МИНОБРНАУКИ РОССИИ

Федеральное государственное бюджетное образовательное учреждение высшего образования

**«МИРЭА – Российский технологический университет»**

**РТУ МИРЭА**

Кафедра: КБ-4 «Киберразведка и противодействие угрозам с применением технологий искусственного интеллекта»

**Лабораторная работа №1**

**по дисциплине**

**«Анализ защищенности систем искусственного интеллекта»**

<div align="right">

Выполнил:

Филимонов И.М.

Группа:

ББМО-01-22, 2 курс


Проверил:

Спирин А.А.

</div>

Москва, 2024 г.

Скопируем проект по ссылке в локальную среду выполнения

```
!git clone https://github.com/ewatson2/EEL6812_DeepFool_Project.git
```

```
Cloning into 'EEL6812_DeepFool_Project'...
remote: Enumerating objects: 96, done.
remote: Counting objects: 100% (3/3), done.
remote: Compressing objects: 100% (2/2), done.
remote: Total 96 (delta 2), reused 1 (delta 1), pack-reused 93
Receiving objects: 100% (96/96), 33.99 MiB | 18.13 MiB/s, done.
Resolving deltas: 100% (27/27), done.
```

Сменим директорию исполнения на вновь созданную папку
"EEL6812_DeepFool_Project" проекта

```
%cd /content/EEL6812_DeepFool_Project
```

```
/content/EEL6812_DeepFool_Project
```

Выполним импорт библиотек

```
import numpy as np
import json, torch
from torch.utils.data import DataLoader, random_split
from torchvision  import datasets, models
from torchvision.transforms import transforms
```

Выполним импорт вспомогательных библиотек из локальных файлов проекта

```
from models.project_models import FC_500_150, LeNet_CIFAR, LeNet_MNIST, Net
from utils.project_utils import get_clip_bounds, evaluate_attack, display_attack
```

Установим случайное рандомное значение в виде переменной rand_seed для
варианта 10

```
rand_seed = 10

# Установим указанное значение для np.random.seed и torch.manual_seed
np.random.seed(rand_seed)
torch.manual_seed(rand_seed)
```

```
<torch._C.Generator at 0x781914375af0>
```

Используем в качестсве устройства видеокарту

```
use_cuda = torch.cuda.is_available()
device = torch.device('cuda' if use_cuda else 'cpu')
```

Загрузим датасет MNIST с параметрами mnist_mean = 0.5, mnist_std = 0.5, mnist_dim = 28

```
mnist_mean = 0.5
mnist_std = 0.5
mnist_dim = 28

mnist_min, mnist_max = get_clip_bounds(mnist_mean, mnist_std, mnist_dim)

mnist_min = mnist_min.to(device)
mnist_max = mnist_max.to(device)

mnist_tf = transforms.Compose([ transforms.ToTensor(), transforms.Normalize( mean=mnist_mean, std=mnist_std)])

mnist_tf_train = transforms.Compose([ transforms.RandomHorizontalFlip(), transforms.ToTensor(), transforms.Normalize( me

mnist_tf_inv = transforms.Compose([ transforms.Normalize( mean=0.0, std=np.divide(1.0, mnist_std)), transforms.Normalize

mnist_temp = datasets.MNIST(root='datasets/mnist', train=True, download=True, transform=mnist_tf_train)
mnist_train, mnist_val = random_split(mnist_temp, [50000, 10000])
mnist_test = datasets.MNIST(root='datasets/mnist', train=False, download=True, transform=mnist_tf)
```

```
Downloading http://yann.lecun.com/exdb/mnist/train-images-idx3-ubyte.gz
Downloading http://yann.lecun.com/exdb/mnist/train-images-idx3-ubyte.gz to datasets/mnist/MNIST/raw/train-images-idx3-ubyt
e.gz
100%|████████| 9912422/9912422 [00:00<00:00, 190126449.62it/s]
Extracting datasets/mnist/MNIST/raw/train-images-idx3-ubyte.gz to datasets/mnist/MNIST/raw

Downloading http://yann.lecun.com/exdb/mnist/train-labels-idx1-ubyte.gz
Downloading http://yann.lecun.com/exdb/mnist/train-labels-idx1-ubyte.gz to datasets/mnist/MNIST/raw/train-labels-idx1-ubyt
e.gz
100%|████████| 28881/28881 [00:00<00:00, 39821069.63it/s]
Extracting datasets/mnist/MNIST/raw/train-labels-idx1-ubyte.gz to datasets/mnist/MNIST/raw

Downloading http://yann.lecun.com/exdb/mnist/t10k-images-idx3-ubyte.gz
Downloading http://yann.lecun.com/exdb/mnist/t10k-images-idx3-ubyte.gz to datasets/mnist/MNIST/raw/t10k-images-idx3-ubyte.g
z
100%|████████| 1648877/1648877 [00:00<00:00, 75522434.28it/s]
Extracting datasets/mnist/MNIST/raw/t10k-images-idx3-ubyte.gz to datasets/mnist/MNIST/raw

Downloading http://yann.lecun.com/exdb/mnist/t10k-labels-idx1-ubyte.gz
Downloading http://yann.lecun.com/exdb/mnist/t10k-labels-idx1-ubyte.gz to datasets/mnist/MNIST/raw/t10k-labels-idx1-ubyte.g
z
100%|████████| 4542/4542 [00:00<00:00, 21821911.53it/s]
Extracting datasets/mnist/MNIST/raw/t10k-labels-idx1-ubyte.gz to datasets/mnist/MNIST/raw
```

