

---

# Lab: Fast Gradient Signed Method

## 1 Introduction

### 1.1 Adversarial learning

Adversarial learning is a machine learning method that aims to trick machine learning models by providing deceptive input. Hence, it includes both the generation and detection of adversarial examples, which are inputs specially created to deceive classifiers.

Adversarial examples are inputs intentionally designed to be in close resemblance with samples from the distribution  $D$ , but cause a misclassification.

### 1.2 Fast Gradient Signed Method (FGSM)

There are two types of adversarial attacks: White-box attacks and Black-box attacks. White-box attacks have the complete knowledge of the targeted model, including its parameter values, architecture, training method, and in some cases its training data as well. Black-box attacks feed a targeted model with the adversarial examples (during testing) that are generated without the knowledge of that model. By observing the output of the targeted model and adjusting the parameters, black-box attacks could also attack many neural network based model.

As FGSM is a type of white box attack, the attacker needs to have knowledges of the pretrained model, different model will have different parameters and loss function.

$$\tilde{x} = x + \eta \quad (1)$$

In Equation 1,  $\tilde{x}$  is the adversarial sample,  $x$  is the original input, the  $\eta$  is the adversarial perturbation added to the original input.

$$\eta = \epsilon \times \text{sign}(\nabla_x J(\theta, x, y)) \quad (2)$$

In Equation 2,  $y$  is the targets associated with  $x$ . In face recognition scenario,  $x$ ,  $y$  are face images, which  $x$  is our input images and  $y$  is the victim image( the targeting image).  $J$  is the loss function of the current model,  $\theta$  is the set of parameters in this model. In this case the coefficient  $\epsilon$  indicates the size of the step that perturbation moving backwards at the direction of the gradients. By changing the  $\epsilon$  we can control how much perturbation adding to the input images. And this “step” can also divided into several smaller steps, which is referred as granularity of the perturbation.

## 2 Environment Setup

### 2.1 Google Colab

Colab is a hosted Jupyter Notebook service that requires no setup to use and provides free access to computing resources, including GPUs and TPUs. Colab is especially well suited to machine learning, data science, and education.

To setup the Google colab:

1. Go to Google Colab using your web browser.
2. Sign in to your Google account .
3. Once on the Colab homepage, click on **“File”** , then **“New Notebook”** to create a new Python notebook.
4. Rename your notebook by clicking on the default name (e.g., “Untitled.ipynb”) at the top and enter desired name.
5. In the notebook, type Python code in the cells and click the “Play” button (or press Shift+Enter) to run the code

## 2.2 First CNN

A sample for Python code...

```
1 import numpy as np
2
3 def incmatrix(genl1,genl2):
4     m = len(genl1)
5     n = len(genl2)
6     M = None #to become the incidence matrix
7     VT = np.zeros((n*m,1), int) #dummy variable
8
9     #compute the bitwise xor matrix
10    M1 = bitxormatrix(genl1)
11    M2 = np.triu(bitxormatrix(genl2),1)
12
13    for i in range(m-1):
14        for j in range(i+1, m):
15            [r,c] = np.where(M2 == M1[i,j])
16            for k in range(len(r)):
17                VT[(i)*n + r[k]] = 1;
18                VT[(i)*n + c[k]] = 1;
19                VT[(j)*n + r[k]] = 1;
20                VT[(j)*n + c[k]] = 1;
21
22            if M is None:
23                M = np.copy(VT)
24            else:
25                M = np.concatenate((M, VT), 1)
26
27            VT = np.zeros((n*m,1), int)
28
29    return M
```

## 3 Run FGSM attack experiment

### 3.1 Libraries setup

install the **cleverhans** library, which provides tools to benchmark the vulnerability of machine learning models to adversarial examples. Open a new code cell in your Colab notebook and run

```
1 !pip install cleverhans
2 import torch
3 import torchvision.transforms as transforms
4 import torchvision.models as models
5 from PIL import Image
6 import numpy as np
7 from cleverhans.torch.attacks.fast_gradient_method import fast_gradient_method
8 import matplotlib.pyplot as plt
```

### 3.2 Pretrained model setup

- **Load ResNet-18 Model:** `model = models.resnet18(pretrained=True)`  
Initializes ResNet-18 with pre-trained weights.
- **Set to Evaluation Mode:** `model.eval()`  
Configures the model for inference by disabling dropout and batch normalization updates.

```
1 model = models.resnet18(pretrained=True)
2 model.eval()
```

### 3.3 Image Transformation Pipeline

- **Resize Image:** `transforms.Resize(256)`  
Resizes the image to 256 pixels on the shorter side, maintaining aspect ratio.
- **Center Crop Image:** `transforms.CenterCrop(224)`  
Crops a central 224x224 region from the resized image.
- **Convert to Tensor:** `transforms.ToTensor()`  
Converts the image to a tensor and scales pixel values from [0, 255] to [0.0, 1.0].

```
1 transform = transforms.Compose([
2     transforms.Resize(256),
3     transforms.CenterCrop(224),
4     transforms.ToTensor(),
5 ])
```

