

Problem Description

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

Set-Up

Training models on encrypted images, on the datasets
mnist and fashion-mnist, see figure 1.

Encryption techniques:

- Permutation (on the pixels)
- AES in ECB, CBC and CTR modes of operation

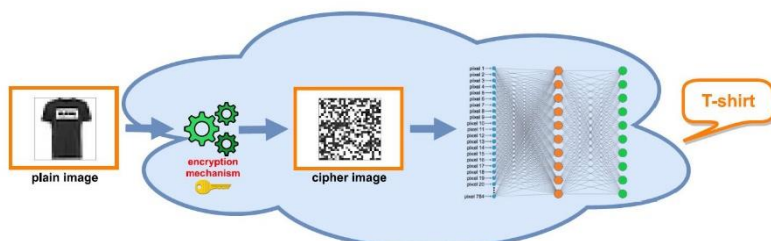


Figure 1: Architecture of the models

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>

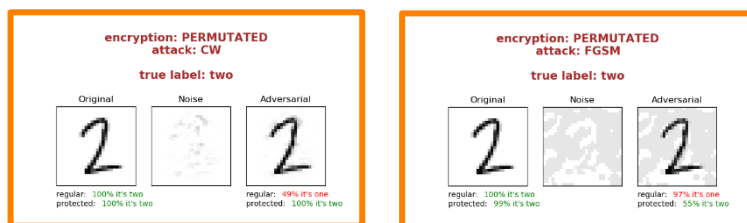


Figure 2: visualization of an attack

Results

the project is based on the article "Bridging machine
learning and cryptography in defence against adversarial
attacks", in the article they used only permutation, we
got quite surprising results with AES in CTR mode
encryption as well.

See figure 4 for the detailed results.

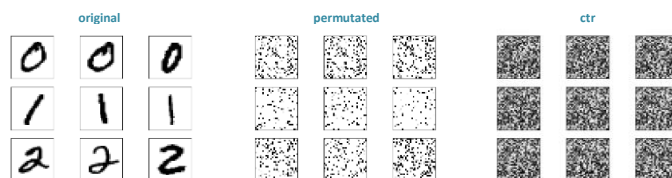


Figure 3: Sample of the encrypted images (permuted and aes-ctr). Interesting to
see how for the human eyes it's impossible to distinguish between various classes
but a DNN model classifies quite well, see figure 4 for accuracies

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 owner of 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

Classification error (%) on the first 1000 test samples							
		original images	mnist		original images	fashion_mnist	
			adversarial images attack	gray-box		adversarial images attack	gray-box
UNENCRYPTED	A	1.49	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
	B	2.10	FGSM	39.50	9.50	FGSM	77.20
PERMUTATED	A	3.70	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
	B	4.20	FGSM	8.60	12.00	FGSM	29.80
AES · CTR	A	3.70	CW I ₂	4.20	17.40	CW I ₂	17.20
	B	2.70	FGSM	4.90	16.70	FGSM	26.50

Figure 4: Table of accuracies