

Appendix D

Competition: Resilience to Adversarial Attack

This is a student competition to address two key issues in modern deep learning, i.e.,

- O1 how to find better adversarial attacks, and
- O2 how to train a deep learning model with better robustness to the adversarial attacks.

We provide a template code (**Competition/Competition.py**), where there are two code blocks corresponding to the training and the attack, respectively. The two code blocks are filled with the simplest implementations representing the baseline methods, and the participators are expected to replace the baseline methods with their own implementations, in order to achieve better performance regarding the above O1 and O2.

D.1 Submissions

In the end, we will collect submissions from the students and rank them according to a pre-specified metric taking into consideration both O1 and O2. Assume that we have n students participating in this competition, and we have a set S of submissions.

Every student with student number i will submit a package $i.zip$, which includes two files:

1. $i.pt$, which is the file to save the trained model, and
2. $competition_i.py$, which is your script after updating the two code blocks in **Competition.py** with your implementations.

NB: Please carefully follow the naming convention as indicated above, and we will not accept submissions which do not follow the naming convention.

D.2 Source Code

The template source code of the competition is available at

[https://github.com/xiaoweih/
AISafetyLectureNotes/tree/main/Competition](https://github.com/xiaoweih/AISafetyLectureNotes/tree/main/Competition)

In the following, we will explain each part of the code.

D.2.1 Load Packages

First of all, the following code piece imports a few packages that are needed.

```

1 import numpy as np
2 import pandas as pd
3 import torch
4 import torch.nn as nn
5 import torch.nn.functional as F
6 from torch.utils.data import Dataset, DataLoader
7 import torch.optim as optim
8 import torchvision
9 from torchvision import transforms
10 from torch.autograd import Variable
11 import argparse
12 import time
13 import copy

```

Note: You can add necessary packages for your implementation.

D.2.2 Define Competition ID

The below line of code defines the student number. By replacing it with your own student number, it will automatically output the file *i.pt* once you trained a model.

```

1 # input id
2 id_ = 1000

```

D.2.3 Set Training Parameters

The following is to set the hyper-parameters for training. It considers e.g., batch size, number of epochs, whether to use CUDA, learning rate, and random seed. You may change them if needed.

```

1 # setup training parameters
2 parser = argparse.ArgumentParser(description='PyTorch MNIST
   Training')
3 parser.add_argument('--batch-size', type=int, default=128,
   metavar='N',
4                       help='input batch size for training (default:
   128)')
5 parser.add_argument('--test-batch-size', type=int, default=128,
   metavar='N',
6                       help='input batch size for testing (default:
   128)')
7 parser.add_argument('--epochs', type=int, default=10, metavar='N'
   ,
8                       help='number of epochs to train')
9 parser.add_argument('--lr', type=float, default=0.01, metavar='LR'
   ,
10                      help='learning rate')
11 parser.add_argument('--no-cuda', action='store_true', default=
   False,
12                      help='disables CUDA training')
13 parser.add_argument('--seed', type=int, default=1, metavar='S',
14                      help='random seed (default: 1)')
15 args = parser.parse_args(args=[])

```

D.2.4 Toggle GPU/CPU

Depending on whether you have GPU in your computer, you may toggle between devices with the below code. Just to remark that, for this competition, the usual CPU is sufficient and a GPU is not needed.

```

1 # judge cuda is available or not
2 use_cuda = not args.no_cuda and torch.cuda.is_available()
3 #device = torch.device("cuda" if use_cuda else "cpu")
4 device = torch.device("cpu")
5
6 torch.manual_seed(args.seed)
7 kwargs = {'num_workers': 1, 'pin_memory': True} if use_cuda else
   {}

```

D.2.5 Loading Dataset and Define Network Structure

In this competition, we use the same dataset (FashionMNIST) and the same network architecture. The following code specify how to load dataset and how to construct a 3-layer neural network. Please do not change this part of code.

```

1 #####don't change
  the below code
  #####
2
3 train_set = torchvision.datasets.FashionMNIST(root='data', train=
  True, download=True, transform=transforms.Compose([transforms
  .ToTensor()]))
4 train_loader = DataLoader(train_set, batch_size=args.batch_size,
  shuffle=True)
5
6 test_set = torchvision.datasets.FashionMNIST(root='data', train=
  False, download=True, transform=transforms.Compose([
  transforms.ToTensor()]))
7 test_loader = DataLoader(test_set, batch_size=args.batch_size,
  shuffle=True)
8
9 # define fully connected network
10 class Net(nn.Module):
11     def __init__(self):
12         super(Net, self).__init__()
13         self.fc1 = nn.Linear(28*28, 128)
14         self.fc2 = nn.Linear(128, 64)
15         self.fc3 = nn.Linear(64, 32)
16         self.fc4 = nn.Linear(32, 10)
17
18     def forward(self, x):
19         x = self.fc1(x)
20         x = F.relu(x)
21         x = self.fc2(x)
22         x = F.relu(x)
23         x = self.fc3(x)
24         x = F.relu(x)
25         x = self.fc4(x)
26         output = F.log_softmax(x, dim=1)
27         return output
28
29 #####end of "don't
  change the below code"
  #####

