

DEFENCE AGAINST ADVERSARIAL EXAMPLES

Yishay Asher • Steve Gutfreund
Instructor: Hanan Rosemarin



אוניברסיטת בר-אילן
Bar-Ilan University

Problem Description

Building high accuracy DNN models which are
sufficiently resistant to adversarial attacks

Background

An adversarial example is an instance with small,
intentional feature perturbations that cause a
machine learning model to make a false prediction.
(See figure 1)

Goal

Find a way to train 'secured' models such that this
sort of attacks should not affect them.

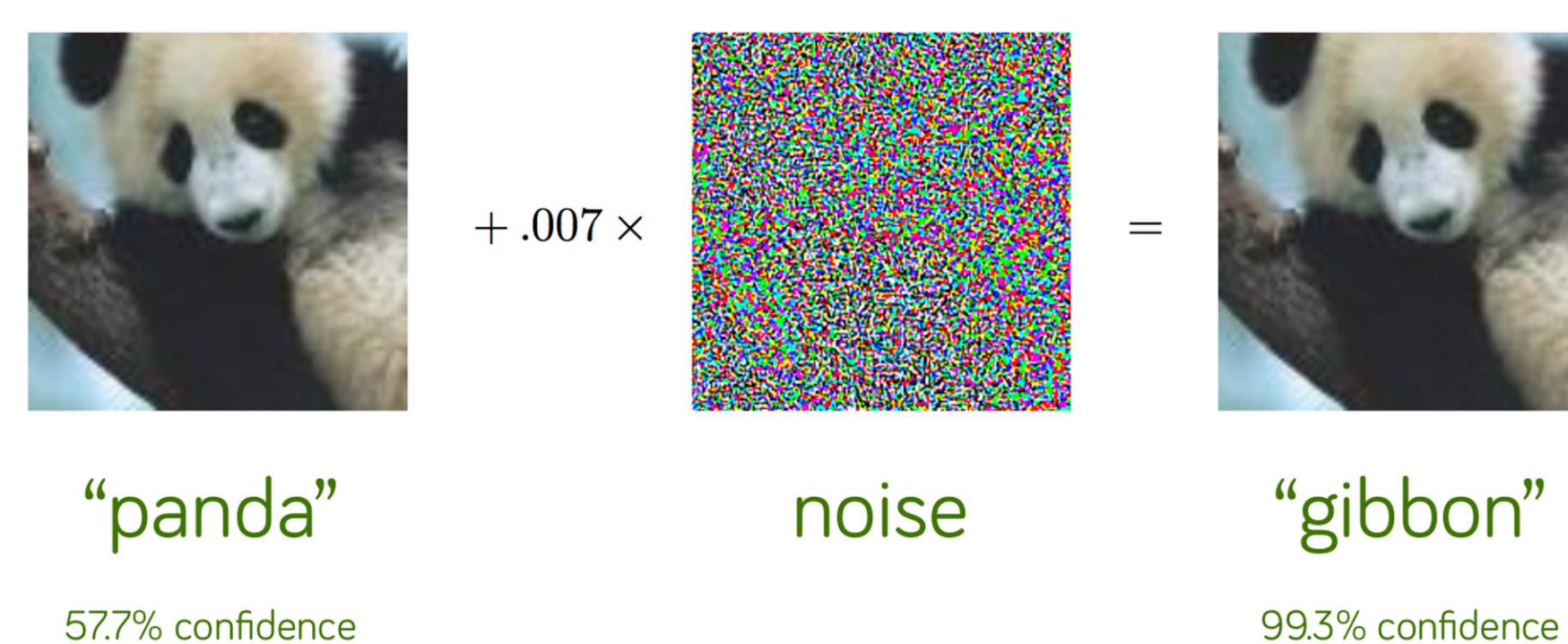


Figure 1: example of an adversarial example

Set-Up

We worked on mnist and fashion-mnist.
We've set neural networks known to be able to
learn these datasets very well.

