. Explain why the better one outperforms the other?

**Answer**

1. In fine-tunning the model was already pre-trained on some task. The model weights contain information required for that task. Information from other task can be utilized for new task ny updating the weights specific to new task. This way model is trained for new task and contains information learned from before. When training from scratch, the model does not

have any prior information and so the weights are set to either 0 or some value based on different normalization techniques. The model is trained specifically for the task and takes time to get same accuracy as compared to fine-tuning.

2. The pre-trained model performs better than the model trained from scratch. The pre-trained model can utilize information from ither tasks and apply that information to similar task or domain and thus gives better results.

```
1 example_model = models.alexnet(pretrained=True)
2 print(example_model)
```

Your task is two fold:

1. Implement the method **activate training layers** which sets the requires_grad of relevant parameters to True, so that training can occur. You can iterate over training parameters with

```
for name, param in self.conv_model.named_parameters():
```

and

```
for name, param in self.fc_model.named_parameters():
```

For the conv layers, every param should have requires_grad set to false except for the last layer (10).

For the linear layers, all layers must be trainable aka requires_grad must be set to True.

2. Implement the forward pass in the following order: self.conv_model, self.final_max_pool, self.avg_pool, self.fc_model

3. Remember to fill in your comments under the question mentioned above.

```
 1 class GradAlexNet(nn.Module):
 2
 3     def __init__(self):
 4         super().__init__()
 5
 6         self.base_alex_net = models.alexnet(pretrained=True)
 7
 8
 9         self.conv_model = self.get_conv_layers()
10         self.final_max_pool = self.final_pool_layer()
11         self.avg_pool = self.transition_layer()
12         self.fc_model = self.get_fc_layers()
13
14         self.activate_training_layers()
15
16     def activate_training_layers(self):
17
18         #TODO Fill out the function below
19         for name, param in self.conv_model.named_parameters():
20
21             print(name)
22             #this is the number of every convolutional layer. From what model print
23             #the last convolutional layer?
24             # ================================================
25             # Layer 10
26             # ================================================
27             number = int(name.split('.')[0])
28
29             # TODO: for all layers except the last layer set param.requires_grad =
30             if number < 10:
31                 param.requires_grad = False
32
33         for name, param in self.fc_model.named_parameters():
34
35             # for all of these layers set param.requires_grad as True
36             param.requires_grad = True
37
```

```
38
39      def get_conv_layers(self):
40
41          return self.base_alex_net.features[:12]
42
43      def final_pool_layer(self):
44          return nn.MaxPool2d(kernel_size=3, stride=2, padding=0, dilation=1, ceil_mc
45
46      def transition_layer(self):
47          return nn.AdaptiveAvgPool2d(output_size=(6, 6))
48
49      def get_fc_layers(self):
50          return nn.Sequential(
51              nn.Dropout(p=0.5, inplace=False),
52              nn.Linear(in_features=9216, out_features=4096, bias=True),
53              nn.ReLU(inplace=True),
54              nn.Dropout(p=0.5, inplace=False),
55              nn.Linear(in_features=4096, out_features=4096, bias=True),
56              nn.ReLU(inplace=True),
57              nn.Linear(in_features=4096, out_features=1000, bias=True),
58              nn.ReLU(inplace=True),
59              nn.Linear(in_features=1000, out_features=10, bias=True),
60          )
61
62      def register_grad_hook(self, grad):
63          self.grad = grad
64
65      def forward(self, x):
66
67          #TODO fill out the forward pass
68          x = self.conv_model(x)
69          h = x.register_hook(self.register_grad_hook)
70          x = self.final_max_pool(x)
71          x = self.avg_pool(x)
72          x = torch.flatten(x,start_dim=1)
73          x = self.fc_model(x)
74
75          return x
76
77
78      def get_gradient_activations(self):
79          return self.grad
80
81      def get_final_conv_layer(self, x):
82          return self.conv_model(x)
```

## Training the Model and Evaluating the Model's Performance

Performing network training is intuitive, i.e., feed the training data into the network and then use the groundtruth label to guide the training of the network parameters. In practice, we need to measure the inconsistency between the network output and groundtruth label and then update the network through gradient back propagation. As such, during a network training, a few components are necessary for the implementataion.

