

▼ Assignment 3 Part 3

```
1 import os
2 import numpy as np
3 import torch
4 import torch.nn as nn
5 import torchvision
6
7 from torchvision import transforms, datasets
8 # from PIL import Image, ImageDraw
9 import matplotlib.pyplot as plt
10 from torchvision.datasets import ImageFolder
11 from torch.utils.data import DataLoader
12 import torchvision.models as models
13 import torch.nn.functional as F
14
15 %matplotlib inline
16 %load_ext autoreload
17 %autoreload 2
```

```
1 from gan.utils import sample_noise, show_images, deprocess_img, preprocess_img
```

▼ Custom Loss Functions

From [Wanglei](#)

```
1 def linear_combination(x, y, epsilon):
2     return epsilon * x + (1 - epsilon) * y
3
4
5 def reduce_loss(loss, reduction='mean'):
6     return loss.mean() if reduction == 'mean' else loss.sum() if reduction == 'sum' else
7
8
9 class LabelSmoothingCrossEntropy(nn.Module):
10     def __init__(self, epsilon: float = 0.1, reduction='mean'):
11         super().__init__()
12         self.epsilon = epsilon
13         self.reduction = reduction
14
15     def forward(self, preds, target):
16         n = preds.size()[-1]
17         log_preds = F.log_softmax(preds, dim=-1)
18         loss = reduce_loss(-log_preds.sum(dim=-1), self.reduction)
19         nll = F.nll_loss(log_preds, target, reduction=self.reduction)
20         return linear_combination(loss / n, nll, self.epsilon)
21
```

 0s completed at 8:05 PM 

In the following cell we will load the training data and also apply some transforms to the data.

```
1 # get the imagenette dataset
2 !wget https://s3.amazonaws.com/fast-ai-imageclas/imagenette2.tgz
3 !tar -xzf imagenette2.tgz
4 !rm imagenette2.tgz
```

Streaming output truncated to the last 5000 lines.

```
imagenette2/train/n03888257/n03888257_16077.JPEG
imagenette2/train/n03888257/n03888257_23339.JPEG
imagenette2/train/n03888257/n03888257_44204.JPEG
imagenette2/train/n03888257/n03888257_61633.JPEG
imagenette2/train/n03888257/n03888257_15067.JPEG
imagenette2/train/n03888257/n03888257_75365.JPEG
imagenette2/train/n03888257/n03888257_63966.JPEG
imagenette2/train/n03888257/n03888257_3927.JPEG
imagenette2/train/n03888257/n03888257_20684.JPEG
imagenette2/train/n03888257/ILSVRC2012_val_00047778.JPEG
imagenette2/train/n03888257/n03888257_14016.JPEG
imagenette2/train/n03888257/n03888257_37776.JPEG
imagenette2/train/n03888257/ILSVRC2012_val_00041706.JPEG
imagenette2/train/n03888257/n03888257_17513.JPEG
imagenette2/train/n03888257/n03888257_17143.JPEG
imagenette2/train/n03888257/n03888257_6738.JPEG
imagenette2/train/n03888257/n03888257_4355.JPEG
imagenette2/train/n03888257/ILSVRC2012_val_00030583.JPEG
imagenette2/train/n03888257/n03888257_18127.JPEG
imagenette2/train/n03888257/n03888257_10844.JPEG
imagenette2/train/n03888257/n03888257_6503.JPEG
imagenette2/train/n03888257/n03888257_12195.JPEG
imagenette2/train/n03888257/n03888257_14368.JPEG
imagenette2/train/n03888257/n03888257_11886.JPEG
imagenette2/train/n03888257/n03888257_5293.JPEG
imagenette2/train/n03888257/n03888257_10014.JPEG
imagenette2/train/n03888257/n03888257_30158.JPEG
imagenette2/train/n03888257/n03888257_11255.JPEG
imagenette2/train/n03888257/n03888257_38929.JPEG
imagenette2/train/n03888257/n03888257_10106.JPEG
imagenette2/train/n03888257/n03888257_36937.JPEG
imagenette2/train/n03888257/n03888257_128.JPEG
imagenette2/train/n03888257/n03888257_23723.JPEG
imagenette2/train/n03888257/n03888257_18023.JPEG
imagenette2/train/n03888257/n03888257_25024.JPEG
imagenette2/train/n03888257/n03888257_8225.JPEG
imagenette2/train/n03888257/n03888257_20676.JPEG
imagenette2/train/n03888257/n03888257_4154.JPEG
imagenette2/train/n03888257/n03888257_18664.JPEG
imagenette2/train/n03888257/n03888257_29569.JPEG
imagenette2/train/n03888257/n03888257_21159.JPEG
imagenette2/train/n03888257/ILSVRC2012_val_00022496.JPEG
imagenette2/train/n03888257/n03888257_15797.JPEG
imagenette2/train/n03888257/n03888257_75495.JPEG
imagenette2/train/n03888257/n03888257_4553.JPEG
imagenette2/train/n03888257/n03888257_24025.JPEG
```

```
imagenette2/train/n03888257/n03888257_10653.JPEG
imagenette2/train/n03888257/n03888257_50108.JPEG
imagenette2/train/n03888257/n03888257_10079.JPEG
imagenette2/train/n03888257/n03888257_8423.JPEG
imagenette2/train/n03888257/n03888257_41099.JPEG
imagenette2/train/n03888257/n03888257_27166.JPEG
imagenette2/train/n03888257/n03888257_10645.JPEG
```

```
1 nb_epochs = 45
2 batch_size = 64
3 learning_rate = 0.001
4 scale_size = 256
5
6 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
7
8 normalize = transforms.Normalize(mean=[0.485, 0.456, 0.406],
9                                   std= [0.229, 0.224, 0.225])
10 data_dir = './imagenette2'
11 # print(os.listdir(data_dir))
12 train_dataset = ImageFolder(root=data_dir+'/train',
13                             transform=transforms.Compose([
14                                 transforms.Resize((scale_size, scale_size)),
15                                 transforms.CenterCrop(256),
16                                 transforms.ToTensor(),
17                                 normalize,
18                             ]))
19 val_dataset = ImageFolder(root=data_dir+'/val',
20                           transform=transforms.Compose([
21                               transforms.Resize((scale_size, scale_size)),
22                               transforms.CenterCrop(256),
23                               transforms.ToTensor(),
24                               normalize,
25                           ]))
26
27 train_loader_im = DataLoader(train_dataset, batch_size=batch_size, shuffle=True, num_wo
28 val_loader_im = DataLoader(val_dataset , batch_size=batch_size, num_workers=1, pin_memo
29
30 imgs = train_loader_im.__iter__().next()[0].numpy().squeeze()
```