1. Securing Models

We followed the approach presented in the article
'[Bridging machine learning and cryptography in
defence against adversarial attacks](#)':

training models on encrypted images, see figure 2.
Encryption techniques:

- Permutation on the pixels, as done in the article
- AES in ECB, CBC and CTR modes of operation

2. Cutting loose ends

Before performing an attack, we eliminated the
models that did not learn well.

As can be seen in figure 3, learning is not so
intuitive.

For this reason, AES in ECB and CBC mode were
irrelevant to continue with.

3. Attacking

Attacking the sufficiently accurate models with the
following attacks:

- Carlini & Wagner, CW
https://github.com/carlini/nn_robust_attacks
- Fast Gradient Sign Method, FGSM
<https://github.com/tensorflow/cleverhans>

We focused on the 'gray-box' scenario, i.e. the
attacker knows the architecture of the model but
has no access to the private key.
See figure 4 for a visualization.

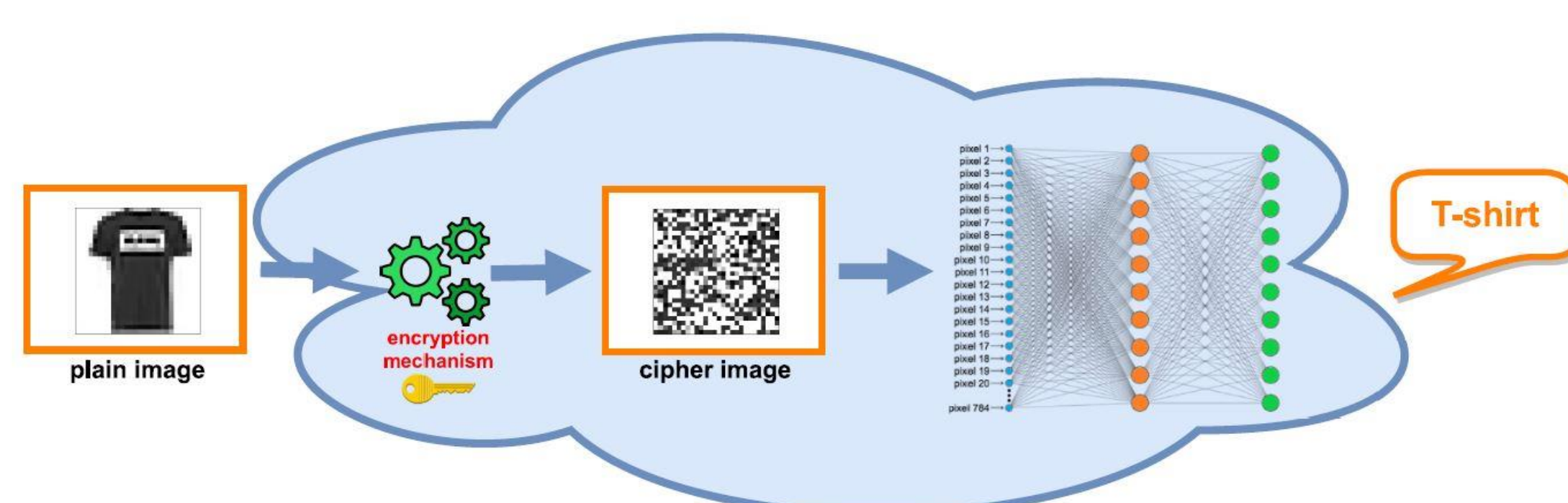


Figure 2: architecture of the secure models

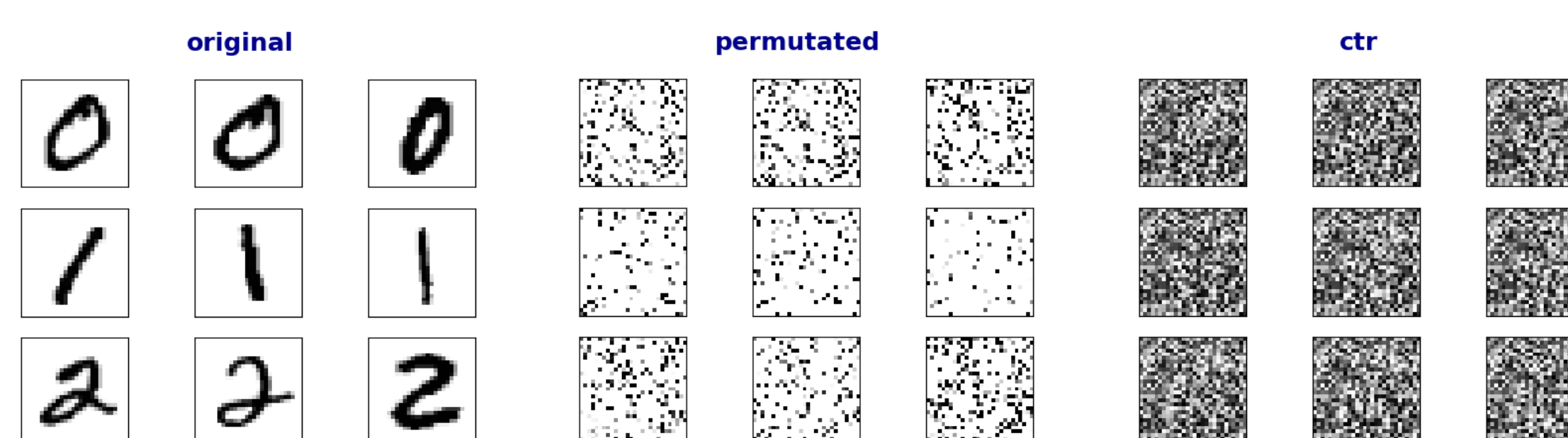


Figure 3: Sample of the encrypted images (permuted and
aes-ctr). Interesting to see how for the human eyes it's
impossible to distinguish between various classes but a DNN
model classifies quite well, see figure 5 for accuracies