Загрузим датасет CIFAR-10 с параметрами cifar_mean = [0.491, 0.482, 0.447] cifar_std = [0.202, 0.199, 0.201] cifar_dim

```
cifar_mean = [0.491, 0.482, 0.447]
cifar_std = [0.202, 0.199, 0.201]
cifar_dim = 32

cifar_min, cifar_max = get_clip_bounds(cifar_mean, cifar_std, cifar_dim)

cifar_min = cifar_min.to(device)
cifar_max = cifar_max.to(device)

cifar_tf = transforms.Compose([ transforms.ToTensor(), transforms.Normalize( mean=cifar_mean, std=cifar_std)])

cifar_tf_train = transforms.Compose([ transforms.RandomCrop( size=cifar_dim, padding=4), transforms.RandomHorizontalFlip

cifar_tf_inv = transforms.Compose([ transforms.Normalize( mean=[0.0, 0.0, 0.0], std=np.divide(1.0, cifar_std)), transform

cifar_temp = datasets.CIFAR10(root='datasets/cifar-10', train=True, download=True, transform=cifar_tf_train)

cifar_train, cifar_val = random_split(cifar_temp, [40000, 10000])
cifar_test = datasets.CIFAR10(root='datasets/cifar-10', train=False, download=True, transform=cifar_tf)
cifar_classes = ['airplane', 'automobile', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship', 'truck']
```

```
Downloading https://www.cs.toronto.edu/~kriz/cifar-10-python.tar.gz to datasets/cifar-10/cifar-10-python.tar.gz
100%|███████████| 170498071/170498071 [00:05<00:00, 29134501.88it/s]
Extracting datasets/cifar-10/cifar-10-python.tar.gz to datasets/cifar-10
Files already downloaded and verified
```

Выполним настройку и загрузку DataLoader batch_size = 64 workers = 4

```
batch_size = 64
workers = 4

mnist_loader_train = DataLoader(mnist_train, batch_size=batch_size, shuffle=True, num_workers=workers)
mnist_loader_val = DataLoader(mnist_val, batch_size=batch_size, shuffle=False, num_workers=workers)
mnist_loader_test = DataLoader(mnist_test, batch_size=batch_size, shuffle=False, num_workers=workers)

cifar_loader_train = DataLoader(cifar_train, batch_size=batch_size, shuffle=True, num_workers=workers)
cifar_loader_val = DataLoader(cifar_val, batch_size=batch_size, shuffle=False, num_workers=workers)
cifar_loader_test = DataLoader(cifar_test, batch_size=batch_size, shuffle=False, num_workers=workers)
```

```
/usr/local/lib/python3.10/dist-packages/torch/utils/data/dataloader.py:557: UserWarning: This DataLoader will create 4 work
er processes in total. Our suggested max number of worker in current system is 2, which is smaller than what this DataLoade
r is going to create. Please be aware that excessive worker creation might get DataLoader running slow or even freeze, lowe
r the worker number to avoid potential slowness/freeze if necessary.
  warnings.warn(_create_warning_msg(
```

Настроим параметры для обучения

```
batch_size = 10
num_classes = 10
overshoot = 0.02
max_iters = 50
deep_args = [batch_size, num_classes, overshoot, max_iters]
```

Загрузим и оценим стойкость модели Network-In-Network Model к FGSM и DeepFool атакам на основе датасета CIFAR-10

```
fgsm_eps = 0.2
model = Net().to(device)
model.load_state_dict(torch.load('weights/clean/cifar_nin.pth', map_location=torch.device('cpu')))

evaluate_attack('cifar_nin_fgsm.csv', 'results', device, model, cifar_loader_test, cifar_min, cifar_max, fgsm_eps, is_fg
print('')
evaluate_attack('cifar_nin_deepfool.csv', 'results', device, model, cifar_loader_test, cifar_min, cifar_max, deep_args,

if device.type == 'cuda': torch.cuda.empty_cache()
```

```
FGSM Test Error : 81.29%
FGSM Robustness : 1.77e-01
FGSM Time (All Images) : 0.67 s
FGSM Time (Per Image) : 67.07 us

DeepFool Test Error : 93.76%
DeepFool Robustness : 2.12e-02
DeepFool Time (All Images) : 185.12 s
DeepFool Time (Per Image) : 18.51 ms
```

