# IT581 Adversarial Machine Learning

## Lab Assignment - 5

## Group 05

**2021 17013: Shaikh Faizan Ahmed**
**2021 17014: Sonam Bharti**

## Question 1.

In this exercise use **Uniform noise** and Correlated Gaussian noise and compare the probability of a successful attack for all three cases.
Write code for this and show result for the sample images. You can use any classification technique and data to perform this task.
Comment your observation.

**Answer**

**code**

```
1  # Import libraries
2  !pip install adversarial-robustness-toolbox
3
4  import cv2
5  import numpy as np
6  import matplotlib.image as mpimg
7  import matplotlib.pyplot as plt
8  import os
9  import sys
10 import PIL
11 from PIL import Image
12 from matplotlib.pyplot import figure
13
14 import torch.nn as nn
15 from PIL import Image
16 from torchvision import transforms
17 import torch.optim as optim
18 import torchvision.models as models
19 from art.estimators.classification import PyTorchClassifier
20 import warnings
21 warnings.filterwarnings('ignore')
22 import scipy.signal
23
24 from tensorflow.keras.applications.resnet50 import decode_predictions
25
26 # Softmax function
27 def softmax_activation(inputs):
28     inputs = inputs.tolist()
29     exp_values = np.exp(inputs - np.max(inputs))
30
31     # Normalize
32     probabilities = exp_values / np.sum(exp_values)
33     return probabilities
34
35 # image transformation
36
37 preprocess = transforms.Compose([
38     transforms.Resize(224),
39     transforms.CenterCrop(224),
40     transforms.ToTensor()])
41
42 model_resnet50 = models.resnet50(pretrained=True)
43
44 criterion = nn.CrossEntropyLoss()
```

```
45
46 # Create the ART classifier
47
48 classifier = PyTorchClassifier(
49     model=model_resnet50,
50     loss=criterion,
51     input_shape=(3, 224, 224),
52     nb_classes=1000,
53     device_type='gpu')
```

Now We will classify the original Image.

```
1  img = cv2.imread('/content/panda.jpg')
2
3  color_coverted = cv2.cvtColor(img, cv2.COLOR_BGR2RGB)
4  original_image = Image.fromarray(color_coverted)
5
6  input_tensor = preprocess(original_image)
7  input_batch = input_tensor.unsqueeze(0).numpy().astype(np.float32)
8  input = input_batch[0].transpose((1,2,0))
9
10
11 label = decode_predictions(classifier.predict(input_batch))
12 label = label[0][0]
13
14 preds = classifier.predict(input_batch)
15
16 accuracy = np.max(softmax_activation(preds), axis=1)
17 accuracy = round(accuracy[0], 2)
18
19 plt.imshow(input)
20 plt.title(label[1] + " " + format(accuracy * 100) + "%", fontsize=20)
21 plt.axis('off')
22 plt.show()
```

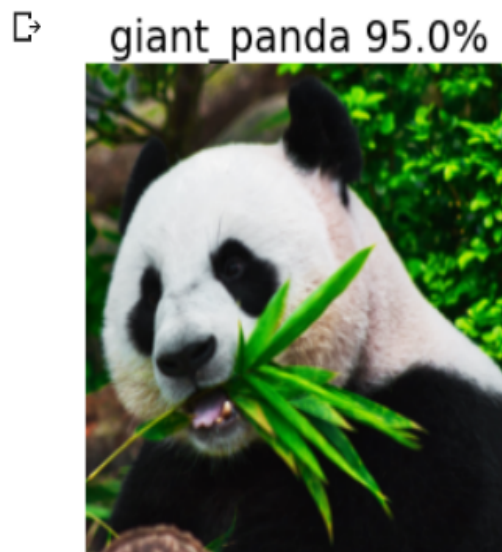The given below figure represents the classification of original image or say, clean image.



Figure 1:

**Gaussian Noise Attack**

```
1  figure(figsize=(12, 10), dpi=80)
2
3  plt.subplot(331)
4  plt.imshow(img)
5  plt.title('Image')
6  plt.axis('off')
7
8
9  # Generate Gaussian noise
10 gauss = np.random.normal(1,20,img.size)
11 gauss = gauss.reshape(img.shape[0],img.shape[1],img.shape[2]).astype('uint8')
12
13 plt.subplot(332)
14 plt.imshow(gauss)
15 plt.title('Gaussian Noise')
16 plt.axis('off')
17
18
19 # Add the Gaussian noise to the image
20 img_gauss = cv2.add(img,gauss)
21
22 plt.subplot(333)
23 plt.imshow(img_gauss)
24 plt.title('Attacked Img')
25 plt.axis('off')
```
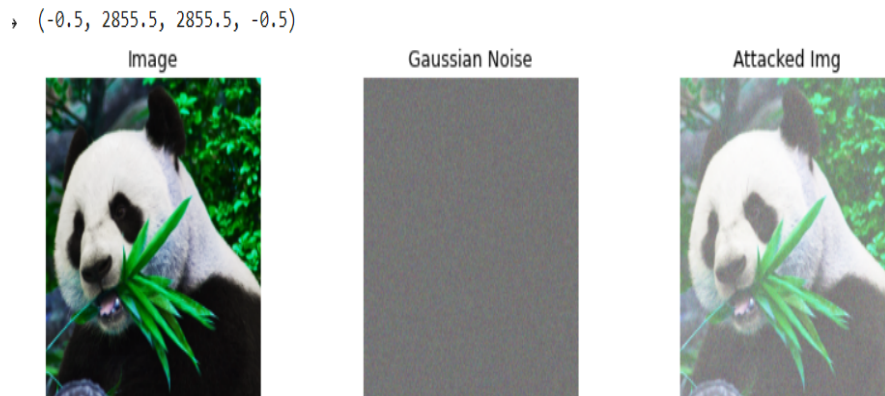
**Output**



Figure 2:

Now Classifying the Adversarial Image, we get the below output.

```
1  color_coverted = cv2.cvtColor(img_gauss, cv2.COLOR_BGR2RGB)
2  adv_image = Image.fromarray(color_coverted)
3
4  input_tensor = preprocess(adv_image)
5  input_batch1 = input_tensor.unsqueeze(0).numpy().astype(np.float32)
6  adv_input = input_batch1[0].transpose((1,2,0))
7
8  label = decode_predictions(classifier.predict(input_batch1))
9  label = label[0][0]
10
11 preds = classifier.predict(input_batch1)
12
13 accuracy = np.max(softmax_activation(preds), axis=1)
14 accuracy = round(accuracy[0], 2)
15
16 plt.imshow(adv_input)
17 plt.title(label[1] + " " + format(accuracy * 100) + "%", fontsize=20)
18 plt.axis('off')
19 plt.show()
```

tray 17.0%

Figure 3:

## Uniform Noise Attack

```
1  figure(figsize=(12, 10), dpi=80)
2
3  plt.subplot(331)
4  plt.imshow(img)
5  plt.title('Image')
6  plt.axis('off')
7
8  #plt.resize(img,height=33,width=45)
9
10 uniform = np.random.uniform(100,250,img.size)
11 uniform_noise = uniform.reshape(img.shape[0],img.shape[1],img.shape[2]).astype('uint8')
12 #uniform_noise = cv2.randu(uniform_noise,0,255)
13 plt.subplot(332)
14 plt.imshow(uniform_noise)
15 plt.title('Unifrom Noise')
16 plt.axis('off')
17
18 #uniform_noise = (uniform_noise*0.5).astype(np.uint8)
19 noisy_image2 = cv2.add(img,uniform_noise)
20
21 plt.subplot(333)
22 plt.imshow(noisy_image2)
23 plt.title('Attacked Img')
24 plt.axis('off')
```

(-0.5, 2855.5, 2855.5, -0.5)



Figure 4:

4

Now Classifying the Adversarial Image, we get the below output.

```
1  color_coverted = cv2.cvtColor(noisy_image2, cv2.COLOR_BGR2RGB)
2  adv_image = Image.fromarray(color_coverted)
3
4  input_tensor = preprocess(adv_image)
5  input_batch2 = input_tensor.unsqueeze(0).numpy().astype(np.float32)
6  adv_input = input_batch2[0].transpose((1,2,0))
7
8  label = decode_predictions(classifier.predict(input_batch2))
9  label = label[0][0]
10
11 preds = classifier.predict(input_batch2)
12
13 accuracy = np.max(softmax_activation(preds), axis=1)
14 accuracy = round(accuracy[0], 2)
15
16 plt.imshow(adv_input)
17 plt.title(label[1] + " " + format(accuracy * 100) + "%", fontsize=20)
18 plt.axis('off')
19 plt.show()
```



Figure 5:

**Correlated Gaussian noise**

```
1  figure(figsize=(12, 10), dpi=80)
2
3  plt.subplot(331)
4  plt.imshow(img)
5  plt.title('Image')
6  plt.axis('off')
7
8
9  # Generate correlated Gaussian noise
10 mean = [40, 1000, 90]
11 cov = [[10, 0, 0], [0, 3000, 0],[0, 0, 4000]]
12 cov_gauss = np.random.multivariate_normal(mean, cov,8156736)
13
14 cgauss_noise = cov_gauss.reshape(2856,2856,3).astype('uint8')
15 plt.subplot(332)
16 plt.imshow(cgauss_noise)
17 plt.title('Correlated Gaussian Noise')
18 plt.axis('off')
19
20 # Add the correlated Gaussian noise to the image
21 noisy_image3 = cv2.add(img,cgauss_noise)
22 plt.subplot(333)
```

```
23 plt.imshow(noisy_image3)
24 plt.title('Attacked Img')
25 plt.axis('off')
```
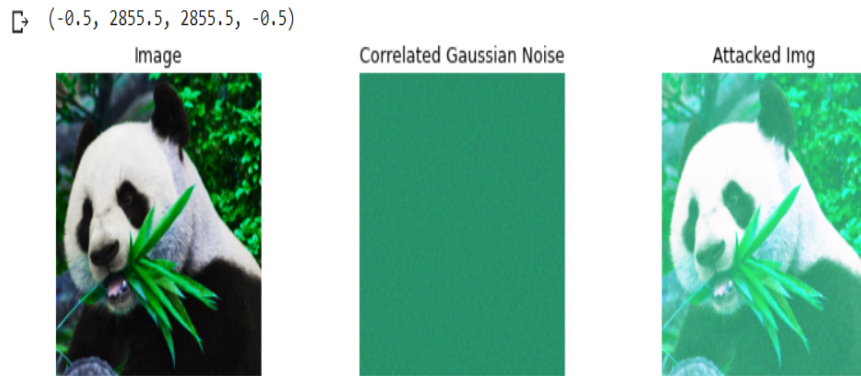

(-0.5, 2855.5, 2855.5, -0.5)

Figure 6:

Now Classifying the Adversarial Image, we get the below output.

```
1  color_coverted = cv2.cvtColor(noisy_image3, cv2.COLOR_BGR2RGB)
2  adv_image = Image.fromarray(color_coverted)
3
4  input_tensor = preprocess(adv_image)
5  input_batch3 = input_tensor.unsqueeze(0).numpy().astype(np.float32)
6  adv_input = input_batch3[0].transpose((1,2,0))
7
8  label = decode_predictions(classifier.predict(input_batch3))
9  label = label[0][0]
10
11 preds = classifier.predict(input_batch3)
12
13 accuracy = np.max(softmax_activation(preds), axis=1)
14 accuracy = round(accuracy[0], 2)
15
16 plt.imshow(adv_input)
17 plt.title(label[1] + " " + format(accuracy * 100) + "%", fontsize=20)
18 plt.axis('off')
19 plt.show()
```


tray 14.000000000000002%

Figure 7:

6

**Observation**

For same dimension of same image (say, our first image of Panda) the all three attack have different probability to misclassify. Uniform noise has the least probability while Gaussian noise has the highest probability to misclassify the image and the correlated gaussian noise has average probability.

Here, we have observed that for different dimensions of images accuracy/probability is also changing. When we checked it with lower dimension image (such as eagle) and the output is given below.
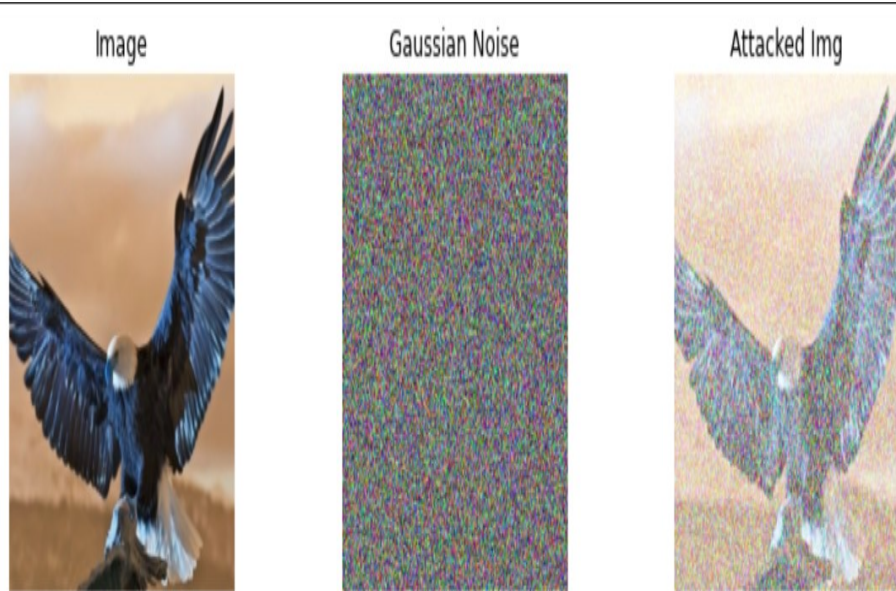


Figure 8:



Figure 9:

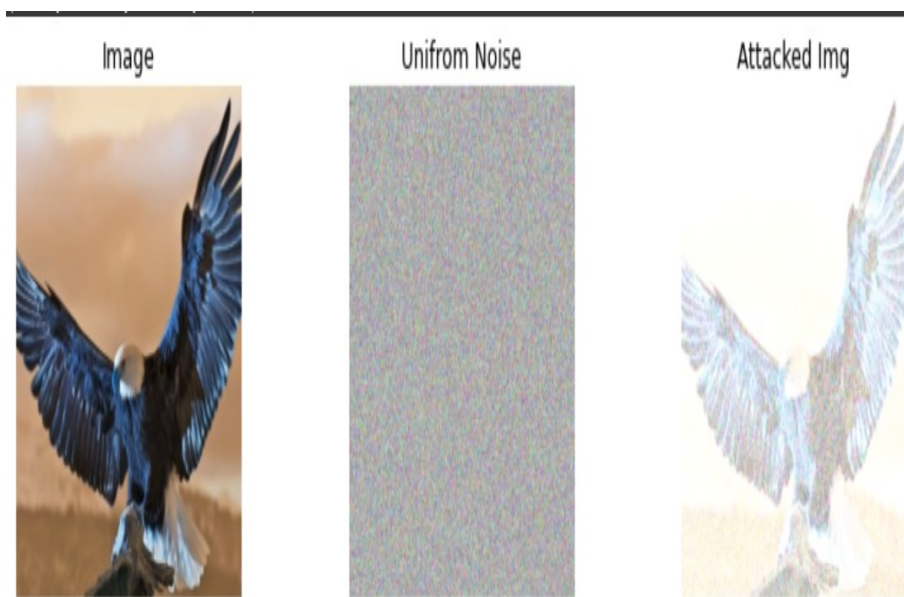Hence, We can conclude that for higher dimensions the probability is decreasing to perturbe the image.
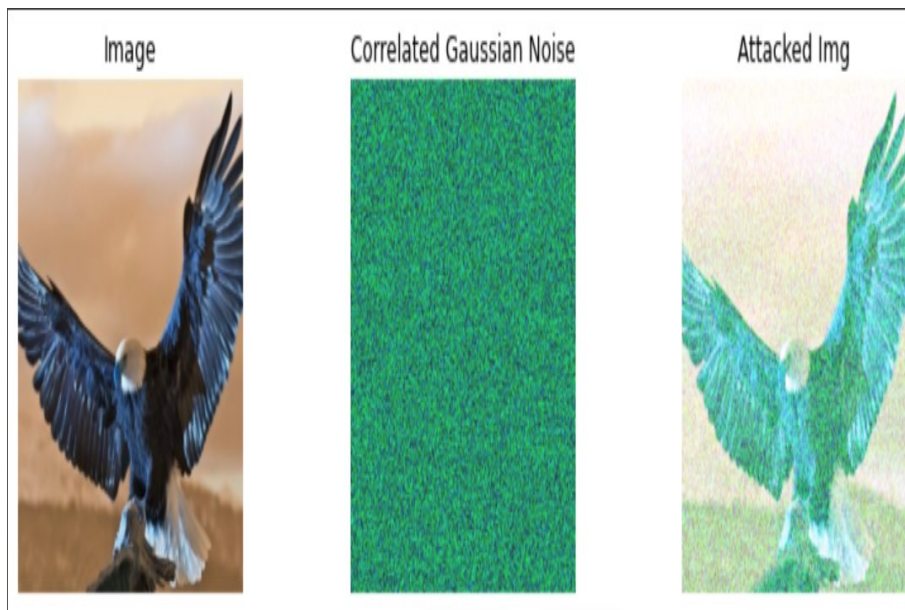
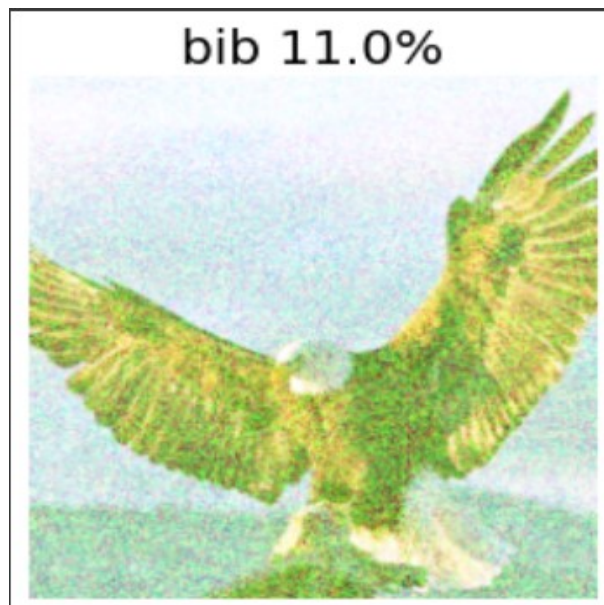Figure 10:



Figure 11:

Figure 12:



Figure 13: