

## ▼ Lab2: Adversarial Attacks on Deep Neural Networks

### ▼ Load Dataset

```

1  import tensorflow as tf
2  from tensorflow.keras.datasets import mnist, cifar10, cifar100
3  from tensorflow.keras import Sequential
4  from tensorflow.keras.callbacks import LambdaCallback
5  from tensorflow.keras.layers import Conv2D, MaxPooling2D, Dropout, Dense, Flat
6  import numpy as np
7  import pandas as pd
8  import random
9  import matplotlib.pyplot as plt

```

```

1  #loading the MNIST
2  fashion_mnist = tf.keras.datasets.fashion_mnist
3  (xTrain, yTrain), (xTest, yTest) = fashion_mnist.load_data()

```

↳ Downloading data from <https://storage.googleapis.com/tensorflow/tf-keras-datasets>  
 32768/29515 [=====] - 0s 0us/step  
 40960/29515 [=====] - 0s 0us/step  
 Downloading data from <https://storage.googleapis.com/tensorflow/tf-keras-datasets>  
 26427392/26421880 [=====] - 0s 0us/step  
 26435584/26421880 [=====] - 0s 0us/step  
 Downloading data from <https://storage.googleapis.com/tensorflow/tf-keras-datasets>  
 16384/5148 [=====]  
 Downloading data from <https://storage.googleapis.com/tensorflow/tf-keras-datasets>  
 4423680/4422102 [=====] - 0s 0us/step  
 4431872/4422102 [=====] - 0s 0us/step

```

1  # Plot picture
2  plt.figure()
3  plt.imshow(xTrain[0])
4  plt.colorbar()
5  plt.grid(False)
6  plt.show()

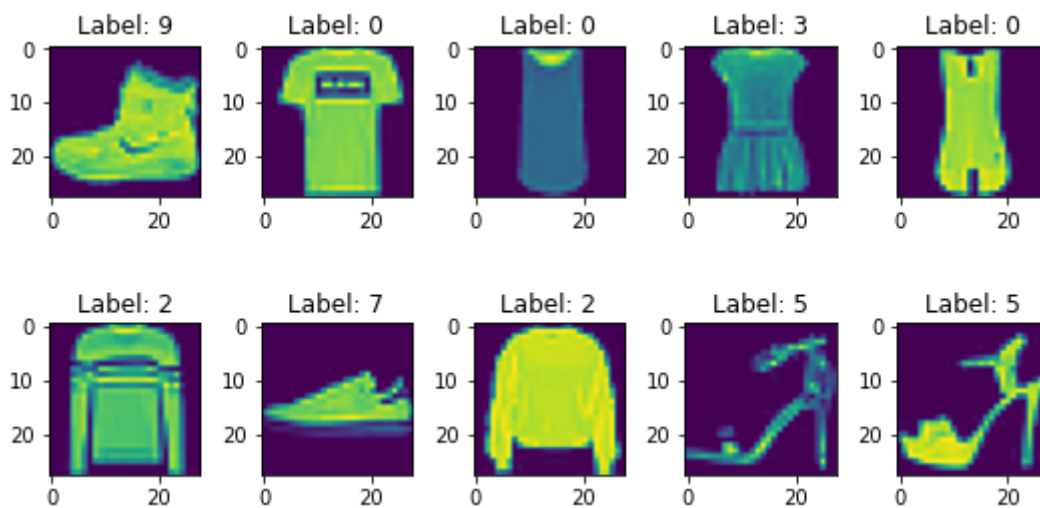
```



```

1 num = 10
2 images = xTrain[:num]
3 labels = yTrain[:num]
4
5 num_row = 2
6 num_col = 5 # plot images
7 fig, axes = plt.subplots(num_row, num_col, figsize=(1.5*num_col,2*num_row))
8 for i in range(num):
9     ax = axes[i//num_col, i%num_col]
10    ax.imshow(images[i])
11    ax.set_title('Label: {}'.format(labels[i]))
12 plt.tight_layout()
13 plt.show()