Загрузим и оценим стойкость модели LeNet к FGSM и DeepFool атакам на основе датасета CIFAR-10

```
fgsm_eps = 0.1
model = LeNet_CIFAR().to(device)
model.load_state_dict(torch.load('weights/clean/cifar_lenet.pth', map_location=torch.device('cpu')))

evaluate_attack('cifar_lenet_fgsm.csv', 'results', device, model, cifar_loader_test, cifar_min, cifar_max, fgsm_eps, is_
print('')
evaluate_attack('cifar_lenet_deepfool.csv', 'results', device, model, cifar_loader_test, cifar_min, cifar_max, deep_args

if device.type == 'cuda': torch.cuda.empty_cache()
```

```
FGSM Test Error : 91.71%
FGSM Robustness : 8.90e-02
FGSM Time (All Images) : 0.40 s
FGSM Time (Per Image) : 40.08 us

DeepFool Test Error : 87.81%
DeepFool Robustness : 1.78e-02
DeepFool Time (All Images) : 73.27 s
DeepFool Time (Per Image) : 7.33 ms
```

Выполним оценку атакующих примеров для сетей:

## LeNet на MNIST

```
fgsm_eps = 0.6
model = LeNet_MNIST().to(device)
model.load_state_dict(torch.load('weights/clean/mnist_lenet.pth', map_location=device))

display_attack(device, model, mnist_test, mnist_tf_inv, mnist_min, mnist_max, fgsm_eps, deep_args, has_labels=False, l2_

if device.type == 'cuda': torch.cuda.empty_cache()
```



## FCNet на MNIST

```
fgsm_eps = 0.2
model = FC_500_150().to(device)
model.load_state_dict(torch.load('weights/clean/mnist_fc.pth'))

display_attack(device, model, mnist_test, mnist_tf_inv, mnist_min, mnist_max, fgsm_eps, deep_args, has_labels=False, l2_

if device.type == 'cuda': torch.cuda.empty_cache()
```

## Network-in-Network на CIFAR

```python
fgsm_eps = 0.2
model = Net().to(device)
model.load_state_dict(torch.load('weights/clean/cifar_nin.pth'))

display_attack(device, model, cifar_test, cifar_tf_inv, cifar_min, cifar_max, fgsm_eps, deep_args, has_labels=False, l2_

if device.type == 'cuda': torch.cuda.empty_cache()
```
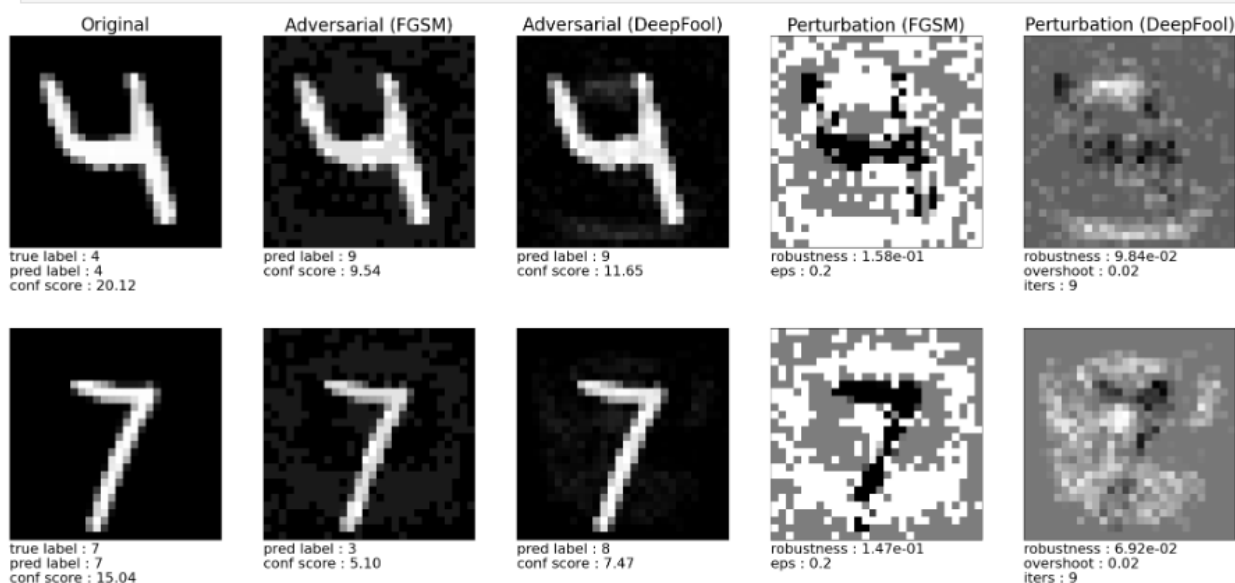


| Original | Adversarial (FGSM) | Adversarial (DeepFool) | Perturbation (FGSM) | Perturbation (DeepFool) |
|---|---|---|---|---|
| true label : bird<br>pred label : bird<br>conf score : 21.25 | pred label : frog<br>conf score : 26.86 | pred label : frog<br>conf score : 15.81 | robustness : 1.63e-01<br>eps : 0.2 | robustness : 1.34e-02<br>overshoot : 0.02<br>iters : 2 |
| true label : bird<br>pred label : bird<br>conf score : 20.41 | pred label : cat<br>conf score : 21.60 | pred label : cat<br>conf score : 18.37 | robustness : 3.40e-01<br>eps : 0.2 | robustness : 1.06e-02<br>overshoot : 0.02<br>iters : 2 |

