

Detecting Adversarial Examples via Neural Fingerprinting: ICLR Reproducibility Challenge 2019

Jonathan Besomi, Pietro Carta, Pawel Golinski
Department of Computer Science, EPFL, Switzerland

Abstract—Neural fingerprinting [1] has been proposed as a method to detect adversarial examples [3]. The paper has been submitted for acceptance to ICLR’19, and our work revolves around reproducing the results obtained by its authors. The key idea of *Neural fingerprinting* is to detect adversarial examples by using a form of *consistency check of the model output around the input*. After reproducing the method and experiments independently, we confirm *Neural fingerprinting* as viable, robust defense against adversarial attacks.

I. INTRODUCTION

Neural fingerprinting can be considered as a Neural Network regularization method to curb vulnerability to adversarial examples. We provide a brief recap of the working mechanism of *Neural fingerprinting*. For a more extensive explanation of the topic, refer to the original paper under review [1].

We define a fingerprint as a pair $(\Delta x, \{\Delta y_k\})$, with $k \in \{1, \dots, K\}$ and K the number of classes. Δx is the fingerprint direction, and $\{\Delta y_k\}$ is the set of target fingerprint responses. Δx is chosen by sampling $\mathcal{U}[-\epsilon, \epsilon]^k$, and consist in the direction in which each data point x is to be perturbed. $\{\Delta y_k\}$ is a set of the K target fingerprint responses. Every class has a different target response. When evaluating the network f on $x + \Delta x$ and on x , the *change* in output should be close to Δy_k , if x belongs to class k .

Several definitions of the *change* $F(f(x), f(x + \Delta x))$ can be used. The original authors accept as the change the difference between the normalized logits for x and for $x + \Delta x$.

Therefore, during training, when the class of an example x is known, the network parameters are to be trained not only so that $f(x)$ indicates the correct class k , but also so that requirement that the distance between $F(f(x), f(x + \Delta x))$ and Δy_k is minimized.

In this setting, Δy_k is a K dimensional vector, and the original authors suggest that the value of its coordinates be -0.25 if the coordinate index is different than the class index, or 0.75 if they coincide.

The training phase therefore simply requires an addition of the fingerprint loss L_{fp} the standard loss function L_{data} . The minimization problem with respect to the parameters θ of the model becomes:

$$\min_{\theta} L_{data} + \lambda L_{fp} \quad (1)$$

When evaluating the network on a data point x for which the class is unknown, every class is tried for a valid fingerprint match. The validity of a fingerprint response is given by how far the output change F is from its target Δy . If the distance for every class is bigger than some τ , we discard x as an adversarial example. If the distance is smaller than τ for any class, we accept x as real.

The setting can be generalized without efforts to several fingerprint directions. During training, the contribute of each fingerprint is considered equally and the losses for different directions are summed. During evaluation, the distances of the fingerprint responses from the target responses are averaged.

The original authors call the set of fingerprints χ , and the distance of the responses from the target fingerprints $D(x, k, f, F, \chi)$. If the number of fingerprints is N , D is defined as follows:

$$D(x, f, \chi, k) = \frac{1}{N} \sum_{i=1}^N \|F(x, \Delta x_i) - \Delta y_{i,k}\|_2 \quad (2)$$

input : example x , comparison function D ,
fingerprinting χ , model f
output: accept/reject
if $\exists k$ such that $D(x, f, \chi, k) \leq \tau$ **then**
| Return accept;
else
| Return reject;
end

Algorithm 1: Neural Fingerprint

II. IMPLEMENTATION

The authors of [1] include their code for reproducibility purposes. The original code from the paper is mostly implemented in Pytorch, but uses the Tensorflow library Cleverhans [5] for adversarial examples generation. It also requires a number of other dependencies. However, the library versions were not explicitly specified, and we were unable to run it without errors. We therefore resolved to produce an independent reimplementation of *Neural fingerprinting* within an up-to-date version of Tensorflow.

