

ENSICAEN 2A – Security of Artificial Intelligence

Hands-on Session 2: A Simple Defense: Adversarial Training

Recap and Introduction to Defenses

In the last session, we learned how easy it is to fool a pre-trained model using the FGSM attack. Today, we'll explore a basic defense technique called *adversarial training*.

Objectives:

- Understand the principle of adversarial training.
- Implement adversarial training in PyTorch.
- Evaluate the robustness of the adversarially trained model.

Recap of Session 1:

Quickly review the FGSM attack and its implications. Remind yourselves how a small perturbation can cause a misclassification.

Feature Squeezing (Quantization)

Introduction: Feature squeezing aims to reduce the search space available to an adversary by reducing the degrees of freedom of the input. One way to achieve this is by quantization, reducing the number of bits used to represent each pixel. We will explore how to detect adversarial examples by observing the change in the model's output after quantization.

Implementation:

```
1 def quantize(data, bits=8):
2     # Quantize the image to a specified number of bits
3
4     # to be completed
5
6     return quantized_data
7
8 def score_change(model, data, quantized_data, target):
9     # Calculate the change in the model's output probabilities
10    with torch.no_grad():
11        output_original = model(data)
12    # to be completed
13    return change
14
15 # Example usage:
16 bits = 4 # Reduced bit depth
17 # Load a sample batch (data, target) from testloader
18
19 data_batch, target_batch = next(iter(testloader))
20 data = data_batch[0].unsqueeze(0) # Single image example
21 target = target_batch[0].unsqueeze(0)
22
23 data_quantized = quantize(data, bits=bits)
24 change_scores = score_change(model, data, data_quantized, target)
25
26 print(f"Score change after quantization: {change_scores.item():.4f}")
27
28 # Define a threshold for anomaly detection
29 threshold = 0.1
30
31 # Anomaly Detection Logic
```

```

32 if change_scores > threshold:
33     print("Potential adversarial example detected!")
34 else:
35     print("Likely a clean example.")

```

Question:

- Experiment with different bit depths. How does the ‘bits’ parameter affect the ‘score_change’? At what bit depth does the model’s performance degrade significantly even on clean examples? How does the change score vary between clean and adversarial images? Can you use that variability to define a threshold for adversarial detection?

Input Randomization (Translation)

Introduction: Input randomization injects randomness into the input before feeding it to the model. One simple form is random translation (shifting the image by a random number of pixels). This can disrupt adversarial perturbations that are precisely crafted for the original input.

Implementation :

```

1 import random
2 from torchvision.transforms import functional as F
3
4 def random_translation(data, max_translation=5):
5     # Apply a random translation to the image
6     # to be completed
7     return translated_data
8
9
10 def evaluate_with_randomization(model, dataloader, attack=None, epsilon=0.0,
11                                max_translation=5):
12     correct = 0
13     total = 0
14     with torch.no_grad():
15         for data, target in dataloader:
16             if attack is not None:
17                 data = attack(data, epsilon, model, target) # attack the data
18                 # Randomly translate the input
19                 # to be completed
20                 accuracy = 100 * correct / total
21                 return accuracy
22
23 # Example usage:
24 max_translation = 3 # Maximum translation in pixels
25 randomization_accuracy = evaluate_with_randomization(model, testloader, fgsm_attack,
26                                                       epsilon=0.1, max_translation=max_translation)
27
28 print(f"Accuracy with random translation (max translation={max_translation}): {
29       randomization_accuracy:.2f}%")

```

Question:

- How does the ‘max_translation’ parameter affect the model’s accuracy on clean and adversarial examples? Does increasing the maximum translation always improve robustness? Why or why not? Observe the tradeoff between robustness and accuracy on clean data.

Introducing Adversarial Training:

Adversarial training is a technique where we train a model on a combination of clean (original) examples and adversarial examples. The idea is to make the model more robust by exposing it to the types of inputs it is likely to encounter during an attack.

Adversarial Training Implementation (1 hour)

This is the core of the session. We’ll create a simplified training loop.