1. criterion (loss function): measure the inconsistency between the network output and label
2. optimizer: calculates gradients and then update the network parameters.

Similar to the human vision system which may need to obtain new knowledge by repeating the process of learning and practicing, purely feeding the training data into the model once is far from enough. Meanwhile, we also need to check whether the model overfits to the training data. As such, we monitor the status of the learned network by checking its performance on the test dataset at each training loop.

To note, one training loop is denoted as one epoch and means all of the training data has been used to train the deep neural network once. Now we will fill out the classifier we use to train these networks. Study the general framework for the code below well.

1. In the **init** method fill out the **self.criterion** and **self.optimizer**. Remember this is a classification problem so we will need a cross entropy loss. For the optimizer, I would recommend stochastic gradient descent with a learning rate of 0.001 and momentum of 0.9. The rest of **init** has been filled out for you.
2. Next, fill out the **training loop**. Here you are expected to iterate over **self.dataloaders['train']** and optimize on the loss with the groundtruth labels. The **validation** aspect of the training loop has been provided for you and so has the evaluate method. In the training loop, print the average loss every 1000 images you process. Remeber to zero out gradients on the model before you do loss.backward(), and then only after this backward step, make a step in the right direction using the optimizer.

At this stage ignore all methods, after evaluate. They will be relevant to later sections and you will have to return to them when you have more instructions.

```
1 class Classifier():
2
3     def __init__(self, name, model, dataloaders, class_names, use_cuda=False):
4
5         '''
6         @name: Experiment name. Will define stored results etc.
7         @model: Either a GradBasicNet() or a GradAlexNet()
8         @dataloaders: Dictionary with keys train, val and test and corresponding da
9         @class_names: list of classes, where the idx of class name corresponds to t
10         @use_cuda: whether or not to use cuda
11         '''
12
```

```
13          self.name = name
14          if use_cuda and not torch.cuda.is_available():
15              raise Exception("Asked for CUDA but GPU not found")
16
17          self.use_cuda = use_cuda
18
19          self.model = model.to('cuda' if use_cuda else 'cpu')
20
21          #TODO
22          self.criterion = nn.CrossEntropyLoss()
23          self.optim = optim.SGD(self.model.parameters(), lr=0.001, momentum=0.9)
24
25          self.dataloaders = dataloaders
26          self.class_names = class_names
27          self.activations_path = os.path.join('activations', self.name)
28          self.kernel_path = os.path.join('kernel_viz', self.name)
29
30          save_path = os.path.join(os.getcwd(), 'models', self.name)
31          if not os.path.exists(save_path):
32              os.makedirs(save_path)
33
34          if not os.path.exists(self.activations_path):
35              os.makedirs(self.activations_path)
36
37          if not os.path.exists(self.kernel_path):
38              os.makedirs(self.kernel_path)
39
40          self.save_path = save_path
41
42      def train(self, epochs, save=True):
43          '''
44          @epochs: number of epochs to train
45          @save: whether or not to save the checkpoints
46          '''
47
48          best_val_accuracy = - math.inf
49
50          for epoch in range(epochs):
51              print('Epoch', epoch+1)
52              self.model.train()
53
54              batches_in_pass = len(self.dataloaders['train'])
55
56              #You may comment these two lines if you do not wish to use them
57              loss_total = 0.0 # Record the total loss within a few steps
58              epoch_loss = 0.0 # Record the total loss for each epoch
59
60              # TODO Iterate over the training dataloader (see how it is done for val
61              # to call the optim.zero_grad(), loss.backward() and optim.step()
62
63              counter = 0 # variable to count no of images processed
```

```
64                    for x, y in self.dataloaders['train']:
65                      self.optim.zero_grad()
66
67                      # Print training loss after processsing every 1000 images
68                      counter += len(x)
69                      if counter%10000 == 0:
70                        print('Training loss after every 10000 images', epoch_loss/(counter
71
72                      x = x.to('cuda' if self.use_cuda else 'cpu')
73                      y = y.to('cuda' if self.use_cuda else 'cpu')
74
75                      output = self.model(x)
76                      loss = self.criterion(output, y)
77                      epoch_loss += loss.cpu().detach().numpy()#/len(x)
78                      loss.backward()
79                      self.optim.step()
80
81
82
83              '''Give validation'''
84              epoch_loss /= batches_in_pass
85
86              self.model.eval()
87
88              #DO NOT modify this part
89              correct = 0.0
90              total = 0.0
91              for idx, data in enumerate(self.dataloaders['val']):