## LeNet на CIFAR

```python
fgsm_eps = 0.1
model = LeNet_CIFAR().to(device)
model.load_state_dict(torch.load('weights/clean/cifar_lenet.pth'))

display_attack(device, model, cifar_test, cifar_tf_inv, cifar_min, cifar_max, fgsm_eps, deep_args, has_labels=False, l2_

if device.type == 'cuda': torch.cuda.empty_cache()
```



| Original | Adversarial (FGSM) | Adversarial (DeepFool) | Perturbation (FGSM) | Perturbation (DeepFool) |
|---|---|---|---|---|
| true label : horse<br>pred label : horse<br>conf score : 6.90 | pred label : dog<br>conf score : 8.78 | pred label : dog<br>conf score : 5.46 | robustness : 8.48e-02<br>eps : 0.1 | robustness : 1.19e-02<br>overshoot : 0.02<br>iters : 3 |
| true label : deer<br>pred label : deer<br>conf score : 6.09 | pred label : airplane<br>conf score : 4.36 | pred label : airplane<br>conf score : 4.66 | robustness : 1.44e-01<br>eps : 0.1 | robustness : 1.90e-02<br>overshoot : 0.02<br>iters : 2 |

Отразим отличия для fgsm_eps = (0.001, 0.02, 0.5, 0.9, 10)

```python
fgsm_eps_arr = [0.001, 0.02, 0.5, 0.9, 10]

for fgsm_eps in fgsm_eps_arr:
  model = LeNet_MNIST().to(device)
  model.load_state_dict(torch.load('weights/clean/mnist_lenet.pth'))

  display_attack(device, model, mnist_test, mnist_tf_inv, mnist_min, mnist_max, fgsm_eps, deep_args, has_labels=False, l

  if device.type == 'cuda': torch.cuda.empty_cache()
  print("eps: ", fgsm_eps)
```
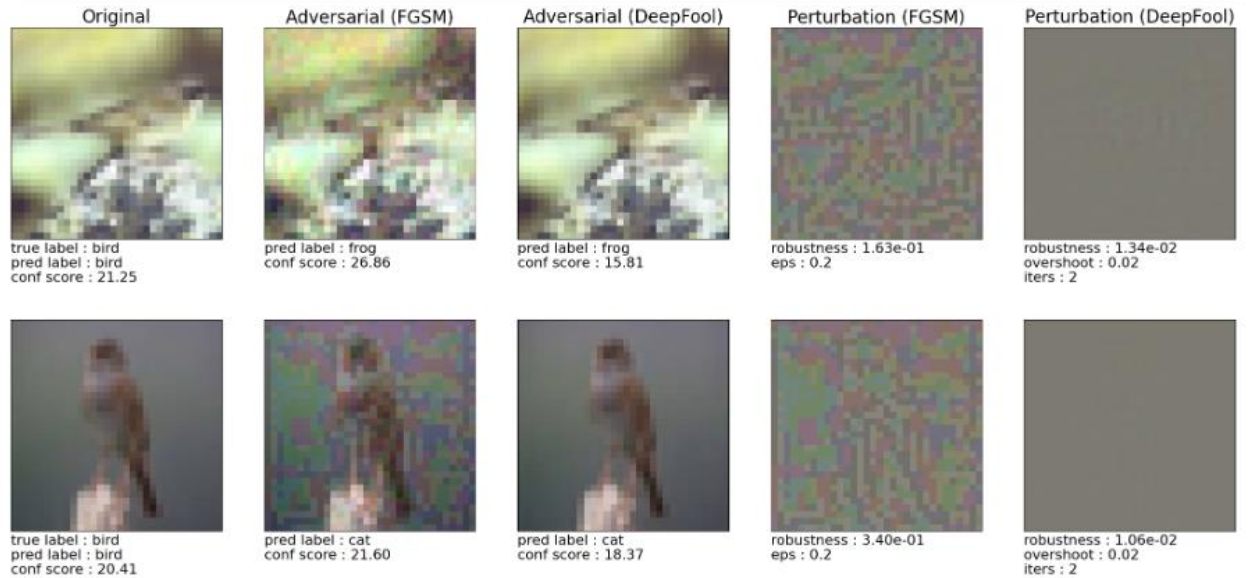
LeNet на MNIST, eps: 0,001