Pretrained Models

In this project, I used the pretrained AlexNet model

```
1 num_epochs = 4
2 test_frequency = 5
3
4 # Load Pretrained AlexNet
5 classifier = torchvision.models.alexnet(pretrained=True)
6 # set a new output size
```

Downloading: "<https://download.pytorch.org/models/alexnet-owt-7be5be79.pth>" to /root/

100%

233M/233M [00:04<00:00, 55.7MB/s]

Training the network

```
1 def train_classifier(train_loader, classifier, criterion, optimizer):
2     classifier.train()
3     loss_ = 0.0
4     losses = []
5     for i, (images, labels) in enumerate(train_loader):
6         images, labels = images.to(device), labels.to(device)
7         optimizer.zero_grad()
8         logits = classifier(images)
9         loss = criterion(logits, labels)
10        loss.backward()
11        optimizer.step()
12        losses.append(loss)
13    return torch.stack(losses).mean().item()

1 def test_classifier(test_loader, classifier, criterion, print_ind_classes=True, print
2     classifier.eval()
3     losses = []
4     correct = 0
5     total = 0
6     with torch.no_grad():
7         for i, (images, labels) in enumerate(test_loader):
8             images, labels = images.to(device), labels.to(device)
9             logits = classifier(images)
10            # loss = F.nll_loss(logits, labels)
11            loss = criterion(logits, labels)
12            losses.append(loss.item())
13            # ignore first class which is background
14            preds = logits.cpu().detach().numpy()
15            correct += (np.argmax(preds, axis=1) == labels.cpu().detach().numpy()).sum()
16            total += test_loader.batch_size
17
18        mAP = float(correct) / float(total)
19        test_loss = np.mean(losses)
20        if print_total:
21            print('Acc: {0:.4f}'.format(mAP))
22            print('Avg loss: {}'.format(test_loss))
23
24    return mAP, test_loss, 0
```

```

8     print("Starting epoch number " + str(epoch))
9     train_loss = train_classifier(train_loader, classifier, criterion, optimizer)
10    train_losses.append(train_loss)
11    print("Loss for Training on Epoch " +str(epoch) + " is "+ str(train_loss))
12    if(epoch%test_frequency==0 or epoch==1):
13        mAP_train, _, _ = test_classifier(train_loader, classifier, criterion, False)
14        train_mAPs.append(mAP_train)
15        mAP_val, val_loss, _ = test_classifier(val_loader, classifier, criterion)
16        print('Evaluating classifier')
17        print("Mean Precision Score for Testing on Epoch " +str(epoch) + " is "+ str(val_loss))
18        val_losses.append(val_loss)
19        val_mAPs.append(mAP_val)
20
21    return classifier, train_losses, val_losses, train_mAPs, val_mAPs

```

```

1 classifier, train_losses, val_losses, train_mAPs, val_mAPs = train(classifier, num_epochs, train_loader, val_loader, criterion, optimizer)

```

```

Starting epoch number 1
/usr/local/lib/python3.7/dist-packages/torch/nn/functional.py:718: UserWarning: Named tensors are not supported by autograd in this mode.
  return torch.max_pool2d(input, kernel_size, stride, padding, dilation, ceil_mode)
Loss for Training on Epoch 1 is 53.87113952636719
Acc: 0.9100
Avg loss: 46.27601300516436
Evaluating classifier
Mean Precision Score for Testing on Epoch 1 is 0.9100302419354839
Starting epoch number 2
Loss for Training on Epoch 2 is 40.755882263183594
Starting epoch number 3
Loss for Training on Epoch 3 is 36.26933288574219
Starting epoch number 4
Loss for Training on Epoch 4 is 34.52720260620117