```

D.2.6 Adversarial Attack

The part is the place needing your implementation, for O1. In the template code, it includes a baseline method which uses random sampling to find adversarial attacks. You can only replace the middle part of the function with your own implementation (as indicated in the code), and are not allowed to change others.

```

1 'generate adversarial data, you can define your adversarial
  method'
2 def adv_attack(model, X, y, device):

```

```

3 X_adv = Variable(X.data)
4
5 #####Note: below is
  the place you need to edit to implement your own attack
  algorithm
6 #####
7
8 random_noise = torch.FloatTensor(*X_adv.shape).uniform_(-0.1,
9 0.1).to(device)
10 X_adv = Variable(X_adv.data + random_noise)
11
12 ##### end of attack
  method
13 #####
14 return X_adv

```

D.2.7 Evaluation Functions

Below are two supplementary functions that return loss and accuracy over test dataset and adversarially attacked test dataset, respectively. We note that the function **adv_attack** is used in the second function. You are not allowed to change these two functions.

```

1 'predict function'
2 def eval_test(model, device, test_loader):
3     model.eval()
4     test_loss = 0
5     correct = 0
6     with torch.no_grad():
7         for data, target in test_loader:
8             data, target = data.to(device), target.to(device)
9             data = data.view(data.size(0), 28*28)
10            output = model(data)
11            test_loss += F.nll_loss(output, target, size_average=
12 False).item()
13            pred = output.max(1, keepdim=True)[1]
14            correct += pred.eq(target.view_as(pred)).sum().item()
15 test_loss /= len(test_loader.dataset)
16 test_accuracy = correct / len(test_loader.dataset)
17 return test_loss, test_accuracy
18
19 def eval_adv_test(model, device, test_loader):
20     model.eval()
21     test_loss = 0
22     correct = 0
23     with torch.no_grad():
24         for data, target in test_loader:
25             data, target = data.to(device), target.to(device)
26             data = data.view(data.size(0), 28*28)

```

```

26     adv_data = adv_attack(model, data, target, device=
device)
27     output = model(adv_data)
28     test_loss += F.nll_loss(output, target, size_average=
False).item()
29     pred = output.max(1, keepdim=True)[1]
30     correct += pred.eq(target.view_as(pred)).sum().item()
31     test_loss /= len(test_loader.dataset)
32     test_accuracy = correct / len(test_loader.dataset)
33     return test_loss, test_accuracy

```

D.2.8 Adversarial Training

Below is the second place needing your implementation, for O2. In the template code, there is a baseline method. You can replace relevant part of the code as indicated in the code.

```

1  #train function, you can use adversarial training
2  def train(args, model, device, train_loader, optimizer, epoch):
3      model.train()
4      for batch_idx, (data, target) in enumerate(train_loader):
5          data, target = data.to(device), target.to(device)
6          data = data.view(data.size(0), 28*28)
7
8          #use adversarial data to train the defense model
9          #adv_data = adv_attack(model, data, target, device=device
)
10
11         #clear gradients
12         optimizer.zero_grad()
13
14         #compute loss
15         #loss = F.nll_loss(model(adv_data), target)
16         loss = F.nll_loss(model(data), target)
17
18         #get gradients and update
19         loss.backward()
20         optimizer.step()
21
22     #main function, train the dataset and print train loss, test loss
for each epoch
23     def train_model():
24         model = Net().to(device)
25
26         #
27         #####
28         ## Note: below is the place you need to edit to implement
your own training algorithm
29         ## You can also edit the functions such as train(...).

```

```

29  #
30  #####
31  optimizer = optim.SGD(model.parameters(), lr=args.lr)
32  for epoch in range(1, args.epochs + 1):
33      start_time = time.time()
34
35      #training
36      train(args, model, device, train_loader, optimizer, epoch
37  )
38
39      #get trnloss and testloss
40      trnloss, trnacc = eval_test(model, device, train_loader)
41      advloss, advacc = eval_adv_test(model, device,
42  train_loader)
43
44      #print trnloss and testloss
45      print('Epoch '+str(epoch)+' : '+str(int(time.time()-
46  start_time))+ 's', end=', ')
47      print('trn_loss: {:.4f}, trn_acc: {:.2f}%'.format(trnloss
48  , 100. * trnacc), end=', ')
49      print('adv_loss: {:.4f}, adv_acc: {:.2f}%'.format(advloss
50  , 100. * advacc))
51
52      adv_tstloss, adv_tstacc = eval_adv_test(model, device,
53  test_loader)
54      print('Your estimated attack ability, by applying your attack
55  method on your own trained model, is: {:.4f}'.format(1/
56  adv_tstacc))
57      print('Your estimated defence ability, by evaluating your own
58  defence model over your attack, is: {:.4f}'.format(
59  adv_tstacc))
60      #####
61      ## end of training method
62      #####
63
64      #save the model
65      torch.save(model.state_dict(), str(id_)+'.pt')
66      return model

