# Лабораторная работа №1

По дисциплине
«Анализ защищенности систем искусственного интеллекта»

Выполнил:

Суслов Антон Константинович

Группа: ББМО-02-22

Москва 2023

Ход работы:

1. Клонирование репозитория

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

```
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 | 9.79 MiB/s, done.
Resolving deltas: 100% (27/27), done.
```

2. Загрузка библиотек и установка значения переменной rand_seed(номер по списку)

```
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
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
```

```
rand_seed = 40
np.random.seed(rand_seed)
torch.manual_seed (rand_seed)
use_cuda = torch.cuda.is_available()
device = torch.device('cuda' if use_cuda else 'cpu')
```

3. Загрузка датасетов MNIST и Cifar-10

```
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-ubyte.gz
100%|██████████| 9912422/9912422 [00:00<00:00, 348637433.71it/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-ubyte.gz
100%|██████████| 28881/28881 [00:00<00:00, 24135424.15it/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.gz
100%|██████████| 1648877/1648877 [00:00<00:00, 162918525.24it/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.gz
100%|██████████| 4542/4542 [00:00<00:00, 20074319.04it/s]
Extracting datasets/mnist/MNIST/raw/t10k-labels-idx1-ubyte.gz to datasets/mnist/MNIST/raw
```

```
        transforms.Normalize(
            mean=np.multiply(-1.0, cifar_mean),
            std=[1.0, 1.0, 1.0])])

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)
```

```
Downloading https://www.cs.toronto.edu/~kriz/cifar-10-python.tar.gz to datasets/cifar-10/cifar-10-python.tar.gz
100%|██████████| 170498071/170498071 [00:14<00:00, 12085696.50it/s]
Extracting datasets/cifar-10/cifar-10-python.tar.gz to datasets/cifar-10
Files already downloaded and verified
```

## 4. Настройка и загрузка DataLoader

```
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 worker proce
  warnings.warn(_create_warning_msg(
```

## 5. Настройка модели

```python
import os
train_model = True

epochs = 50
epochs_nin = 100

lr = 0.004
lr_nin = 0.01
lr_scale = 0.5

momentum = 0.9

print_step = 5

deep_batch_size = 64
deep_num_classes = 10
deep_overshoot = 0.02
deep_max_iters = 50

deep_args = [deep_batch_size, deep_num_classes,
             deep_overshoot, deep_max_iters]

if not os.path.isdir('weights/deepfool'):
    os.makedirs('weights/deepfool', exist_ok=True)

if not os.path.isdir('weights/fgsm'):
    os.makedirs('weights/fgsm', exist_ok=True)
```
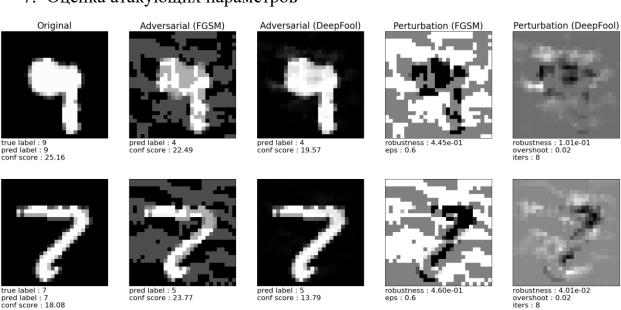
6. Загрузка и оценка стойкости модели LeNet к FGSM и DeepFool атакам, также загрузка и оценка стойкости модели LeNet к FGSM и DeepFool атакам
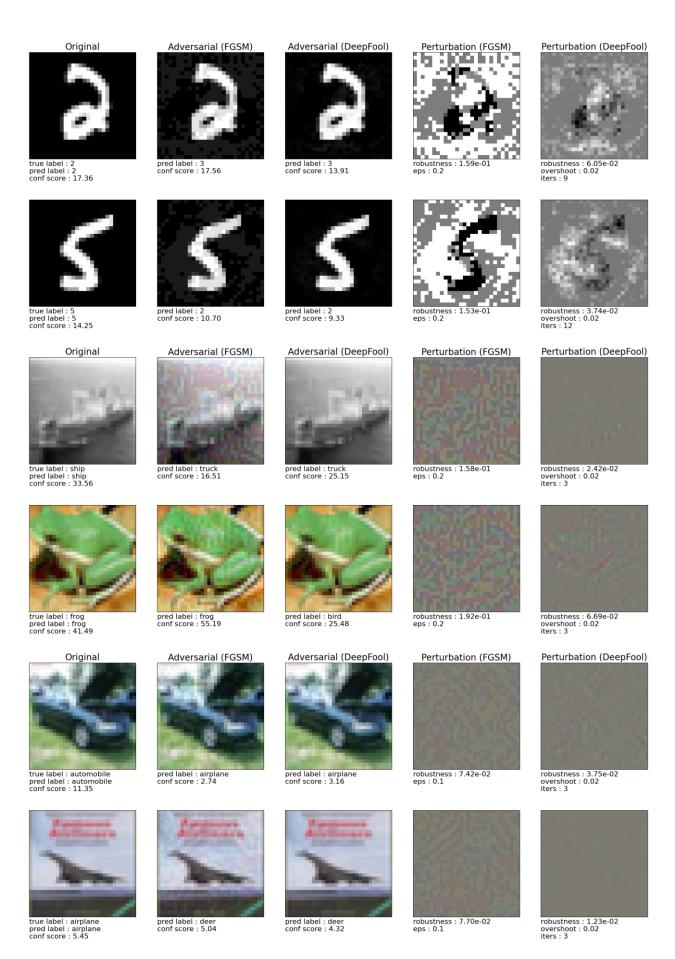
```python
fgsm_eps = 0.6
model = LeNet_MNIST().to(device)
model.load_state_dict(torch.load('weights/clean/mnist_lenet.pth',map_location=torch.device('cpu')))
evaluate_attack('mnist_lenet_fgsm.csv',
                'results', device, model, mnist_loader_test,
                mnist_min, mnist_max,fgsm_eps, is_fgsm=True)
print('')

evaluate_attack('mnist_lenet_deepfool.csv', 'results', device, model,
mnist_loader_test, mnist_min, mnist_max, deep_args, is_fgsm=False)
if device.type == 'cuda': torch.cuda.empty_cache()
```

```
FGSM Test Error : 87.89%
FGSM Robustness : 4.58e-01
FGSM Time (All Images) : 0.29 s
FGSM Time (Per Image) : 28.86 us

DeepFool Test Error : 98.74%
DeepFool Robustness : 9.64e-02
DeepFool Time (All Images) : 193.32 s
DeepFool Time (Per Image) : 19.33 ms
```

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

evaluate_attack('mnist_fc_fgsm.csv', 'results', device, model,
                mnist_loader_test, mnist_min, mnist_max, fgsm_eps, is_fgsm=True)
print('')

evaluate_attack('mnist_fc_deepfool.csv', 'results', device, model,
mnist_loader_test, mnist_min, mnist_max, deep_args, is_fgsm=False)
if device.type == 'cuda': torch.cuda.empty_cache()
```

```
FGSM Test Error : 87.08%
FGSM Robustness : 1.56e-01
FGSM Time (All Images) : 0.15 s
FGSM Time (Per Image) : 14.99 us

DeepFool Test Error : 97.92%
DeepFool Robustness : 6.78e-02
DeepFool Time (All Images) : 141.81 s
DeepFool Time (Per Image) : 14.18 ms
```

## 7. Оценка атакующих параметров



| Original | Adversarial (FGSM) | Adversarial (DeepFool) | Perturbation (FGSM) | Perturbation (DeepFool) |
|---|---|---|---|---|
| true label : 9<br>pred label : 9<br>conf score : 25.16 | pred label : 4<br>conf score : 22.49 | pred label : 4<br>conf score : 19.57 | robustness : 4.45e-01<br>eps : 0.6 | robustness : 1.01e-01<br>overshoot : 0.02<br>iters : 8 |
| true label : 7<br>pred label : 7<br>conf score : 18.08 | pred label : 5<br>conf score : 23.77 | pred label : 5<br>conf score : 13.79 | robustness : 4.60e-01<br>eps : 0.6 | robustness : 4.01e-02<br>overshoot : 0.02<br>iters : 8 |

| Original | Adversarial (FGSM) | Adversarial (DeepFool) | Perturbation (FGSM) | Perturbation (DeepFool) |
|---|---|---|---|---|
| true label : 2<br>pred label : 2<br>conf score : 17.36 | pred label : 3<br>conf score : 17.56 | pred label : 3<br>conf score : 13.91 | robustness : 1.59e-01<br>eps : 0.2 | robustness : 6.05e-02<br>overshoot : 0.02<br>iters : 9 |
| true label : 5<br>pred label : 5<br>conf score : 14.25 | pred label : 2<br>conf score : 10.70 | pred label : 2<br>conf score : 9.33 | robustness : 1.53e-01<br>eps : 0.2 | robustness : 3.74e-02<br>overshoot : 0.02<br>iters : 12 |
| true label : ship<br>pred label : ship<br>conf score : 33.56 | pred label : truck<br>conf score : 16.51 | pred label : truck<br>conf score : 25.15 | robustness : 1.58e-01<br>eps : 0.2 | robustness : 2.42e-02<br>overshoot : 0.02<br>iters : 3 |
| true label : frog<br>pred label : frog<br>conf score : 41.49 | pred label : frog<br>conf score : 55.19 | pred label : bird<br>conf score : 25.48 | robustness : 1.92e-01<br>eps : 0.2 | robustness : 6.69e-02<br>overshoot : 0.02<br>iters : 3 |
| true label : automobile<br>pred label : automobile<br>conf score : 11.35 | pred label : airplane<br>conf score : 2.74 | pred label : airplane<br>conf score : 3.16 | robustness : 7.42e-02<br>eps : 0.1 | robustness : 3.75e-02<br>overshoot : 0.02<br>iters : 3 |
| true label : airplane<br>pred label : airplane<br>conf score : 5.45 | pred label : deer<br>conf score : 5.04 | pred label : deer<br>conf score : 4.32 | robustness : 7.70e-02<br>eps : 0.1 | robustness : 1.23e-02<br>overshoot : 0.02<br>iters : 3 |

8. Влияние параметров fgsm_esp для LeNet на датасетах MNIST и Cifar-10

```
Evaluating FGSM Attack with eps=0.001...
/usr/local/lib/python3.10/dist-packages/torch/utils/data/dataloader.py:557: UserWarning: This DataLoader will create 4 worker processes in total. Our suggested max number
    warnings.warn(_create_warning_msg(
FGSM Batches Complete : (157 / 157)
FGSM Test Error : 3.07%
FGSM Robustness : 8.08e-04
FGSM Time (All Images) : 0.56 s
FGSM Time (Per Image) : 55.87 us
/usr/local/lib/python3.10/dist-packages/torch/utils/data/dataloader.py:557: UserWarning: This DataLoader will create 4 worker processes in total. Our suggested max number
    warnings.warn(_create_warning_msg(
FGSM Batches Complete : (157 / 157)
FGSM Test Error : 10.12%
FGSM Robustness : 8.92e-04
FGSM Time (All Images) : 1.25 s
FGSM Time (Per Image) : 124.84 us
Evaluating FGSM Attack with eps=0.02...
/usr/local/lib/python3.10/dist-packages/torch/utils/data/dataloader.py:557: UserWarning: This DataLoader will create 4 worker processes in total. Our suggested max number
    warnings.warn(_create_warning_msg(
FGSM Batches Complete : (157 / 157)
FGSM Test Error : 5.54%
FGSM Robustness : 1.60e-02
FGSM Time (All Images) : 0.51 s
FGSM Time (Per Image) : 51.38 us
/usr/local/lib/python3.10/dist-packages/torch/utils/data/dataloader.py:557: UserWarning: This DataLoader will create 4 worker processes in total. Our suggested max number
    warnings.warn(_create_warning_msg(
FGSM Batches Complete : (157 / 157)
FGSM Test Error : 30.76%
FGSM Robustness : 1.78e-02
FGSM Time (All Images) : 1.10 s
FGSM Time (Per Image) : 109.76 us
```

Вывод: параметр fgsm_esp влияет на искажение изображения, чем выше параметр, тем искажение больше.