The reimplementation pointed out a couple of minor differences relative to what was written in the original paper. In

section 2.3 of [1], the Δx are said to be chosen by uniform sampling of $[-\epsilon, \epsilon]^l$, where l is the input dimension. However the reference implementation for MNIST fingerprints shared by the paper authors uses a uniform sampling from $[0, \epsilon]^l$.

Another point was specially confusing. *Neural fingerprinting* is said to be robust with respect to the choice of Δy . In practice, we found that targeting with $\delta y_j = -0.25$ if $j = i$ and $\Delta y_j = 0.75$ if $j \neq i$ when the class of the example is i proved computationally unstable, often taking several training epochs before sudden convergence. On the other hand, using $\Delta \hat{y} = -\Delta y$ as target for the fingerprints showed much stabler convergence.

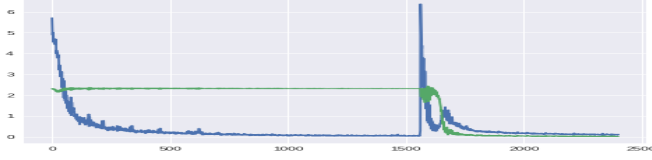


Fig. 1. data (blue) and fingerprint (green) losses when training with proposed Δy .

III. EVALUATION

We were able to reproduce key results from the paper after rewriting the code in Tensorflow. Experiments are available in the repository accompanying this report, in form of a single Jupyter Notebook. The following sections describes our findings from each reproduced experiment.

A. Gaussian balls

The original paper features a section that investigates the effects of *Neural fingerprinting* on a simple classifier to distinguish the points of two two-dimensional Gaussian balls. In the original paper, this example is supposed to provide qualitative insight on the mechanism by which *Neural fingerprinting* works. This example provided us a simple testing framework for our implementation. Since the size of the perturbation Δx was not provided in the original paper, we derived a viable one through trial and error, and came to the understanding that the perturbation size must be in the order of magnitude of the size of the separation between classes.

Our Gaussian balls are centered around $(-3, -3)$ and $(3, 3)$, with identity covariance. When tuning the ϵ parameter responsible for the perturbation size, we noticed that for small epsilons in the order of 0.1, the expected fingerprint responses couldn't be trained, while for larger epsilon in the order of 6 and more the response could be easily trained. This is intuitive when reasoning that small epsilons translate the input x to a nearby $x + \Delta x$. It follows that $f(x)$ and $f(x + \Delta x)$ will tend to be close together, and learning a specific difference might require a huge optimization effort. When Δx is big, on the other hand, the fingerprints can be learned without making the classification function too unstable. We therefore suggest that the ϵ parameter should be in the order of magnitude of the distance of each feature between a class and the others.

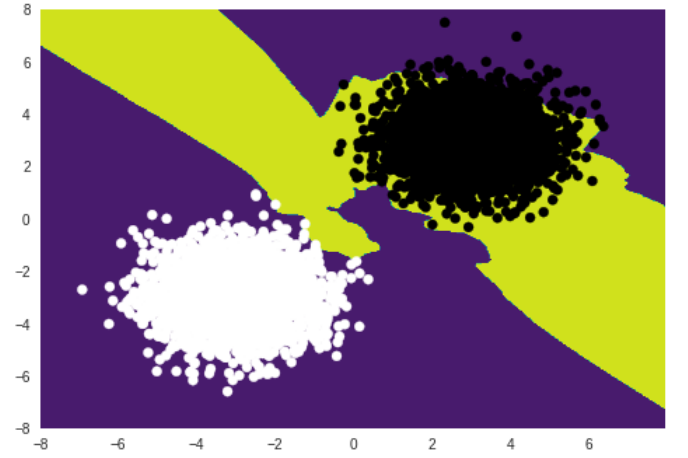


Fig. 2. Including the fingerprint loss in the training generates highly non-linear decision regions