```

D.2.9 Define Distance Metrics

In this competition, we take the L_∞ as the distance measure. You are not allowed to change the code.

```

1  #compute perturbation distance
2  def p_distance(model, train_loader, device):
3      p = []
4      for batch_idx, (data, target) in enumerate(train_loader):
5          data, target = data.to(device), target.to(device)

```

```

6     data = data.view(data.size(0), 28*28)
7     data_ = copy.deepcopy(data.data)
8     adv_data = adv_attack(model, data, target, device=device)
9     p.append(torch.norm(data_-adv_data, float('inf')))
10    print('epsilon p: ', max(p))

```

D.2.10 *Supplementary Code for Test Purpose*

In addition to the above code, we also provide two lines of code for testing purpose. You must comment them out in your submission. The first line is to call the **train_model()** method to train a new model, and the second is to check the quality of attack based on a model.

```

1  #Comment out the following command when you do not want to re-
   #train the model
2  #In that case, it will load a pre-trained model you saved in
   #train_model()
3  model = train_model()
4
5  #Call adv_attack() method on a pre-trained model
6  #The robustness of the model is evaluated against the infinite-
   #norm distance measure
7  #!!! Important: MAKE SURE the infinite-norm distance (epsilon p)
   #less than 0.11 !!!
8  p_distance(model, train_loader, device)

```

D.3 Implementation Actions

Below, we summarise the actions that need to be taken for the completion of a submission:

1. You must assign the variable **id_** with your student ID *i*;
2. You need to update the **adv_attack** function with your adversarial attack method;
3. You may change the hyper-parameters defined in **parser** if needed;
4. You must make sure the perturbation distance less than **0.11**, (which can be computed by **p_distance** function);
5. You need to update the **train_model** function (and some other functions that it called such as **train**) with your own training method;
6. You need to use the line “model = train_model()” to train a model and check whether there is a file **i.pt**, which stores the weights of your trained model;
7. You must submit **i.zip**, which includes two files **i.pt** (saved model) and **competition_i.py** (your script).

D.3.1 Sanity Check

Please make sure that the following constraints are satisfied. Your submission won't be marked if they are not followed.

- Submission file: please follow the naming convention as suggested above.
- Make sure your code can run smoothly.
- Comment out the two lines “model = train_model()” and “p_distance(model, train_loader , device)”, which are for test purpose.

D.4 Evaluation Criteria

Assume that, among the submissions S , we have n submissions that can run smoothly and correctly. We can get model M_i by reading the file $i.pt$.

Then, we collect the following matrix

$$\text{Score} = \begin{matrix} \mathbf{i} = \mathbf{1} \\ \mathbf{i} = \mathbf{2} \\ \dots \\ \mathbf{i} = \mathbf{n} - \mathbf{1} \\ \mathbf{i} = \mathbf{n} \end{matrix} \begin{pmatrix} s_{11} & s_{12} & \dots & s_{1(n-1)} & s_{1n} \\ s_{21} & s_{22} & \dots & s_{2(n-1)} & s_{2n} \\ \dots & \dots & \dots & \dots & \dots \\ s_{(n-1)1} & s_{(n-1)2} & \dots & s_{(n-1)(n-1)} & s_{(n-1)n} \\ s_{n1} & s_{n2} & \dots & s_{n(n-1)} & s_{nn} \end{pmatrix} \begin{matrix} \mathbf{j} = \mathbf{1} & \mathbf{j} = \mathbf{2} & \dots & \mathbf{j} = \mathbf{n} - \mathbf{1} & \mathbf{j} = \mathbf{n} \end{matrix} \tag{D.1}$$

for the mutual evaluation scores of using M_i to evaluate Atk_j (defined in function `adv_attack`). The score s_{ij} is the **test accuracy** obtained by using `adv_attack` function from the file competition `_j.py` to attack the model from $i.pt$. From Eq. (D.1), we get j 's attacking ability by letting

$$\text{AttackAbility}_j = \sum_{i=1}^n \text{Score}_{i,j} \tag{D.2}$$

to be the total of the scores of j -th column. Let **AttackAbility** be the vector of AttackAbility_j . Moreover, we get i 's defence ability by letting

$$\text{DefenceAbility}_i = \sum_{j=1}^n \text{Score}_{i,j}, \tag{D.3}$$