Defining a Simpler Model:

To speed up training, we’ll define a smaller, simpler CNN model. Here’s an example:

```

1 import torch
2 import torch.nn as nn

```

```

3 import torch.optim as optim
4 import torchvision
5 import torchvision.transforms as transforms
6
7 class SimpleCNN(nn.Module):
8     def __init__(self):
9         super(SimpleCNN, self).__init__()
10        self.conv1 = nn.Conv2d(3, 6, 5)
11        self.pool = nn.MaxPool2d(2, 2)
12        self.conv2 = nn.Conv2d(6, 16, 5)
13        self.fc1 = nn.Linear(16 * 5 * 5, 120)
14        self.fc2 = nn.Linear(120, 84)
15        self.fc3 = nn.Linear(84, 10) # Assuming 10 classes (e.g., CIFAR-10)
16
17    def forward(self, x):
18        x = self.pool(torch.nn.functional.relu(self.conv1(x)))
19        x = self.pool(torch.nn.functional.relu(self.conv2(x)))
20        x = x.view(-1, 16 * 5 * 5)
21        x = torch.nn.functional.relu(self.fc1(x))
22        x = torch.nn.functional.relu(self.fc2(x))
23        x = self.fc3(x)
24        return x
25
26 model = SimpleCNN()

```

Loading and Preparing Data:

We'll use a small subset of CIFAR-10 for training.

```

1 transform = transforms.Compose(
2     [transforms.ToTensor(),
3      transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))])
4
5 trainset = torchvision.datasets.CIFAR10(root='./data', train=True,
6                                         download=True, transform=transform)
7 trainloader = torch.utils.data.DataLoader(trainset, batch_size=4,
8                                           shuffle=True, num_workers=2) #use small
9
10 dataset
11
12 testset = torchvision.datasets.CIFAR10(root='./data', train=False,
13                                       download=True, transform=transform)
14 testloader = torch.utils.data.DataLoader(testset, batch_size=4,
15                                         shuffle=False, num_workers=2)
16
17 classes = ('plane', 'car', 'bird', 'cat',
18            'deer', 'dog', 'frog', 'horse', 'ship', 'truck')

```

Defining Optimizer and Loss Function

```

1 criterion = nn.CrossEntropyLoss()
2 optimizer = optim.Adam(model.parameters(), lr=0.001) # learning rate

```

Implementing the Training Loop:

Here's the core training loop. Study it carefully!

```

1 import torch.nn.functional as F
2 # Assume fgsm_attack function from Session 1 is defined
3
4 num_epochs = 5 # keep the training epochs low.
5
6 for epoch in range(num_epochs):
7     for i, (data, target) in enumerate(trainloader): # enumerate over the batches
8         # Generate adversarial examples
9         epsilon = 0.1
10        data_adv = fgsm_attack(data, epsilon, model, target) # generate adversarial
11        examples
12        # Combine original and adversarial data
13        combined_data = torch.cat((data, data_adv), 0)
14        combined_target = torch.cat((target, target), 0) # Duplicate the targets for the
15        adversarial examples
16
17        # Train on the combined batch
18        optimizer.zero_grad() # zero out gradients
19
20        output = model(combined_data) # get the predictions

```

```

19     loss = criterion(output, combined_target) # calculate the loss
20     loss.backward()
21     optimizer.step() # updates weights
22
23     if i % 2000 == 1999: # print every 2000 mini-batches
24         print(f'[{epoch + 1}, {i + 1:5d}] loss: {loss.item():.3f}')

```

Evaluating the Defense

Now, let's see how well our adversarially trained model performs.

Evaluation on Clean and Adversarial Examples:

Write functions to evaluate the model's accuracy on both clean images and adversarial images.

```

1 def evaluate(model, dataloader, attack=None, epsilon=0.0):
2     correct = 0
3     total = 0
4     with torch.no_grad():
5         for data, target in dataloader:
6             if attack is not None:
7                 data = attack(data, epsilon, model, target) # attack the data
8                 outputs = model(data)
9                 _, predicted = torch.max(outputs.data, 1)
10                total += target.size(0)
11                correct += (predicted == target).sum().item()
12
13    accuracy = 100 * correct / total
14    return accuracy
15
16 # Example usage:
17 epsilon = 0.1
18 clean_accuracy = evaluate(model, testloader)
19 adversarial_accuracy = evaluate(model, testloader, fgsm_attack, epsilon)
20
21 print(f"Accuracy on clean examples: {clean_accuracy:.2f}%")
22 print(f"Accuracy on adversarial examples (epsilon={epsilon}): {adversarial_accuracy:.2f}%")

```

Compare the accuracy of the adversarially trained model to the accuracy of the *original* pre-trained model (from Session 1, if you still have it) on both clean and adversarial images.

Experimenting with Epsilon Values:

Try different epsilon values for generating the adversarial examples during evaluation.

- Does training with a larger epsilon improve robustness against larger attacks?

Discussion and Wrap-up

Limitations of Adversarial Training:

Discuss the limitations of this simple adversarial training approach.

- Why is it not a perfect defense?
- Does it affect the accuracy on clean examples?
- Does it generalize to other types of attacks?

Further Exploration:

Briefly mention more advanced defense techniques (e.g., defensive distillation, certified defenses).