eps:  0.001

LeNet на MNIST, eps: 0,02



eps:  0.02

LeNet на MNIST, eps: 0,5



LeNet на MNIST, eps: 0,9

# LeNet на MNIST, eps: 10



# Network-in-Network на CIFAR, eps: 0,001

```
for fgsm_eps in fgsm_eps_arr:
    model = Net().to(device)
    model.load_state_dict(torch.load('weights/clean/cifar_nin.pth'))

    display_attack(device, model, cifar_test, cifar_tf_inv, cifar_min, cifar_max, fgsm_eps, deep_args, has_labels=False, l

    if device.type == 'cuda': torch.cuda.empty_cache()
    print("eps: ", fgsm_eps)
```
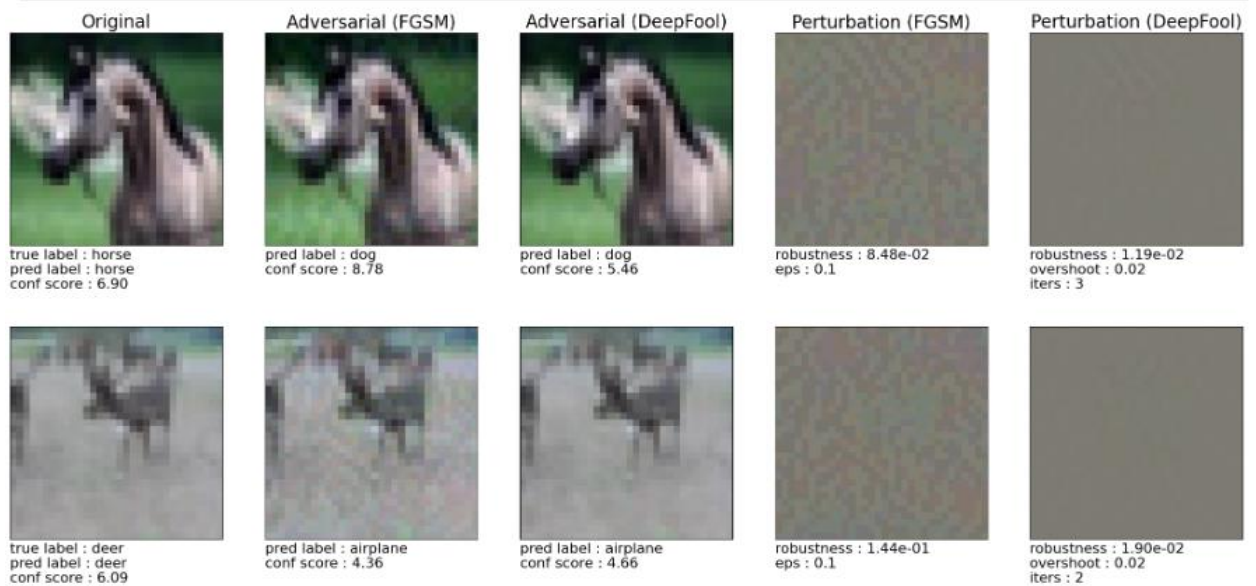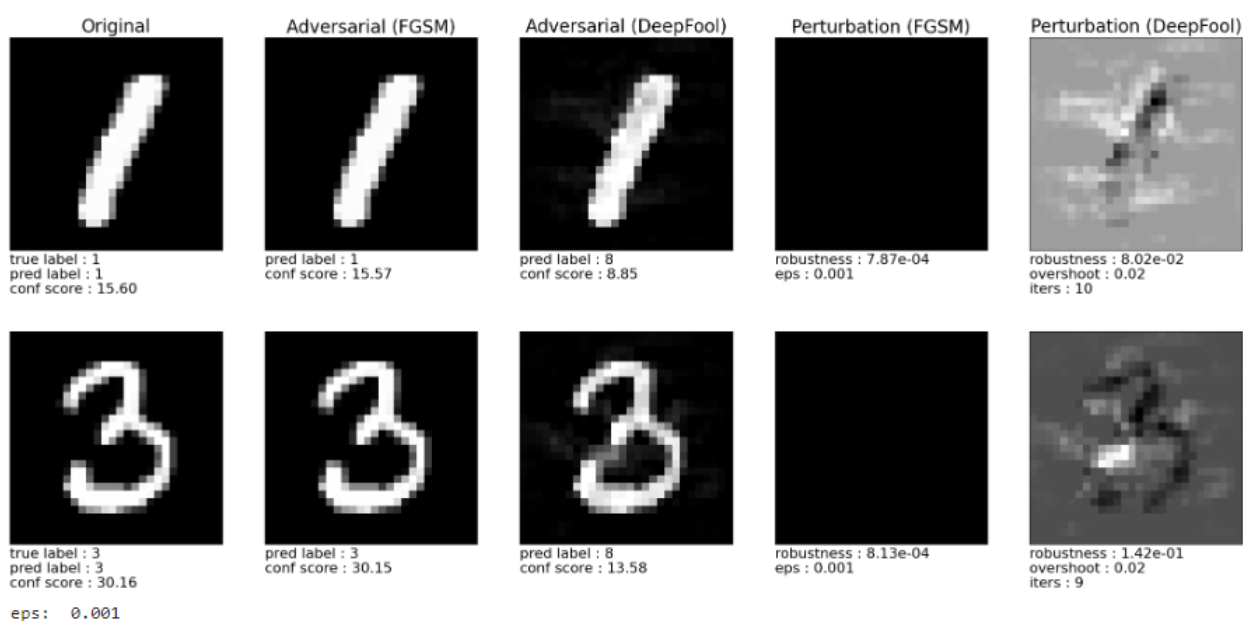
## Network-in-Network на CIFAR, eps: 0,02



