


Assignment 2: Adversarial Attack on NN

(Due: 2023/10/01 Sunday 23:59)

q1. (30%) Please read the **fast gradient sign method** (FGSM) paper and write pseudocode similar to the following format

Pseudocode Format Example :

Algorithm 1 Backpropagation Algorithm	
--	---

```
1: procedure TRAIN
2:    $X \leftarrow$  Training Data Set of size  $m \times n$ 
3:    $y \leftarrow$  Labels for records in  $X$ 
4:    $w \leftarrow$  The weights for respective layers
5:    $l \leftarrow$  The number of layers in the neural network,  $1 \dots L$ 
6:    $D_{ij}^{(l)} \leftarrow$  The error for all  $i, j$ 
7:    $t_{ij}^{(l)} \leftarrow 0$ . For all  $i, j$ 
8:   For  $i = 1$  to  $m$ 
9:      $a^l \leftarrow \text{feedforward}(x^{(i)}, w)$ 
10:     $d^l \leftarrow a(L) - y(i)$ 
11:     $t_{ij}^{(l)} \leftarrow t_{ij}^{(l)} + a_j^{(l)} \cdot t_i^{l+1}$ 
12:    if  $j \neq 0$  then
13:       $D_{ij}^{(l)} \leftarrow \frac{1}{m} t_{ij}^{(l)} + \lambda w_{ij}^{(l)}$ 
14:    else
15:       $D_{ij}^{(l)} \leftarrow \frac{1}{m} t_{ij}^{(l)}$ 
16:    where  $\frac{\partial}{\partial w_{ij}^{(l)}} J(w) = D_{ij}^{(l)}$ 
```

q2. (40%) FGSM Attack

1. Download the model weights and model reading program from [here](#)
2. Attack this model by FGSM
3. Download your testing data from [here](#)
4. To convert testing data to images and save them as JPG files (.jpg).
5. Experiment with FGSM Attack using testing data, and record the results along with the parameters used.

The **noise magnitude** of q2 will **not affect** the score, please provide a detailed description of parameter settings and implementation in **q3**.

q3. (15%) Write down your experiment setting in English. The setting should include but not be limited to (1) hardware specification, (2) package version and (3) all the experiment parameters and details in **q2**.

The font size is 12, and the page limit is 2 pages.

Submission Guideline

Please compress your files named {SID}_a2.zip (SID in upper case) to the COOL System, such as D111111_a2.zip, with two required files

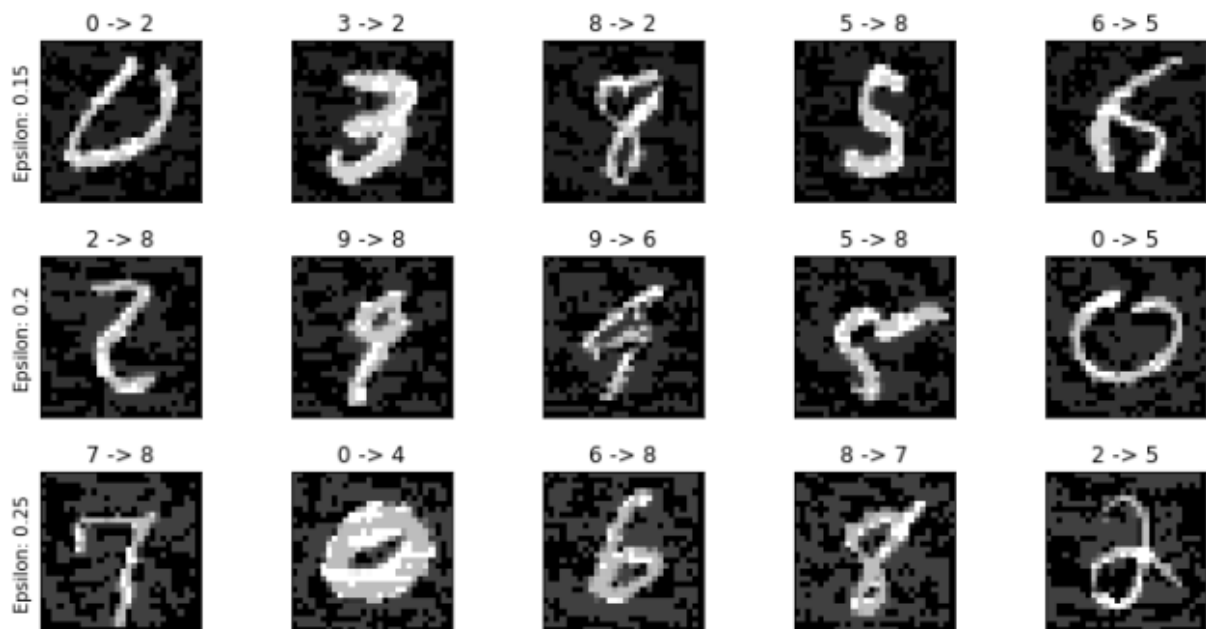
file 1. {SID}_a2.py

Please submit your source code to the COOL system. The following command can exclude the code

The execution results of {SID}_a2.py

The output of your code is the generated attacked images (testing data original image + the noise), ground truth label and predict of FGSM attack.

Output Format Example



files 2. {SID}_a2_report.pdf

Supplementary Materials

PyTorch installation: <https://pytorch.org/>
FGSM paper [1412.6572.pdf \(arxiv.org\)](https://arxiv.org/abs/1412.6572)