```



```
1 xTrain[0].shape
```

```
(28, 28)
```

```

1 # Preprocessing
2 img_rows, img_cols = xTrain[0].shape
3 numOfClasses = 10
4
5 xTrain = xTrain / 255.0
6 xTest = xTest / 255.0

```

## ▼ Baseline DNN

### Q1: Implement baseline DNN

```

1 baseModel = tf.keras.models.Sequential([
2     tf.keras.layers.Flatten(input_shape=(28, 28)),
3     tf.keras.layers.Dense(128, activation='relu'),
4     tf.keras.layers.Dropout(0.2),

```

```

5     tf.keras.layers.Dense(10)
6 ])
7 baseModel.compile(optimizer='adam',
8                   loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
9                   metrics=['accuracy'])
10 baseModel.fit(xTrain, yTrain, epochs=10)

```

```

Epoch 1/10
1875/1875 [=====] - 4s 2ms/step - loss: 0.5351 - accu
Epoch 2/10
1875/1875 [=====] - 4s 2ms/step - loss: 0.4008 - accu
Epoch 3/10
1875/1875 [=====] - 4s 2ms/step - loss: 0.3672 - accu
Epoch 4/10
1875/1875 [=====] - 4s 2ms/step - loss: 0.3444 - accu
Epoch 5/10
1875/1875 [=====] - 4s 2ms/step - loss: 0.3307 - accu
Epoch 6/10
1875/1875 [=====] - 4s 2ms/step - loss: 0.3190 - accu
Epoch 7/10
1875/1875 [=====] - 4s 2ms/step - loss: 0.3075 - accu
Epoch 8/10
1875/1875 [=====] - 4s 2ms/step - loss: 0.2982 - accu
Epoch 9/10
1875/1875 [=====] - 4s 2ms/step - loss: 0.2919 - accu
Epoch 10/10
1875/1875 [=====] - 4s 2ms/step - loss: 0.2821 - accu
<keras.callbacks.History at 0x7f2f91d7f910>

```

```
1 print("Base loss and accuracy on regular images:", baseModel.evaluate(x=xTest, y
```

```
Base loss and accuracy on regular images: [0.3400658667087555, 0.8801000118255
```

## ▼ FGSM based untargeted attacks

### Q2: Implement FGS untargetd attack

```

1 def get_prtbtns(X, label, loss_func, model):
2     perturbations = np.zeros((X.shape[0], X.shape[1], X.shape[2]))
3     for i in range(X.shape[0]):
4         x,y = X[i:i+1],label[i]
5         with tf.GradientTape() as gt:
6             gt.watch(x)
7             p = model(x)
8             loss = loss_func(y, p)
9             grad = gt.gradient(loss, x)
10            perturbations[i] = tf.sign(grad)
11    return perturbations
12
13 xTest_tensor = tf.convert_to_tensor(xTest)
14 loss_func = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
15 perturbations = get_prtbtns(xTest_tensor, yTest, loss_func, baseModel)

```

```

1 def FGS_untargeted_attack(X, label, model, eps=1):
2     newX = np.zeros((X.shape[0], X.shape[1], X.shape[2]))
3
4     for i in range(X.shape[0]):
5         newX[i] = X[i] + e/255.0 * perturbations[i]
6     newX = tf.clip_by_value(newX, 0, 1)
7     return newX