92
93                  inputs, labels = data
94                  inputs = inputs.to('cuda' if self.use_cuda else 'cpu')
95                  labels = labels.to('cuda' if self.use_cuda else 'cpu')
96
97                  outputs = self.model(inputs)
98                  _, predicted = torch.max(outputs, 1)
99
100                  total += labels.shape[0]
101                  correct += (predicted == labels).sum().item()
102
103              epoch_accuracy = 100 * correct / total
104
105              print(f'Train Epoch Loss (Avg): {epoch_loss}')
106              print(f'Validation Epoch Accuracy:{epoch_accuracy}')
107
108              if save:
109                  #  Make sure that your saving pipeline is working well.
110                  # Is os library working on your file system?
111                  # Is your model being saved and reloaded fine?
112                  # When you do the kernel viz, activation maps,
113                  # and GradCAM you must be using the model you have saved before.
114
```

```
115                   torch.save(self.model.state_dict(), os.path.join(self.save_path, f'
116
117               if epoch_accuracy > best_val_accuracy:
118
119                   torch.save(self.model.state_dict(), os.path.join(self.save_path
120                   best_val_accuracy = epoch_accuracy
121
122       print('Done training!')
123
124
125   def evaluate(self):
126
127       try:
128           assert os.path.exists(os.path.join(self.save_path, 'best.pt'))
129
130       except:
131           print('It appears you are testing the model without training. Please t
132           return
133
134       self.model.load_state_dict(torch.load(os.path.join(self.save_path, 'best.pt
135       self.model.eval()
136
137       #total = len(self.dataloaders['test'])
138
139       correct = 0.0
140       total = 0.0
141       for idx, data in enumerate(self.dataloaders['test']):
142
143               inputs, labels = data
144               inputs = inputs.to('cuda' if self.use_cuda else 'cpu')
145               labels = labels.to('cuda' if self.use_cuda else 'cpu')
146
147               outputs = self.model(inputs)
148               _, predicted = torch.max(outputs, 1)
149
150               total += labels.shape[0]
151               correct += (predicted == labels).sum().item()
152
153       print(f'Accuracy: {100 * correct/total}%')
154
155   def grad_cam_on_input(self, img):
156
157       try:
158           assert os.path.exists(os.path.join(self.save_path, 'best.pt'))
159
160       except:
161           print('It appears you are testing the model without training. Please t
162           return
163
164       self.model.load_state_dict(torch.load(os.path.join(self.save_path, 'best.pt
165
```

```
166
167            self.model.eval()
168            img = img.to('cuda' if self.use_cuda else 'cpu')
169
170
171            out = self.model(img)
172
173            _, pred = torch.max(out, 1)
174
175            predicted_class = self.class_names[int(pred)]
176            print(f'Predicted class was {predicted_class}')
177
178            out[:, pred].backward()
179            gradients = self.model.get_gradient_activations()
180
181            print('Gradients shape: ', f'{gradients.shape}')
182
183            mean_gradients = torch.mean(gradients, [0, 2, 3]).cpu()
184            activations = self.model.get_final_conv_layer(img).detach().cpu()
185
186            print('Activations shape: ', f'{activations.shape}')
187
188            for idx in range(activations.shape[1]):
189                activations[:, idx, :, :] *= mean_gradients[idx]
190
191            final_heatmap = np.maximum(torch.mean(activations, dim=1).squeeze(), 0)
192
193            final_heatmap /= torch.max(final_heatmap)
194
195            return final_heatmap
196
197
198
199
200     def trained_kernel_viz(self):
201
202            all_layers = [0, 3]
203            all_filters = []
204            for layer in all_layers:
205
206                #TODO: blank out first line
207                filters = self.model.conv_model[layer].weight
208                all_filters.append(filters.detach().cpu().clone()[:8, :8, :, :])
209
210            for filter_idx in range(len(all_filters)):
211
212                filter = all_filters[filter_idx]
213                print(filter.shape)
214                filter = filter.contiguous().view(-1, 1, filter.shape[2], filter.shape|
215                image = show_img(make_grid(filter))
216                image = 255 * image
```

```
217                  cv2.imwrite(os.path.join(self.kernel_path, f'filter_layer{all_layers[fi
218
219
220     def activations_on_input(self, img):
221
222          img = img.to('cuda' if self.use_cuda else 'cpu')
223
224          all_layers = [0,3,6,8,10]
225          all_viz = []
226
227          for each in all_layers:
228
229               current_model = self.model.conv_model[:each+1]
230               current_out = current_model(img)
231               all_viz.append(current_out.detach().cpu().clone()[:, :64, :, :])
232
233          for viz_idx in range(len(all_viz)):
234
235               viz = all_viz[viz_idx]
236               viz = viz.view(-1, 1, viz.shape[2], viz.shape[3])
237               image = show_img(make_grid(viz))
238               image = 255 * image
239               cv2.imwrite(os.path.join(self.activations_path, f'sample_layer{all_laye
```

Run the classifier for the code using the basic model by running the following snippets. If all goes well, you should have a test accuracy of about ~60-70% at the end of it. It should take <10 mins to run on a GPU.

```
 1 experiment_name = 'basic_debug'  #Provide name to model experiment
 2 model_name = 'basic' #Choose between [basic, alexnet]
 3 batch_size = 5  #You may not need to change this but incase you do
 4 first_run = True #whether or not first time running it
 5
 6 trainloader, valloader, testloader, class_names = get_dataloader(batch_size=batch_s
 7 dataloaders = {'train': trainloader, 'val' : valloader, 'test': testloader, 'mappir
 8
 9 if model_name == 'basic':
10     model = GradBasicNet()
11 elif model_name == 'alexnet':
12     model = GradAlexNet()
13 else:
14     raise NotImplementedError("This option has not been implemented. Choose between
15
16 classifier = Classifier(experiment_name, model, dataloaders, class_names, use_cuda=
```

```
 1 # When you develop your code, to save your time, you can choose to run the model
 2 # for 5 cpoches. The accuracy after training for 5 epoches has already been high
 3 # and close to the model after 20-epoch traiing.
 4 # (In my case, Validation Epoch Accuracy is above 61 after training 5 epoches)
 5
 6 # To note, for your final submission, you should make sure to train the model
 7 # for 20 epoches and analysis that model in the later sections.
 8
 9 # classifier.train(epochs=5) # For your reference
10 classifier.train(epochs=20)
11 classifier.evaluate()
```

```
Epoch 1
Training loss after every 10000 images 2.108399686694145
Training loss after every 10000 images 1.9231158244758844
Training loss after every 10000 images 1.7977678915659587
Training loss after every 10000 images 1.7106661086902022
Train Epoch Loss (Avg): 1.6749626362025738
Validation Epoch Accuracy:49.8
Epoch 2
Training loss after every 10000 images 1.3291555665433408
Training loss after every 10000 images 1.296781690724194
Training loss after every 10000 images 1.2748593087643385
Training loss after every 10000 images 1.2552443586792796
Train Epoch Loss (Avg): 1.2429706958399878
Validation Epoch Accuracy:59.12
Epoch 3
Training loss after every 10000 images 1.0712261559907348
Training loss after every 10000 images 1.058656976292841
Training loss after every 10000 images 1.0512348008410384
Training loss after every 10000 images 1.0415475951433182
Train Epoch Loss (Avg): 1.0386115505840214
Validation Epoch Accuracy:62.68
Epoch 4
Training loss after every 10000 images 0.89866407478787
Training loss after every 10000 images 0.9052423220551573
Training loss after every 10000 images 0.9070602214358127
Training loss after every 10000 images 0.9040456310610753
Train Epoch Loss (Avg): 0.8998311887752886
Validation Epoch Accuracy:65.44
Epoch 5
Training loss after every 10000 images 0.778684928946197
Training loss after every 10000 images 0.7730257469220086
Training loss after every 10000 images 0.7793479903711317
Training loss after every 10000 images 0.7871425542671932
Train Epoch Loss (Avg): 0.7896533831039236
Validation Epoch Accuracy:69.4
Epoch 6
Training loss after every 10000 images 0.6645055955364368
Training loss after every 10000 images 0.6779277559234761
Training loss after every 10000 images 0.6881485834918761
Training loss after every 10000 images 0.6970168631075067
Train Epoch Loss (Avg): 0.6997378139392887
```

```
Validation Epoch Accuracy:67.52
Epoch 7
Training loss after every 10000 images 0.5708401405599434
Training loss after every 10000 images 0.5853927557197166
Training loss after every 10000 images 0.5991539110465286
Training loss after every 10000 images 0.6090084843362856
Train Epoch Loss (Avg): 0.6134476623121882
Validation Epoch Accuracy:69.82
Epoch 8
Training loss after every 10000 images 0.48343467791966394
Training loss after every 10000 images 0.5050446154602832
Training loss after every 10000 images 0.5172256239859562
Training loss after every 10000 images 0.5321084329626683
Train Epoch Loss (Avg): 0.5375097855089116
Validation Epoch Accuracy:69.08
Epoch 9
Training loss after every 10000 images 0.4272262834580615
```

Now run the classifier for the code using the alexnet model specified above. You should notice a notable performance increase. On a GPU, this trained for about 30 mins.

```
 1 experiment_name = 'alexnet_debug'  #Provide name to model experiment
 2 model_name = 'alexnet' #Choose between [basic, alexnet]
 3 batch_size = 5  #You may not need to change this but incase you do
 4 first_run = True #whether or not first time running it
 5
 6 trainloader, valloader, testloader, class_names = get_dataloader(batch_size=batch_s
 7 dataloaders = {'train': trainloader, 'val' : valloader, 'test': testloader, 'mappir
 8
 9 #model = models.alexnet(pretrained=True)
10 if model_name == 'basic':
11
12     model = GradBasicNet()
13
14 elif model_name == 'alexnet':
15
16     model = GradAlexNet()
17
18 else:
19     raise NotImplementedError("This option has not been implemented. Choose between
20
21 classifier = Classifier(experiment_name, model, dataloaders, class_names, use_cuda=
```

```
Files already downloaded and verified
Files already downloaded and verified
0.weight
0.bias
3.weight
3.bias
6.weight
6.bias
8.weight
```

```
      8.bias
      10.weight
      10.bias
```

```python
 1 # When you develop your code, to save your time, you can choose to run the model
 2 # for 3 cpoches. The accuracy after training for 3 epoches has already been high
 3 # and close to the model after 20-epoch traiing.
 4
 5 # To note, for your final submission, it is recommended to run the model for 20 epo
 6 # However, if it takes time, you should at least run the model for 5 epochs.
 7
 8 # classifier.train(epochs=3) # For your reference
 9 classifier.train(epochs=20)
10 classifier.evaluate()
```

```
    Epoch 1
    Training loss after every 10000 images 1.2949434548374266
    Training loss after every 10000 images 1.0529366769925692
    Training loss after every 10000 images 0.9353411313917023
    Training loss after every 10000 images 0.8654944244204962
    Train Epoch Loss (Avg): 0.8405958373999586
    Validation Epoch Accuracy:81.86
    Epoch 2
    Training loss after every 10000 images 0.5345164470453747
    Training loss after every 10000 images 0.5325533928129589
    Training loss after every 10000 images 0.5316038423580273
    Training loss after every 10000 images 0.5265015859799678
    Train Epoch Loss (Avg): 0.5236099226783522
    Validation Epoch Accuracy:84.02
    Epoch 3
    Training loss after every 10000 images 0.3926609784188622
    Training loss after every 10000 images 0.40889036247255717
    Training loss after every 10000 images 0.41188149072230107
    Training loss after every 10000 images 0.4115336829901371
    Train Epoch Loss (Avg): 0.4143479204564347
    Validation Epoch Accuracy:85.24
    Epoch 4
    Training loss after every 10000 images 0.31610381750648686
    Training loss after every 10000 images 0.32527802798452193
    Training loss after every 10000 images 0.33471010669763807
    Training loss after every 10000 images 0.3401378570418028
    Train Epoch Loss (Avg): 0.34142804408026134
    Validation Epoch Accuracy:86.34
    Epoch 5
    Training loss after every 10000 images 0.2554356923725884
    Training loss after every 10000 images 0.2656640505340719
    Training loss after every 10000 images 0.2747126869445243
    Training loss after every 10000 images 0.2775496001906796
    Train Epoch Loss (Avg): 0.2799376558050823
    Validation Epoch Accuracy:85.76
    Epoch 6
    Training loss after every 10000 images 0.2158866383147248
    Training loss after every 10000 images 0.22701751376855492
    Training loss after every 10000 images 0.2287121267898962
```

```
    Training loss after every 10000 images 0.23539165419614402
    Train Epoch Loss (Avg): 0.23716592982548101
    Validation Epoch Accuracy:87.02
    Epoch 7
    Training loss after every 10000 images 0.17729332101917794
    Training loss after every 10000 images 0.18280629259671652
    Training loss after every 10000 images 0.19122098856431965
    Training loss after every 10000 images 0.19816327211533166
    Train Epoch Loss (Avg): 0.19938465434015457
    Validation Epoch Accuracy:87.66
    Epoch 8
    Training loss after every 10000 images 0.14328093567377403
    Training loss after every 10000 images 0.1548716495381098
    Training loss after every 10000 images 0.15861636300978707
    Training loss after every 10000 images 0.16520855326191847
    Train Epoch Loss (Avg): 0.16896901944210066
    Validation Epoch Accuracy:86.74
    Epoch 9
    Training loss after every 10000 images 0 13405209202449264
```

**Before you move forward to the next step**: Try and see what classes your model does well on (you can modify the testing code for this). This will help you pick the best visualization to show later.

## Understand the Deep Neural Networks

The deep learning is quite powerful and can achieve great performance after training for a few epochs. However, we still don't know why it works. Thus, we first study activation map to analysis the region that triggers the final classification. Then, we visualze a few kernels to infer what is learned in the neural network. (The activtion map and kernel visualization are two popurlar directions. However, much more effort is still needed to comprehensively understand deep neural networks and reach the ultimate goal.)

### Activation Map

For interpreting our alexnet model, we will be using a simple version of Grad-CAM (Gradient based Class activation mapping) by Selvaraju et al. (https://arxiv.org/abs/1610.02391). This will help us see what region of the input image the output is 'focusing on' while making its key choice. I recommend reading the paper for those interested: however, the following instructions should also suffice. If you go step-by-step then all information is given in the intructions. You will receive one img to the method.

1. First, move the img to cuda if you are using GPU.
2. The code to load the model trained from before has been provided to you. Use self.model to output the predictions to the variable out. Your output should have dimension (1, 10).
3. The predicted class is the index of the highest value of of these 10 values. Use **torch.max** along dim 1 to get the argument of the max value. This method will return two values, the

latter of them is the required argument. The next two lines have been provided to you: they indicate the predicted loss

4. Call the **backward** method on out[:, pred]. Previously during the forward passes, we have applied a gradient hook on the last convolutional layer. This will mean that during the applied backward pass, we will record the value of the gradient for the **maximum predicted class** with respect to the **output** of the last convolutional layer. This will be the same size as the **output** of the last convolutional layer. Do you see why?

5. After the **backward** call, get the value of the gradient above using the **get_gradient_activations** function on **self.model** and store it in gradients. Now use the torch.mean method to get the mean value of gradients across all dimensionss except the channels dimension (number of filters) and store this in **mean_gradients**. Your output should be of shape (1, 64, 1, 1). If your tensors have been on GPU, you should move them to CPU using .detach().cpu()

6. Use the self.model.get_final_conv_layer(img) to store the activations at the final layer to activations. This should be of shape (1, 64, H, W).

7. Now for each of the 64 filters, **scale** the activations at that filter, with the corresponding **mean_gradients** value for that filter. Your output should have size (1, 64, H, W) Iterate over the 64 filters in the following way: `for idx in range(activations.shape[1]):`

