

Lecture 15 -- Assignment Briefing and Adversarial Attack & Defence on Deep Learning Models

Prof Xiaowei Huang

<https://cgi.csc.liv.ac.uk/~xiaowei/>

(Attendance Code: **276285**)

Information



2 (assignments) in 1 (submission)



Deadline:

17:00, 5th December 2023

17:00, 13th December 2023

Firm deadline!



Section D of Textbook: <https://link.springer.com/book/10.1007/978-981-19-6814-3>

Overall Aims

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.

<https://github.com/xiaoweihs/AISafetyLectureNotes/tree/main/Competition>

The two code blocks are filled with the simplest implementations representing the baseline methods, and the participants are expected to replace the baseline methods with their own implementations, to achieve better performance regarding the above O1 and O2.

Submission Guideline

- Must carefully follow the instructions in
 - D.1 Submissions
 - D.3 Implementation Actions
 - D.3.1 Sanity Check
- Your submission will not be marked if you do not follow the instructions and/or the markers cannot run your code.



D.1 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.

Code Template Explanation



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.3.1(1)

D.2.3 Set Training Parameters

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

D.2.5 Loading Dataset and Define Network Structure

```
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

D.3(2)

```
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
6     # the place you need to edit to implement your own attack
7     # algorithm
8     #####
9
10    random_noise = torch.FloatTensor(*X_adv.shape).uniform_(-0.1,
11                                    0.1).to(device)
12    X_adv = Variable(X_adv.data + random_noise)
13
14    ##### end of attack
15    #####
16
17    return X_adv
```

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.

D.2.8 Adversarial Training

D.3(5)

Can be
something
to consider

```
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
12             #clear gradients
13             optimizer.zero_grad()
14
15             #compute loss
16             #loss = F.nll_loss(model(adv_data), target)
17             loss = F.nll_loss(model(data), target)
18
19             #get gradients and update
20             loss.backward()
21             optimizer.step()
22
23 #main function, train the dataset and print train loss, test loss
24 #    for each epoch
25 def train_model():
26     model = Net().to(device)
27
28     #
29 ##### Note: below is the place you need to edit to implement
30 ## your own training algorithm
31 ## You can also edit the functions such as train(...).
```

```
29 #
30 #####
31
32 optimizer = optim.SGD(model.parameters(), lr=args.lr)
33 for epoch in range(1, args.epochs + 1):
34     start_time = time.time()
35
36     #training
37     train(args, model, device, train_loader, optimizer, epoch
38     )
39
40     #get trnloss and testloss
41     trnloss, trnacc = eval_test(model, device, train_loader)
42     advloss, advacc = eval_adv_test(model, device,
43     train_loader)
44
45     #print trnloss and testloss
46     print('Epoch '+str(epoch)+': '+str(int(time.time())-
47     start_time))+'s', end=' ', )
48     print('trn_loss: {:.4f}, trn_acc: {:.2f}%'.format(trnloss
49     , 100. * trnacc), end=' ', )
50     print('adv_loss: {:.4f}, adv_acc: {:.2f}%'.format(advloss
51     , 100. * advacc))
52
53     adv_tstloss, adv_tstacc = eval_adv_test(model, device,
54     test_loader)
55
56     print('Your estimated attack ability, by applying your attack
57     method on your own trained model, is: {:.4f}'.format(1/
58     adv_tstacc))
59     print('Your estimated defence ability, by evaluating your own
60     defence model over your attack, is: {:.4f}'.format(
61     adv_tstacc))
62 #####
63 ## end of training method
64 #####
65
66     #save the model
67     torch.save(model.state_dict(), str(id_)+'.pt')
68
69     return model
```

D3(6)

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.

Sanity Check –
Comment out the
two lines in your
submission.

```
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).

Do not
change

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.