```

**Q3:** Calculate the fraction of test images that were correctly classified by the baseline DNN that are mis-classified after adversarial perturbation, as a function of  $\epsilon$ , and plot the result.

```

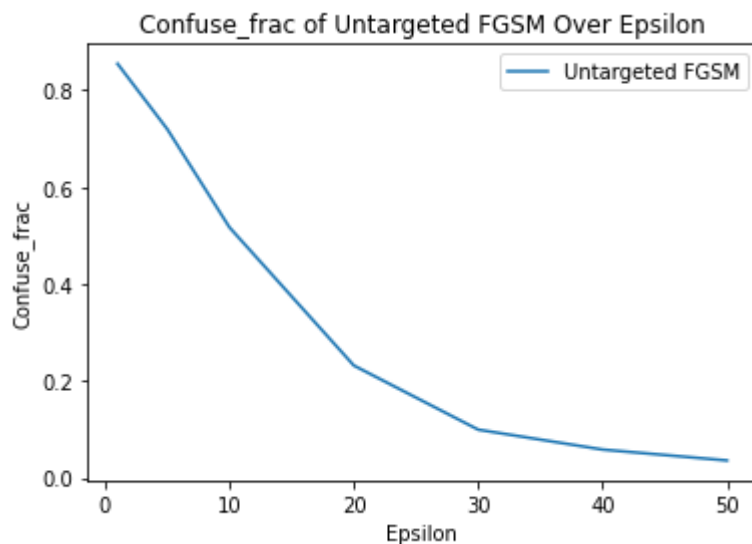
1 epss = [1, 5, 10, 20, 30, 40, 50]
2 rpt = pd.DataFrame(index=epss, columns=['Confuse_frac'])
3
4 for e in epss:
5     xTestAftAttack = FGS_untargeted_attack(xTest, yTest, baseModel, eps=e)
6
7     loss, acc = baseModel.evaluate(x=xTestAftAttack, y=yTest, verbose=0)
8     rpt.loc[e, 'Confuse_frac'] = acc
9

```

```

1 plt.plot(epss, rpt['Confuse_frac'], label = "Untargeted FGSM")
2 plt.xlabel('Epsilon')
3 plt.ylabel('Confuse_frac')
4 plt.title('Confuse_frac of Untargeted FGSM Over Epsilon')
5 plt.legend()
6 plt.show()

```



## ▼ FGSM based targeted attacks

**Q4:** Implement FGS targetd attack

```

1 yTest_target = (yTest + 1) % 10
2
3 xTest_tensor = tf.convert_to_tensor(xTest)
4 loss_func = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
5 perturbations_target = get_prtbtns(xTest_tensor, yTest_target, loss_func, baseMo

```

```

1 def FGS_targeted_attack(X, label, model, eps=1):
2     newX = np.zeros((X.shape[0], X.shape[1], X.shape[2]))
3
4     for i in range(X.shape[0]):
5         newX[i] = X[i] - e/255.0 * perturbations_target[i]
6     newX = tf.clip_by_value(newX, 0, 1)
7     return newX

```

**Q5:** Report the attack's success rate as a function of parameter  $\epsilon$  where success rate is defined as the fraction of test images that were correctly classified by the baseline DNN that are misclassified after adversarial perturbations with label  $(i+1)\%10$ .

```

1 epss = [1, 5, 10, 20, 30, 40, 50]
2 rpt = pd.DataFrame(index=epss, columns=['Mislead_frac'])
3
4 for e in epss:
5     xTestAftAttack = FGS_targeted_attack(xTest, yTest_target, baseModel, eps=e)
6
7     loss, acc = baseModel.evaluate(x=xTestAftAttack, y=yTest, verbose=0)
8     rpt.loc[e, 'Mislead_frac'] = acc
9 rpt

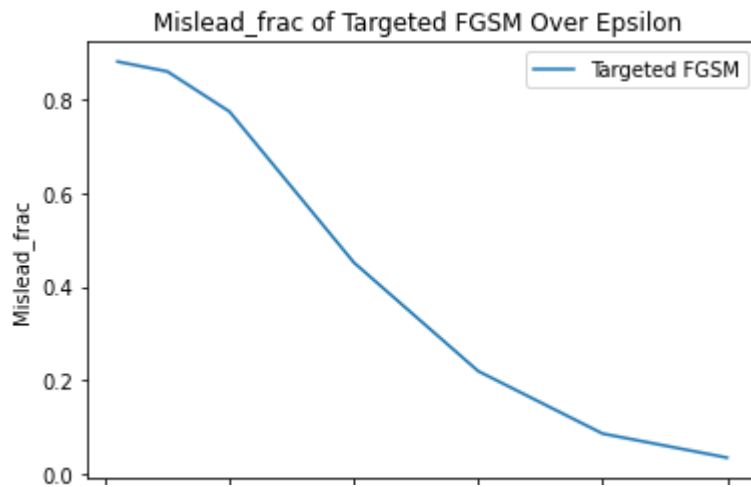
```

	Mislead_frac
1	0.88
5	0.859
10	0.7734
20	0.4515
30	0.2198
40	0.0866
50	0.0354

```

1 plt.plot(epss, rpt['Mislead_frac'], label = "Targeted FGSM")
2 plt.xlabel('Epsilon')
3 plt.ylabel('Mislead_frac')
4 plt.title('Mislead_frac of Targeted FGSM Over Epsilon ')
5 plt.legend()
6 plt.show()

```



## ▼ Adversarial Retraining against Untargeted FGSM Attacks:

**Q6:** Append the adversarially perturbed images to training set, but using their correct labels.

```
1 e = 10
2 loss_func = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
3
4 def attack_train(X, train_target, model, eps=10):
5     newX = np.zeros((X.shape[0], X.shape[1], X.shape[2]))
6     for i in range(X.shape[0]):
7         newX[i] = X[i] + e/255.0 * train_target[i]
8     newX = tf.clip_by_value(newX, 0, 1)
9     return newX
10
11 xTrain_tensor = tf.convert_to_tensor(xTrain)
12 perturbations_train_target = get_prtbtns(xTrain_tensor, yTrain, loss_func, baseM
13 xTrain_AftAttack = attack_train(xTrain, perturbations_train_target, baseModel, e

1 newX = np.concatenate((xTrain, xTrain_AftAttack), axis=0)
2 newY = np.concatenate((yTrain, yTrain), axis=0)
```

```
1 advModel = tf.keras.models.Sequential([
2     tf.keras.layers.Flatten(input_shape=(28, 28)),
3     tf.keras.layers.Dense(128, activation='relu'),
4     tf.keras.layers.Dropout(0.2),
5     tf.keras.layers.Dense(10)
6 ])
7 advModel.compile(optimizer='adam',
8                 loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=Tru
9                 metrics=['accuracy'])
10 advModel.fit(newX, newY, epochs=10)
```

Epoch 1/10

3750/3750 [=====] - 9s 2ms/step - loss: 0.4575 - accu

Epoch 2/10

3750/3750 [=====] - 9s 2ms/step - loss: 0.3103 - accu

Epoch 3/10

```

3750/3750 [=====] - 9s 2ms/step - loss: 0.2660 - accu
Epoch 4/10
3750/3750 [=====] - 9s 2ms/step - loss: 0.2451 - accu
Epoch 5/10
3750/3750 [=====] - 9s 2ms/step - loss: 0.2247 - accu
Epoch 6/10
3750/3750 [=====] - 9s 2ms/step - loss: 0.2123 - accu
Epoch 7/10
3750/3750 [=====] - 9s 2ms/step - loss: 0.2031 - accu
Epoch 8/10
3750/3750 [=====] - 9s 2ms/step - loss: 0.1912 - accu
Epoch 9/10
3750/3750 [=====] - 9s 2ms/step - loss: 0.1867 - accu
Epoch 10/10
3750/3750 [=====] - 9s 2ms/step - loss: 0.1812 - accu
<keras.callbacks.History at 0x7f30140b8450>

```

**Q7:** Report the classification accuracy of the adversarially retrained DNN on the original test dataset that contains only clean inputs. Implement FGSM based untargeted attacks using images from the clean test set on the adversarially retrained DNN. Report the success rate of your attack. Is the adversarially retrained DNN robust against adversarial perturbations?

```

1 new_loss, new_acc = advModel.evaluate(x=xTest, y=yTest, verbose=0)
2 print("The new accuracy for clean test input: ", new_acc)