| Original | Adversarial (FGSM) | Adversarial (DeepFool) | Perturbation (FGSM) | Perturbation (DeepFool) |
|---|---|---|---|---|
| true label : dog<br>pred label : dog<br>conf score : 19.64 | pred label : frog<br>conf score : 19.56 | pred label : frog<br>conf score : 18.31 | robustness : 2.21e-02<br>eps : 0.02 | robustness : 1.01e-02<br>overshoot : 0.02<br>iters : 2 |
| true label : dog<br>pred label : dog<br>conf score : 27.50 | pred label : cat<br>conf score : 26.52 | pred label : cat<br>conf score : 25.71 | robustness : 1.41e-02<br>eps : 0.02 | robustness : 4.71e-03<br>overshoot : 0.02<br>iters : 2 |

eps:  0.02

## Network-in-Network на CIFAR, eps: 0,05



| Original | Adversarial (FGSM) | Adversarial (DeepFool) | Perturbation (FGSM) | Perturbation (DeepFool) |
|---|---|---|---|---|
| true label : ship<br>pred label : ship<br>conf score : 31.34 | pred label : frog<br>conf score : 17.26 | pred label : airplane<br>conf score : 23.31 | robustness : 3.95e-01<br>eps : 0.5 | robustness : 2.05e-02<br>overshoot : 0.02<br>iters : 3 |
| true label : airplane<br>pred label : airplane<br>conf score : 41.57 | pred label : bird<br>conf score : 17.84 | pred label : bird<br>conf score : 29.17 | robustness : 4.22e-01<br>eps : 0.5 | robustness : 6.76e-02<br>overshoot : 0.02<br>iters : 5 |

eps:  0.5

# Network-in-Network на CIFAR, eps: 0,09



| Original | Adversarial (FGSM) | Adversarial (DeepFool) | Perturbation (FGSM) | Perturbation (DeepFool) |
|---|---|---|---|---|
| true label : airplane<br>pred label : airplane<br>conf score : 14.73 | pred label : bird<br>conf score : 12.18 | pred label : deer<br>conf score : 13.59 | robustness : 5.84e-01<br>eps : 0.9 | robustness : 7.48e-03<br>overshoot : 0.02<br>iters : 2 |
| true label : deer<br>pred label : deer<br>conf score : 32.04 | pred label : frog<br>conf score : 14.11 | pred label : bird<br>conf score : 23.79 | robustness : 1.44e+00<br>eps : 0.9 | robustness : 2.43e-02<br>overshoot : 0.02<br>iters : 2 |

eps:  0.9

# Network-in-Network на CIFAR, eps: 10



| Original | Adversarial (FGSM) | Adversarial (DeepFool) | Perturbation (FGSM) | Perturbation (DeepFool) |
|---|---|---|---|---|
| true label : frog<br>pred label : frog<br>conf score : 32.37 | pred label : frog<br>conf score : 20.51 | pred label : deer<br>conf score : 20.98 | robustness : 3.24e+00<br>eps : 10 | robustness : 3.15e-02<br>overshoot : 0.02<br>iters : 4 |
| true label : dog<br>pred label : dog<br>conf score : 50.73 | pred label : frog<br>conf score : 22.93 | pred label : cat<br>conf score : 30.60 | robustness : 1.99e+00<br>eps : 10 | robustness : 3.11e-02<br>overshoot : 0.02<br>iters : 3 |