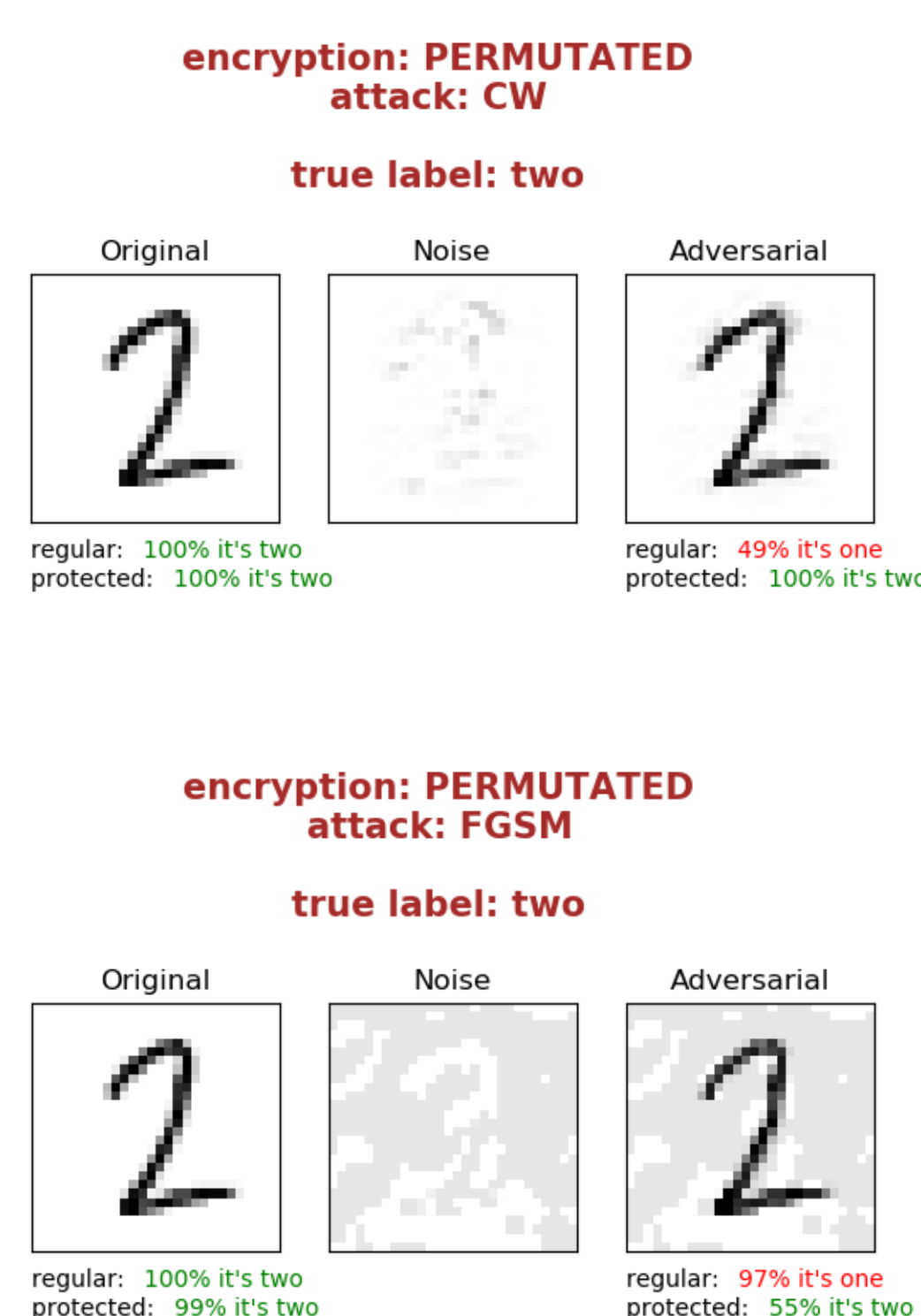


Figure 4: visualization of a CW and FGSM attack

Results

There's a tradeoff between accuracy on the original
images and the accuracy on the adversarial images.

- permutation: although the error rate on the
originals increased by a little bit (1% on mnist
and 4% on fashion-mnist), the error rate on the
adversarial decreased significantly
- AES-CTR: on mnist it performs better than
permutation but yet on fashion mnist the error
rate on the adversarials is slightly higher
See figure 5 for the detailed results.

Classification error (%) on the first 1000 test samples						
model	original images	mnist		fashion_mnist		
		adversarial images	gray-box	original images	adversarial images	gray-box
UNENCRYPTED	A	CW I ₂	100.00	8.30	CW I ₂	100.00
		CW I ₀	100.00		CW I ₀	100.00
		CW I _∞	100.00		CW I _∞	100.00
PERMUTATED	A	CW I ₂	4.50	12.30	CW I ₂	12.70
		CW I ₀	7.30		CW I ₀	12.50
		CW I _∞	5.40		CW I _∞	12.90
AES · ECB	B	FGSM	8.60	12.00	FGSM	29.80
AES · CBC	A	CW I ₂	irrelevant	71.50	CW I ₂	irrelevant
AES · CTR	A	CW I ₂	4.20	17.40	CW I ₂	17.20

Figure 5: table containing all the results

Future Work

- improve accuracy on AES-ECB model (we got
error rate of 19% on mnist and 55% on fashion-
mnist)
- we contacted Nicholas Carlini (the 'C' in CW
attack) and he believes we still might defeat
these defenses
- try some other datasets; i.e. cifar-10, its images
are 3 layered (rgb) and might be more difficult to
learn encrypted images