### 3.4 Load the Test Image

- **Load Image:** `img = Image.open('/content/dog.jpg')`  
Opens the image located at the specified path.
- **Apply Transformations and Move to Device:** `input_tensor = transform(img).unsqueeze(0).to(device)`  
Applies transformations, adds a batch dimension, and moves the tensor to the selected device(CPU or GPU).

```
1 img = Image.open('/content/dog.jpg')
2 input_tensor = transform(img).unsqueeze(0).to(device)
```

### 3.5 Normalization and Prediction

- **Define Normalization Function:** Defines a function to normalize images using pre-set mean and standard deviation values.
- **Define Model Function for Attack:** Defines a function that normalizes input images and returns model logits for attacks.
- **Get Model Prediction:** Computes and prints the model's prediction on the normalized image.

```
1
2 def normalize(x):
3     mean = torch.tensor([0.485, 0.456, 0.406]).to(device).view(1, 3, 1, 1)
4     std = torch.tensor([0.229, 0.224, 0.225]).to(device).view(1, 3, 1, 1)
5     return (x - mean) / std
6
7 with torch.no_grad():
8     output = model(normalize(input_tensor))
9     _, pred_label = torch.max(output, 1)
10    print(f'Original Prediction: {pred_label.item()}')
11
12 def model_fn(x):
13     x_norm = normalize(x)
14     logits = model(x_norm)
15     return logits
```

### 3.6 Attack Setup and Execution

- **Set Attack Parameters:** Specifies the perturbation magnitude and norm for the attack.
- **Run Fast Gradient Method Attack:** Executes the attack to generate an adversarial example using the defined parameters.
- **Get Prediction on Adversarial Example:** Computes and prints the model's prediction on the adversarial example.

```
1 # Set the attack parameters
2 eps = 0.3 # Perturbation magnitude (adjust as needed)
3 norm = np.inf # Use the infinity norm
4
5 # Run the Fast Gradient Method attack
6 adv_tensor = fast_gradient_method(
7     model_fn=model_fn,
8     x=input_tensor,
9     eps=eps,
10    norm=norm,
11    clip_min=0.0,
12    clip_max=1.0,
13    y=pred_label,
14    targeted=False,
15 )
16
17 # Get the model's prediction on the adversarial example
18 with torch.no_grad():
19     adv_output = model(normalize(adv_tensor))
20     _, adv_pred_label = torch.max(adv_output, 1)
21     print(f'Adversarial Prediction: {adv_pred_label.item()}')
```

### 3.7 Results

- Define epsilon values for perturbation strengths. The model predicts the original label and class name from the input image.
- Display the original and adversarial images with corresponding labels for different epsilon values.
- Generate adversarial images using the Fast Gradient Method for comparison.

```
1
2 epsilons = [0.02, 0.05, 0.1, 0.15, 0.2, 0.25]
3 original_label = torch.max(model(normalize(input_tensor)), 1)[1].item()
4 original_class_name = class_names[original_label]
5
6 fig, ax = plt.subplots(1, len(epsilons) + 1, figsize=(18, 6))
7
8 # Display the original image
9 ax[0].imshow(img)
10 ax[0].set_title(f'Original Image\nLabel: {original_label} ({original_class_name})')
11 ax[0].axis('off')
12
13 # Generate and display adversarial images for each epsilon
14 for i, eps in enumerate(epsilons):
15     adv_tensor = fast_gradient_method(
16         model_fn=model_fn,
17         x=input_tensor,
18         eps=eps,
19         norm=np.inf,
20         clip_min=0.0,
21         clip_max=1.0,
22         y=torch.tensor([original_label]).to(input_tensor.device),
23         targeted=False,
24     )
25     adv_img = transforms.ToPILImage()(adv_tensor.squeeze(0).cpu())
26     adv_label = torch.max(model(normalize(adv_tensor)), 1)[1].item()
27     adv_class_name = class_names[adv_label]
28
29     ax[i + 1].imshow(adv_img)
30     ax[i + 1].set_title(f'Epsilon: {eps}\nLabel: {adv_label} ({adv_class_name})')
31     ax[i + 1].axis('off')
32
33 plt.show()
```



Figure 1: Change of Labels with variation in Epsilon  $\epsilon$  value

## 4 References

- [1] Deep Learning by Ian Goodfellow et al
- [2] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199, 2013.
- [3] <https://dl.acm.org/doi/pdf/10.1145/3398394>.
- [4] Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z Berkay Celik, and Ananthram Swami. The limitations of deep learning in adversarial settings. In 2016 IEEE European Symposium on Security and Privacy (EuroS&P), pages 372–387. IEEE, 2016.
- [5] [https://d2l.ai/chapter\\_convolutional\\_neural\\_networks/lenet.html](https://d2l.ai/chapter_convolutional_neural_networks/lenet.html)
- [6] Liu, Yujie, et al. "Sensitivity of adversarial perturbation in fast gradient sign method." 2019 IEEE symposium series on computational intelligence (SSCI). IEEE, 2019.