eps:  10

# FCNet на MNIST, eps: 0,001

```
for fgsm_eps in fgsm_eps_arr:
    model = FC_500_150().to(device)
    model.load_state_dict(torch.load('weights/clean/mnist_fc.pth'))

    display_attack(device, model, mnist_test, mnist_tf_inv, mnist_min, mnist_max, fgsm_eps, deep_args, has_labels=False, l

    if device.type == 'cuda': torch.cuda.empty_cache()
    print("eps: ", fgsm_eps)
```
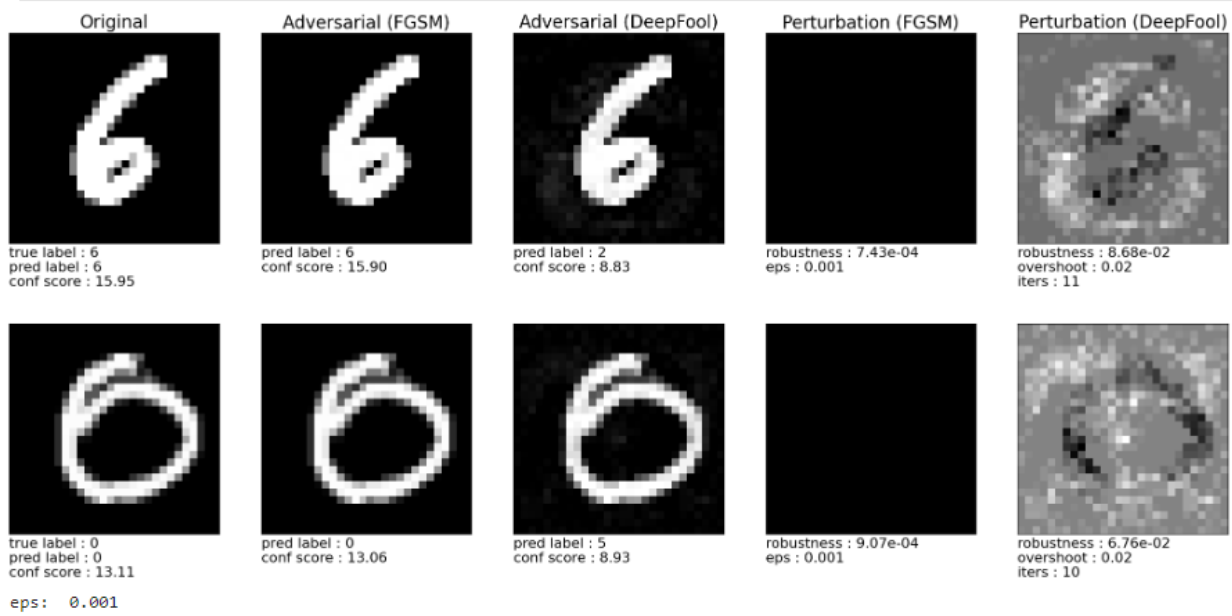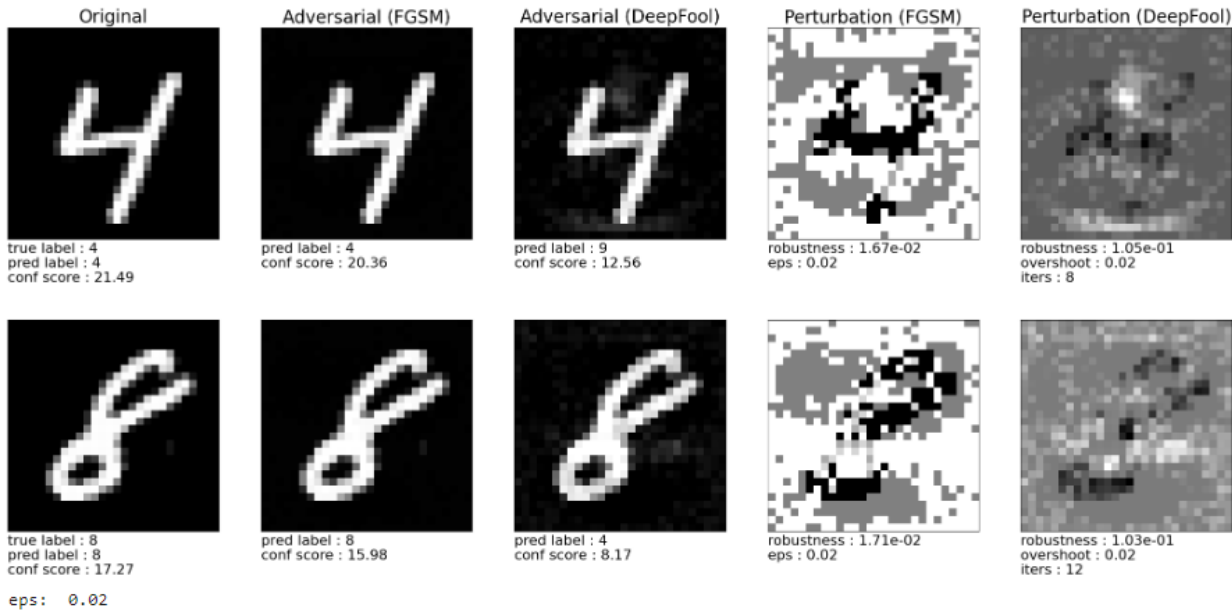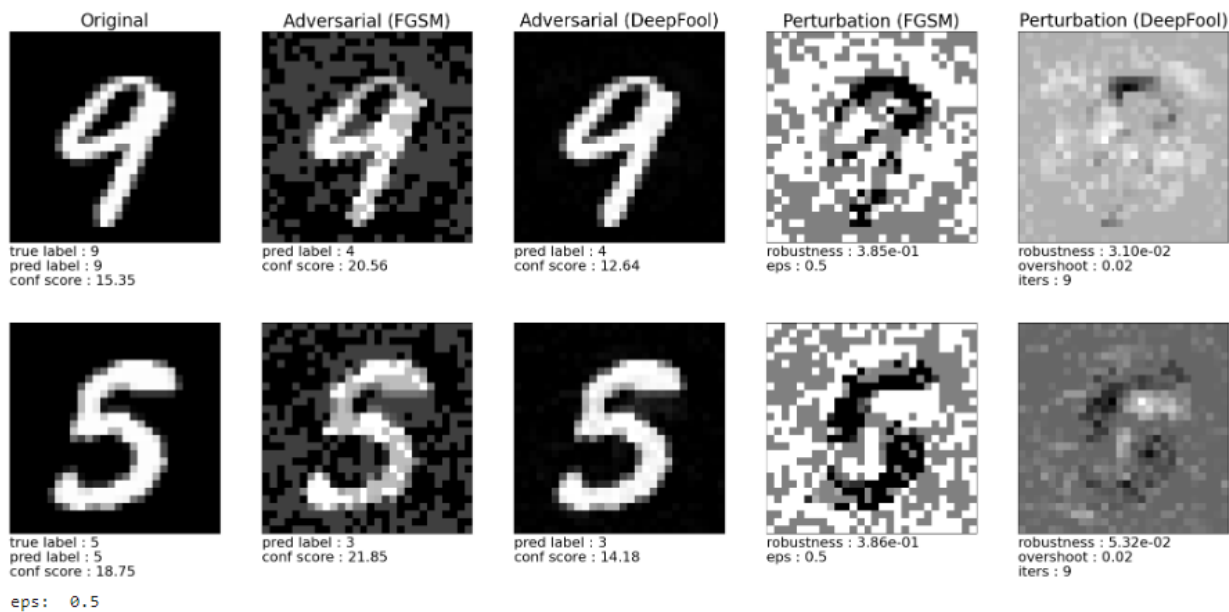


| Original | Adversarial (FGSM) | Adversarial (DeepFool) | Perturbation (FGSM) | Perturbation (DeepFool) |

true label : 6
pred label : 6
conf score : 15.95

pred label : 6
conf score : 15.90

pred label : 2
conf score : 8.83

robustness : 7.43e-04
eps : 0.001

robustness : 8.68e-02
overshoot : 0.02
iters : 11

true label : 0
pred label : 0
conf score : 13.11

pred label : 0
conf score : 13.06

pred label : 5
conf score : 8.93

robustness : 9.07e-04
eps : 0.001

robustness : 6.76e-02
overshoot : 0.02
iters : 10

eps:  0.001

# FCNet на MNIST, eps: 0,02



| Original | Adversarial (FGSM) | Adversarial (DeepFool) | Perturbation (FGSM) | Perturbation (DeepFool) |

true label : 4
pred label : 4
conf score : 21.49

pred label : 4
conf score : 20.36

pred label : 9
conf score : 12.56

robustness : 1.67e-02
eps : 0.02

robustness : 1.05e-01
overshoot : 0.02
iters : 8

true label : 8
pred label : 8
conf score : 17.27

pred label : 8
conf score : 15.98

pred label : 4
conf score : 8.17

robustness : 1.71e-02
eps : 0.02

robustness : 1.03e-01
overshoot : 0.02
iters : 12

eps:  0.02

FCNet на MNIST, eps: 0,5



FCNet на MNIST, eps: 0,9

FCNet на MNIST, eps: 10



| Original | Adversarial (FGSM) | Adversarial (DeepFool) | Perturbation (FGSM) | Perturbation (DeepFool) |
|---|---|---|---|---|
| true label : 0<br>pred label : 0<br>conf score : 11.73 | pred label : 2<br>conf score : 16.79 | pred label : 4<br>conf score : 7.84 | robustness : 1.50e+00<br>eps : 10 | robustness : 8.10e-02<br>overshoot : 0.02<br>iters : 9 |
| true label : 9<br>pred label : 9<br>conf score : 9.92 | pred label : 7<br>conf score : 37.05 | pred label : 7<br>conf score : 7.76 | robustness : 1.50e+00<br>eps : 10 | robustness : 2.79e-02<br>overshoot : 0.02<br>iters : 10 |

eps: 10

**Вывод.**

Наблюдая результат эксперимента, можно сделать вывод, что при увеличении параметра fgsm_eps увеличивается количество шума на изображениях. Это говорит о том, что модель становиться более подвержена ошибкам во время работы, а ее степень устойчивости к атакам меньше, чем если бы был более низкий показатель fgsm_eps.