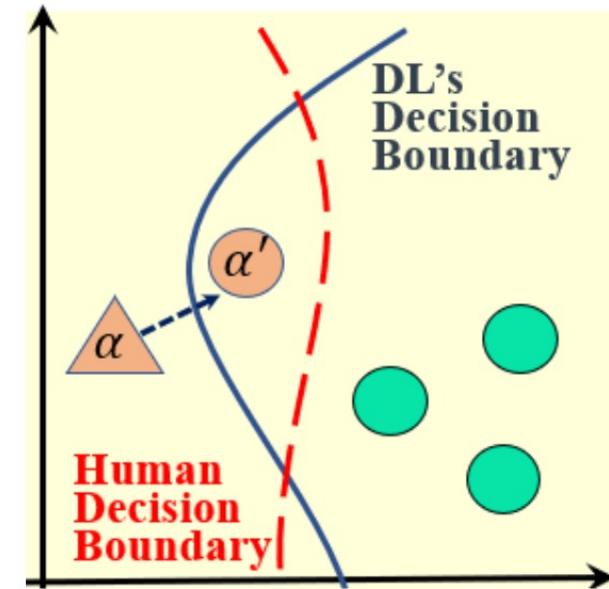
Your submission will not be marked if you do not follow the instructions and/or the markers cannot run your code.

Theoretical Knowledge on Attack and Defence of Deep Learning

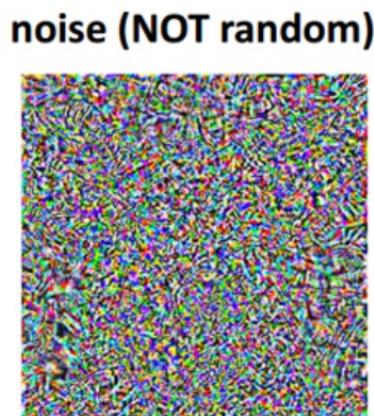


What is Adversarial Example

- ▶ Input: DL model f
A correctly-classified, genuine example α
- ▶ Aim: find a perturbed example α' , such that
 - ▶ f produces a different decision on α'
 - ▶ Human will produce a same decision on α and α'



+ 0.005 x



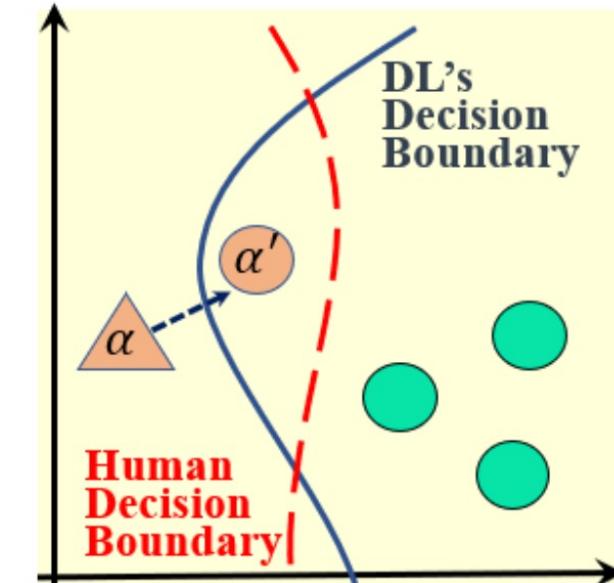
=



What is Adversarial Example

Problem: Given a DL model f and genuine example α , find α' such that

- ▶ $f(\alpha') \neq f(\alpha)$
- ▶ $D_{Human}(\alpha') = D_{Human}(\alpha)$



How to approximate human decision?

- ▶ Use certain distance metric to assure α' and α are small enough

How to search α' such that $f(\alpha') \neq f(\alpha)$?

- ▶ Design certain objective functions for minimization

What kind of information is required from DL model f ?

- ▶ White-box v.s. Black-box

Topics

Attacks

Defence

Attack Definition



One of earliest adversarial attack: optimization based formulation with L_2 -norm metric

- ▶ Model $f : \mathbb{R}^{s_1} \rightarrow \{1 \dots s_K\}$ with s_K labels
- ▶ $x \in \mathbb{R}^{s_1} = [0, 1]^{s_1}$ is an input
- ▶ $t \in \{1 \dots s_K\}$ is a target misclassification label

Find the adversarial perturbation r via

$$\begin{aligned} & \min ||r||_2 \quad \text{assure human-decision unchanged} \\ \text{s.t. } & \arg \max_l f_l(x + r) = t \quad \text{assure misclassification} \\ & x + r \in \mathbb{R}^{s_1} \quad \text{assure perturbed image feasible} \end{aligned} \tag{1}$$

- Solved by L-BFGS, Establish this direction

Attack by Adding Noise (from assignment code template)

```
'generate adversarial data, you can define your adversarial
method'
def adv_attack(model, x, y, device):
    X_adv = Variable(x.data)

    #####Note: below is
    the place you need to edit to implement your own attack
    algorithm
    #####
    random_noise = torch.FloatTensor(*x_adv.shape).uniform_(-0.1,
        0.1).to(device)
    X_adv = Variable(x_adv.data + random_noise)

    ##### end of attack
    #####
    return X_adv
```