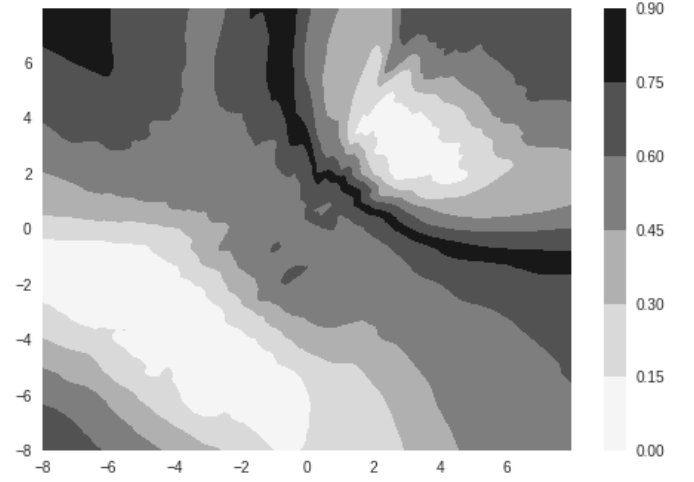


Fig. 3. The fingerprint loss for the gaussian case. Darker regions are likely to be classified as fake.

B. MNIST

As in the original paper, we trained a simple 4-layer network to classify handwritten digits using MNIST as dataset. We implemented the same neural network architecture as specified in table 8 of the original paper.

While training the network, we used five fingerprints, selecting (Δx) perturbations from uniform distribution on $[0, 0.1]$. Δy has been chosen for each fingerprint i and class j had the form of $\Delta y_{n \neq j}^{i,k} = 0.25$, $\Delta y_{j=n}^{i,k} = -0.75$. Our network was trained on weighted sum of fingerprints and cross entropy loss, were both losses contributed equally to the sum. We used the Adam optimizer, with a learning rate of 0.001 in order to achieve convergence. After 10 epochs, the network yields 99.2% accuracy on test data.

Once our model was working properly we started attacking it. Our first experiment has been to compare the distribution of fingerprints loss between original and adversarial examples. Table 2 and Figure 4 from the original paper. Due to

computational constraints, we were able to reproduce only FGSM[3] and BIM[4] attack results, although original paper reports that Neural Fingerprints behave similarly even in case of more sophisticated JSMA [6] and CW[2] attack methods. The results of this experiment are described on Figures 4 and 5.

Next, we examined how different settings of threshold parameters correspond to true and false positive rate. The standard metric to examine that is ROC curve. The result is described on Figure 6 and corresponds to Table 4 and Figure 7c from the original paper.

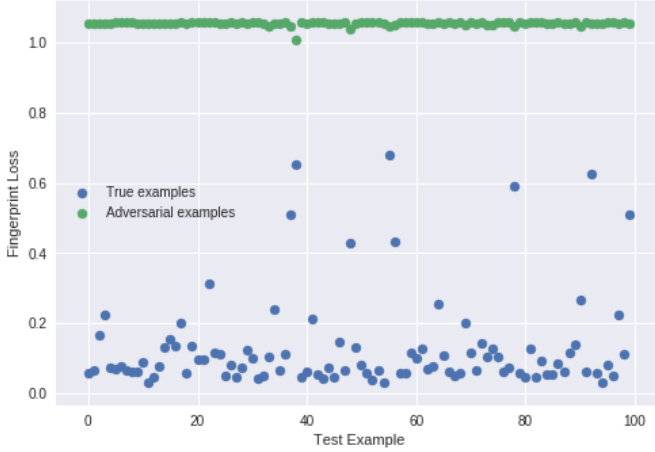


Fig. 4. In case of FGSM attack, there is a clear distinction between true test examples (blue) and adversarial examples (green), when comparing fingerprint loss function value. The accuracy of network on adversarial examples generated by FGSM is 5%.