8. As a final step, we will take mean across the 64 filters and do a ReLU to get rid of negative activations before normalzing one last time to get the heatmap. This has been done for you.

```python
1  '''Sample an image from the test set'''
2
3  #You may change the sampling code to sample an image as you desire.
4  #Make sure to NOT move the sampling code to a different cell.
5
6  img_batch, labels_batch = next(iter(testloader))
7  img = img_batch[3]
8  img = img.unsqueeze(0)
9
10 classifier = Classifier(experiment_name, model, dataloaders, class_names, use_cuda=
11 heatmap = classifier.grad_cam_on_input(img)
12
13 def visualize(img, heatmap):
14
15
16    heatmap = heatmap.cpu().numpy()
17
18    img = show_img(img)
19    img = np.uint8(255 * img)
20    img = cv2.cvtColor(img, cv2.COLOR_BGR2RGB)
21    print(img.shape)
22    print(heatmap.shape)
```

```
23      heatmap = cv2.resize(heatmap, (img.shape[1], img.shape[0]))
24      heatmap = np.uint8(255 * heatmap)
25      heatmap = cv2.applyColorMap(heatmap, cv2.COLORMAP_JET)
26      heatmap = cv2.cvtColor(heatmap, cv2.COLOR_BGR2RGB)
27
28      combine = 0.5 * heatmap + img
29
30      #if not os.path.exists(write_path):
31      #    os.makedirs(write_path)
32
33      plt.imshow(combine/255)
34      plt.show()
35
36
37 visualize(img, heatmap)
```

What do you observe? Show an example of an image where the method is looking at the object in question and another where it appear to be completely unrelated. In the latter case, it might have learnt a spurious correlation- aka a bias in the data which always appears to be correlated with a given label. For the ship class, this **might** be the surrounding water or for a **horse** it might be the surrounding grass. In such cases, do you think the model would predict correctly for a ship on sand or a horse in the air? Answer in a text snippet below. **Causally trained neural networks** (https://www.cmu.edu/dietrich/causality/neurips20ws/) are an exciting direction to solve this problem.

```
1 # Your code here to show an failure case.
2 # You can refer the steps and functions implemented in the previous cell and reuse
3 img_batch, labels_batch = next(iter(testloader))
4 img = img_batch[4]
5 img = img.unsqueeze(0)
6
7 heatmap = classifier.grad_cam_on_input(img)
8 visualize(img, heatmap)
```

```
1 # Your code here to show an failure case.
2 # You can refer the steps and functions implemented in the previous cell and reuse
3 img_batch, labels_batch = next(iter(testloader))
4 img = img_batch[1]
5 img = img.unsqueeze(0)
6
7 heatmap = classifier.grad_cam_on_input(img)
8 visualize(img, heatmap)
```

**Answer**: In the first example, the model is looking at the object which is the frog. In the second example, the model is looking at the surroundings of the object i.e. water near the ship to identify the ship.

It's possible that the model might fail to identify a ship on sand or horse in air. Let's say for a ship on sand, as method is focussing on nearby surroundings, the sand has a different colour and visible paatern than water. However, it is possible that it might be able to identify a ship if it was trained on black and white image and was tested on the same. For the example, horse in air, the grass has a different look and colour than air, so the method has high chances of not being able to identify it. Assumption about Black and white training is not the same as in previous case.

## Kernel and Activation Visualizations

We will visualize some learned convolutional kernels for two layers in the conv-net. Study the code provided for **trained_kernel_viz** carefully. You only have to fill out the line for **filter**. All you are expected to do is to access the relevantt layer from **self.conv_model** and set **filter** equal to its weight parameter.

Call this function on the alexnet classifier. What do you observe?

```
1 classifier.trained_kernel_viz()
```

Now with the kernel viz filled out, write the method for activation visualizations **activations on input**. The structure for the code is very similar to the kernel viz, except that we are actually viewing the

output of the model at each stage and not for the kernel at that stage. Once filled, please call this method on a few sample images. What do you observe? Answer generally in a text snippet below.

**Answer**: We observe from the first image that the lower kernel/layer picks up on low level features like edges and lines while the other learns high level features like shapes.

```
1 '''Sample an image from the test set'''
2 #You may change the sampling code to sample an image as you desire.
3 #Make sure to NOT move the sampling code to a different cell.
4 img_batch, labels_batch = next(iter(testloader))
5 img = img_batch[3]
6 img = img.unsqueeze(0)
7
8 classifier = Classifier(experiment_name, model, dataloaders, class_names, use_cuda=
9 classifier.activations_on_input(img)
```

What do you observe about early layers v. later layers? Answer in a text snippet below.

**Answer**: The observation is same as before. First image contains mostly information about edges whereas the other capture more features like shapes

1