```

```

1 # Test on clean data
2 mAP_test, test_loss, test_aps = test_classifier(val_loader_im, classifier, criterion)
3 print("Test accuracy: ", mAP_test)

Acc: 0.9378
Avg loss: 40.88100867117605
Test accuracy: 0.9377520161290323

```

Advererial Attack

```

1 # FGSM attack code
2 def fgsm_attack(image, epsilon, data_grad):

```

```
1 # Iterative projected gradient Method
2 def igm_attack(model, X, y, thres, alpha, num_iter):
3     """ Construct igm adversarial examples on the examples X"""
4     delta = torch.zeros_like(X, requires_grad=True)
5     for t in range(num_iter):
6         loss = F.nll_loss(model(X + delta), y)
7         # loss = nn.CrossEntropyLoss()(model(X + delta), y)
8         loss.backward()
9         delta.data = (delta + alpha*delta.grad.detach().sign()).clamp(-thres, thres)
10        delta.grad.zero_()
11    return delta.detach() + X
12

1 #Evaluate results on adversarially perturbed
2 def eval_advAttack(model=None, test_loader=None, thres=0.1, num_iter=40):
3     total = 0
4     correct = 0
5     print("Evaluating single model results on adv data")
6     for data, target in test_loader:
7         if torch.cuda.is_available():
8             data, target = data.cuda(), target.cuda()
9             # print('new data -----')
10            # disp_fake_images = deprocess_img(data.data) # denormalize
11            # imgs_numpy = (disp_fake_images).cpu().numpy()
12            # show_images(imgs_numpy[0:4], color=True)
13            # plt.show()
14            # print()
15
16            # Call IGM Attack
17            perturbed_data = igm_attack(model, data, target, thres=thres, alpha=1e-3, num_ite
18
19            # disp_fake_images = deprocess_img(perturbed_data.data) # denormalize
20            # imgs_numpy = (disp_fake_images).cpu().numpy()
21            # show_images(imgs_numpy[0:4], color=True)
22            # plt.show()
23            # print()
24
25            # Re-classify the perturbed image
26            output = model(perturbed_data)
27            # Check for success
28            preds_np = output.cpu().detach().numpy()
29            # finalPred = np.argmax(preds_np, axis=1)
30            correct += (np.argmax(preds_np, axis=1) == target.cpu().detach().numpy()).sum()
31            # correct += (finalPred == target.cpu().detach().numpy()).sum()
32            total += test_loader.batch_size
```

Defense Mechanism

Firstly, we will use training the model with the adversarial data

```
1  #Adversarial Training
2  def adv_train(model, train_loader, criterion, thres=0.1, num_iter=40, nb_epochs=nb_ep
3      batch_size=batch_size, train_end=-1, test_end=-1, learning_rate=learn
4
5      optimizer = torch.optim.Adam(model.parameters(), lr=learning_rate)
6      train_loss = []
7      total = 0
8      correct = 0
9      totalAdv = 0
10     correctAdv = 0
11     step = 0
12     # breakstep = 0
13     for _epoch in range(nb_epochs):
14         for data, target in train_loader:
15             #Normal Training
16             if torch.cuda.is_available():
17                 data1, target = data.cuda(), target.cuda()
18
19             optimizer.zero_grad()
20             preds = model(data1)
21             # Calculate the loss
22             loss = criterion(preds, target1)
23             # loss = F.nll_loss(preds, target)
24             model.zero_grad()
25             loss.backward() # calc gradients
26
27             train_loss.append(loss.data.item())
28             optimizer.step() # update gradients
29
30             preds_np = preds.cpu().detach().numpy()
31             correct += (np.argmax(preds_np, axis=1) == target.cpu().detach().numpy()).sum
32             total += train_loader.batch_size
33
34             #Adversarial Training
35             # Call IGM Attack
36             data_adv = igm_attack(model, data, target, thres=thres, alpha=1e-3, num_iter=
37
38             if torch.cuda.is_available():
```

```
52         if total % 2000 == 0:
53             acc = float(correct) / float(total)
54             print('[%s] Clean Training accuracy: %.2f%%' % (step, acc * 100))
55             total = 0
56             correct = 0
57             accAdv = float(correctAdv) / float(totalAdv)
58             print('[%s] Adv Training accuracy: %.2f%%' % (step, accAdv * 100))
59             totalAdv = 0
60             correctAdv = 0
61
1  print("Training on Adversarial Samples")
2  adv_train(classifier, train_loader_im, criterion, thres=0.07, num_iter=20)
3
1  #Evaluating Again
2  mAP_test, test_loss, test_aps = test_classifier(val_loader_im, classifier, criterion)
3  print("Test accuracy: ", mAP_test)
4  eval_advAttack(classifier, val_loader_im, thres=0.01, num_iter=15)
5
mAP: 0.7926
Avg loss: 65.56647073068926
Test accuracy: 0.7925907258064516
Evaluating single model results on adv data
Test Accuracy = 2778 / 3968 = 70.010%
```

1