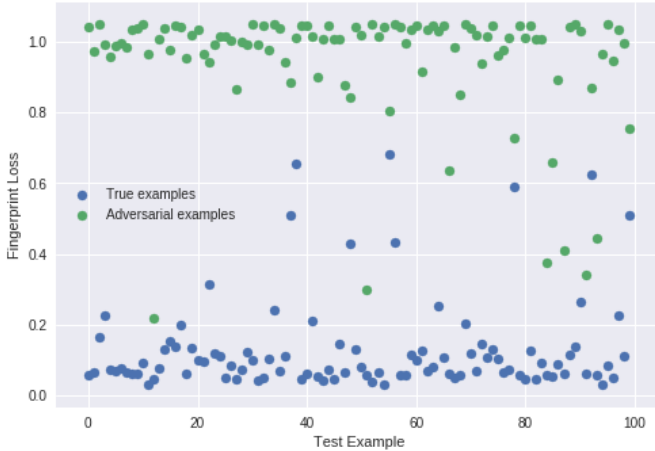


Fig. 5. In case of BIM attack, the distinction is blurrier and there are some outliers. The accuracy of network on adversarial examples generated by BIM is 2%.

C. CIFAR-10

For the CIFAR-10 problem, we obtained similar results as to MNIST. We used the architecture described in table 9 of the paper. More details are available in Jupyter notebook accompanying this report.

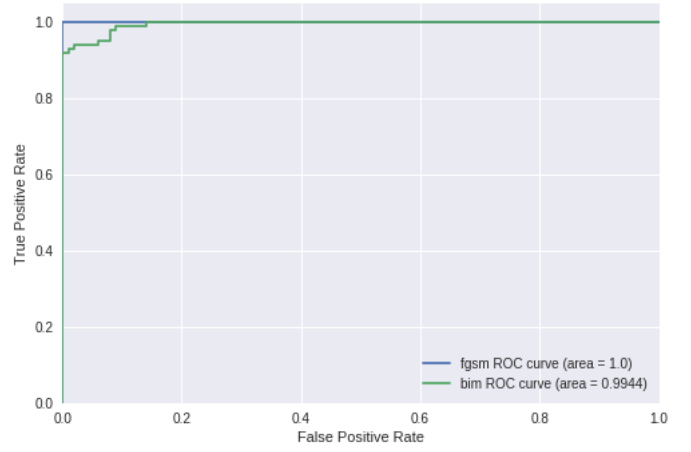


Fig. 6. High AUC-ROC indicates great robustness of Neural Fingerprints against both FGSM and BIM attacks.

IV. VARYING NUMBER OF FINGERPRINTS

We examined how varying number of fingerprints affects final network accuracy, execution time and robustness. All figures are available in Jupyter notebook in its final section, and correspond to figure 7b from original paper. Our conclusion from these experiments is similar to the ones that were presented in the paper - increasing number of fingerprints doesn't affect accuracy, but increases robustness. Moreover, we examined how adding fingerprints affects execution time. As expected, increasing fingerprint number scales training and evaluation time linearly - every additional fingerprint requires an additional network inference.

V. CLEVERHANS ATTACK

In order to test the efficiency and the statement of the paper, with the support of Tensorflow library for adversarial attack Cleverhans we prepared some attacks. Figure 7 show three example of attack generated by Cleverhans.

VI. MACHINE LEARNING STACK

During the whole project we tried different platforms and solution to work with the data. Because the code is available on Github, we attempted to execute it online on Google Cloud Platform. Due to misstated dependencies and the use of obsolete library versions, we opted for a reimplementation, We started experimenting locally within Jupyter notebooks. The problem with such approach is that the execution is slow, since our machines don't feature GPUs. Moreover it was complicated to share resources and code among the team member. A good compromise turned out to be Google Colab Notebook. The advantages of this new services is that it offers for free (as of December 2018) collaborative notebooks accelerated by Tesla K80 GPUs.