```

The new accuracy for clean test input: 0.8755000233650208

```

1 e=10
2 perturbations = get_prtbtns(xTest_tensor, yTest, loss_func, advModel)
3 xTestAftAttack = FGS_untargeted_attack(xTest, yTest, advModel, eps=e)
4 new_attack_loss, new_attack_acc = advModel.evaluate(x=xTestAftAttack, y=yTest, v
5 print("The new accuracy for untargeted attack input: ", new_attack_acc)
6 # costPerPixel = # Your Code

```

The new accuracy for untargeted attack input: 0.7906000018119812

## ▼ Challenge

In this section, I'm going to try IFGSM.

```

1 eps = 5 / 255.0
2 steps = 20
3 step_alpha = 0.005

```

```

1 def IFGS_untargeted_attack(x, y, loss_func, model):
2     perturbation = np.zeros(x.shape)
3     for step in range(steps):
4         with tf.GradientTape() as gt:
5             gt.watch(x)

```

```

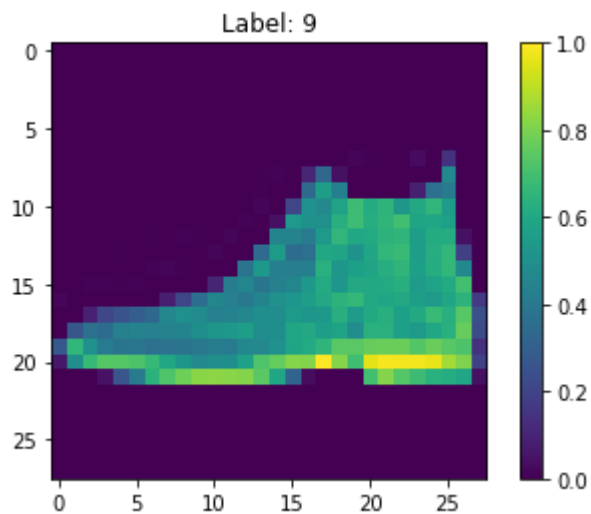
6         p = model(x)
7         loss = loss_func(y, p)
8         grad = gt.gradient(loss, x)
9         perturbation = tf.sign(grad)
10        x += step_alpha * perturbation
11        x = tf.clip_by_value(x, 0, 1)
12    x += eps * perturbation
13    x = tf.clip_by_value(x, 0, 1)
14    # print(x)
15    return x

```

```

1 # Plot picture
2 plt.figure()
3 plt.imshow(xTest[0])
4 plt.title('Label: {}'.format(yTest[0]))
5 plt.colorbar()
6 plt.grid(False)
7 plt.show()

```



```

1 xTest_tensor = tf.convert_to_tensor(xTest)
2 loss_func = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
3 new_img = IFGS_untargeted_attack(xTest_tensor[0:1], yTest[0], loss_func, baseMod

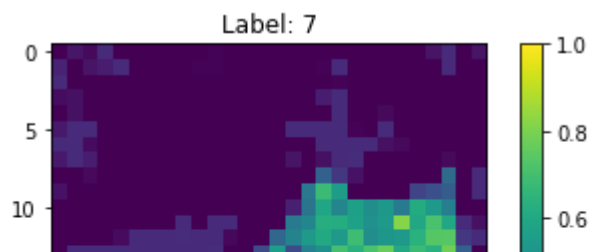
```

```

1 new_y = tf.argmax(baseModel(new_img)[0])
2 plt.figure()
3 plt.imshow(new_img[0])
4 plt.title('Label: {}'.format(new_y))
5 plt.colorbar()
6 plt.grid(False)
7 plt.show()

```





```
1 save_new_img = np.array(new_img)[0]*255
2 # ssave_new_img = save_new_img.astype(int)
```



```
1 from PIL import Image
2 im = Image.fromarray(save_new_img)
3 im.convert('RGB').save("Challenge_picture.jpeg")
```

✓ 0s completed at 11:31 AM