FGSM Attack

- Fast Gradient Sign Method can find adversarial perturbations with a fixed L^∞ -norm constraint very efficiently
 - ▶ θ : the model parameters,
 - ▶ x, y : the input and the label
 - ▶ $J(\theta, x, y)$: the loss function

Find adversarial perturbation r by linearizing the loss function around the current value of θ ,

$$r = \epsilon \operatorname{sign} (\nabla_x J(\theta, x, y)) \quad (2)$$

- A **one-step modification** to all pixel values to increase the loss function with a L_∞ -norm constraint ϵ

PGD Attack

Kurakin et al (2016).
Adversarial machine learning
at scale. arXiv:1611.01236

Idea: Iterative gradient ascent to generate strongest adversarial examples, followed by projection back to the feasible region

Updating rule:

$$x^{t+1} = \Pi_{S_x}(x^t + \alpha \cdot \text{sgn}(\nabla_x L(x^t, y; \theta))),$$

where $\Pi_{S_x}(\cdot)$ projects the inputs onto the region S_x .

Training (from assignment code template)

```
#main function, train the dataset and print train loss, test loss
# for each epoch
def train_model():
    model = Net().to(device)

    #
##### Note: below is the place you need to edit to implement
## your own training algorithm
##     You can also edit the functions such as train(...).
#
#####

optimizer = optim.SGD(model.parameters(), lr=args.lr)
for epoch in range(1, args.epochs + 1):
    start_time = time.time()

    #training
    train(args, model, device, train_loader, optimizer, epoch)

    #get trnloss and testloss
    trnloss, trnacc = eval_test(model, device, train_loader)
    advloss, advacc = eval_adv_test(model, device,
train_loader)
```

```
#train function, you can use adversarial training
def train(args, model, device, train_loader, optimizer, epoch):
    model.train()
    for batch_idx, (data, target) in enumerate(train_loader):
        data, target = data.to(device), target.to(device)
        data = data.view(data.size(0),28*28)

        #use adversarial data to train the defense model
        #adv_data = adv_attack(model, data, target, device=device
        )

        #clear gradients
        optimizer.zero_grad()

        #compute loss
        #loss = F.nll_loss(model(adv_data), target)
        loss = F.nll_loss(model(data), target)

        #get gradients and update
        loss.backward()
        optimizer.step()
```

Original Training

Solve a minimization objective through SGD training

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} L(x, y; \theta)$$

Adversarial Training via Robust Optimisation

Madry et al (ICLR 2018).
Towards Deep Learning Models Resistant to Adversarial Attacks

Idea: solving a **minimax optimisation** problem through SGD training

$$\min_{\theta} \{ \mathbb{E}_{(x,y) \sim \mathcal{D}} [\max_{x' \in S_x} L(x', y; \theta)] \},$$

- ▶ (x, y) - clean training data samples $x \in \mathbb{R}^n$ with labels $y \in [k]$ drawn from the dataset \mathcal{D}
- ▶ $L(\cdot)$ - loss function with model parameter $\theta \in \mathbb{R}^m$
- ▶ $x' \in \mathbb{R}^n$ - perturbed samples in a feasible region
 $S_x \triangleq \{z : z \in B(x, \epsilon) \cap [-1.0, 1.0]^n\}$
- ▶ e.g., $B(z, \epsilon) \triangleq \{z : \|x - z\|_p \leq \epsilon\}$ - the ℓ_p -ball at center x with radius ϵ .

Adversarial Attacks as Inner Maximisation

- ▶ Outer minimisation can be simulated by SGD training

$$\min_{\theta} \left\{ \frac{1}{N} \sum_{i=1}^N \left[\max_{x' \in S_x} L(x', y; \theta) \right] \right\},$$

- ▶ How to compute gradient of a maximisation?
 - ▶ Danskin's Theorem

$$\nabla_{\theta} \max L(x', y; \theta) = \nabla_{\theta} L(x^*, y; \theta)$$

where $x^* = \arg \max L(x', y; \theta)$

- ▶ Inner maximisation $\max_{x' \in S_x} L(x', y; \theta)$ can be simulated by finding the worst-case **adversarial attacks**:
 - ▶ Fast Gradient Method (FGM)
 - ▶ Projected Gradient Method (PGM):