VII. DIFFICULTIES ENCOUNTERED

The original paper effectively explains the concept of *Neural fingerprinting* and provides all the tools for reimplementation. We still incurred some issues while implementing

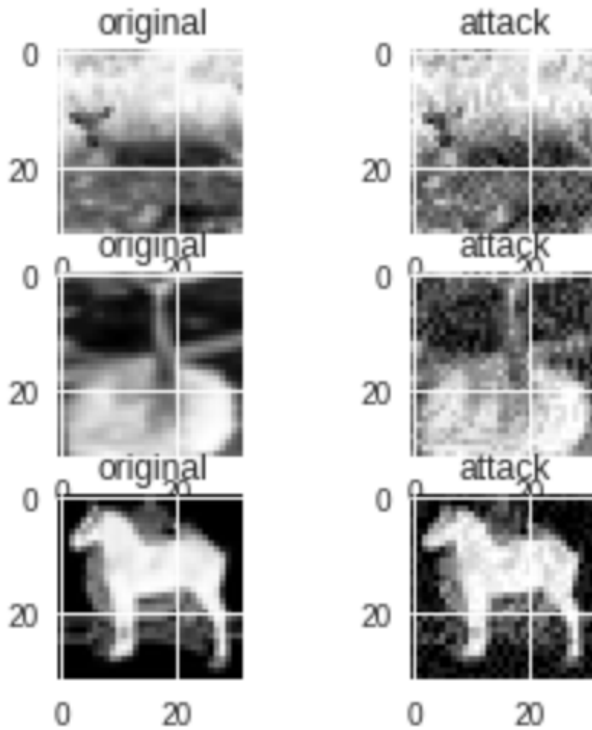


Fig. 7. Cleverhans attack on CIFAR-10

the Gaussian example, since the hyperparameters were not provided. We suspect that due to the low dimensionality of the problem, the optimization was specially sensible to the choice of parameters. The biggest issue remained the trouble when running the provided code. Libraries in use are out of date and since the code is distributed without any form of containerization, sourcing every piece of the puzzle can be very time consuming. We opted to reimplement it ourselves both for the gained flexibility and to learn more specific details about the method. However, our schedule was impacted by this choice, and we couldn't reach our most ambitious goals (replication of adaptive whitebox attacks, implementation of attacks specific for *Neural fingerprinting*). Reproducibility in Machine Learning research remains a moving target, due to the fast moving ecosystem and the necessity for flexibility which is fundamental for researchers. A satisfying middle point could however be achieved by use of up to date libraries and more cautious development practices

VIII. CONCLUSION

Neural fingerprinting as described in the original paper effectively detects all the classes of adversarial attacks that we have tested. Unfortunately, we were far from exhaustive, and in computer security absence of evidence of vulnerabilities doesn't constitute evidence of their absence. We are however confident that this method deserves the attention of the research community, and we impatiently wait to see it faced to more and better attacks.

REFERENCES

- [1] Anonymous. Detecting adversarial examples via neural fingerprinting. In *Submitted to International Conference on Learning Representations*, 2019. under review.
- [2] A. Athalye, N. Carlini, and D. Wagner. Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples. *arXiv e-prints*, page arXiv:1802.00420, Feb. 2018.
- [3] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and Harnessing Adversarial Examples. *arXiv e-prints*, page arXiv:1412.6572, Dec. 2014.
- [4] A. Kurakin, I. Goodfellow, and S. Bengio. Adversarial examples in the physical world. *arXiv e-prints*, page arXiv:1607.02533, July 2016.
- [5] N. Papernot, F. Faghri, N. Carlini, I. Goodfellow, R. Feinman, A. Kurakin, C. Xie, Y. Sharma, T. Brown, A. Roy, A. Matyasko, V. Behzadan, K. Hambarzumyan, Z. Zhang, Y.-L. Juang, Z. Li, R. Sheatsley, A. Garg, J. Uesato, W. Gierke, Y. Dong, D. Berthelot, P. Hendricks, J. Rauber, R. Long, and P. McDaniel. Technical Report on the CleverHans v2.1.0 Adversarial Examples Library. *arXiv e-prints*, page arXiv:1610.00768, Oct. 2016.
- [6] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. Berkay Celik, and A. Swami. The Limitations of Deep Learning in Adversarial Settings. *arXiv e-prints*, page arXiv:1511.07528, Nov. 2015.