

Towards Information Theoretic Adversarial Examples

Chia-Yi Hsu*, Pin-Yu Chen[†] and Chia-Mu Yu*

*National Chung Hsing University, Taiwan

[†]IBM Research

Abstract—Deep learning has shown impressive performance on wide applications. However, recent research shows that deep neural networks are vulnerable to well-crafted input samples, called adversarial examples. Adversarial examples are indistinguishable to humans but can easily fool deep neural networks. Nowadays, most of attacks measure human’s perception of the image quality with L_2 -norm or L_∞ -norm perturbation constraints. In this paper, we introduce mutual information (MI) to evaluate image quality of adversarial examples instead of L_p -norm measures. With MI as an information theoretic metric, our quantitative and qualitative results show that the resulting adversarial examples are more similar to unperturbed data samples.

I. INTRODUCTION

In recent years, deep learning has made significant progress in many tasks of machine learning such as image classification, language translation, and object detection. However, recent studies demonstrated that well-trained deep neural networks (DNNs) are vulnerable to adversarial examples [1]. The purpose of the attacker is to find adversarial examples and minimize perturbations at the same time.

There have been many efforts to attack DNNs. Carlini and Wanger [2] proposed an optimization-based framework for targeted and untargeted attacks, abbreviated C&W attack. They design a L_2 norm regularized loss function apart from the model prediction loss defined by the logit layer representations in DNNs.

There have been many attacks using L_p -norm to restrict perturbations so that adversarial examples are imperceptible for humans. In this paper, our main contribution is that we are the first replacing statistical distance by information-theoretic metrics. MI is most used in generative adversarial networks (GANs) which measures the similarity between random variables and images.

II. OUR METHOD

Belghazi et al. [3] proposed a method using neural network to estimate mutual information called MINE. The concept is to select \mathcal{F} to be the family of functions $T_\Theta : \mathcal{X} \times \mathcal{Z} \rightarrow \mathbb{R}$ parametrized by a deep neural network with parameters $\theta \in \Theta$. We formalize the bound:

$$I(X, Z) \geq I_\Theta(X, Z),$$

where $I_\Theta(X, Z)$ is the neural information quantity defined as

$$I_\Theta(X, Z) = \sup_{\theta \in \Theta} \mathbb{E}_{\mathbb{P}_{XZ}}[T_\theta] - \log(\mathbb{E}_{\mathbb{P}_X \otimes \mathbb{P}_Z}[e^{T_\theta}]).$$

In our case, we need to compute MI of a single pair images. However, MINE must use batch version so that we make an image to batch version by gaussian random projection.

We use the C&W attack loss without L_2 -norm regularized loss in addition to negative MI maximized by gradient ascent. We formalize our attack as the following optimization problem:

$$\underset{\delta, \Theta}{\text{minimize}} \quad c \cdot f(x + \delta) - \alpha \cdot I_\Theta(x, x + \delta)$$

such that $x + \delta \in [0, 1]^n$ and $\delta \in [-\epsilon, \epsilon]^n$.

with f defined as

$$f(x') = \max \{Z(x')_{l_x} - \max(Z(x')_i : i \neq l_x), -\kappa\}.$$

The loss f is the best objection function found earlier, modified slightly so that we can control the confidence with which the misclassification occurs by adjusting κ .

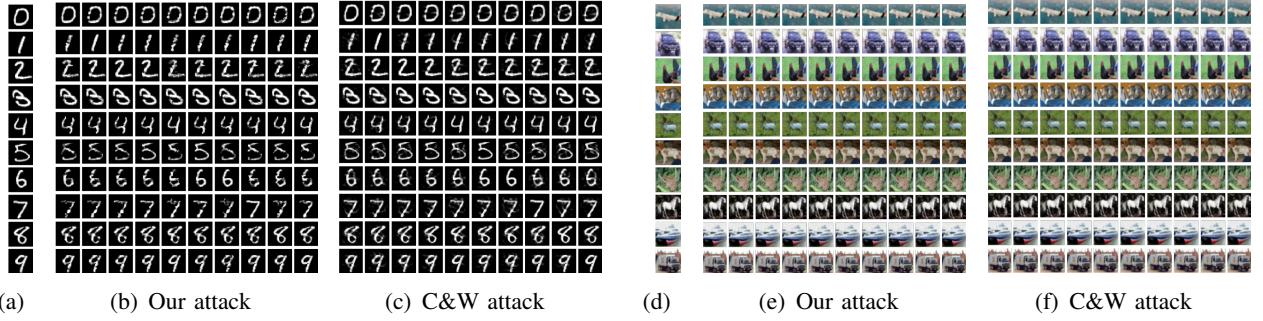


Fig. 1. Visual comparison of successful adversarial examples in MNIST and CIFAR-10. Each row shows crafted adversarial examples from the sampled images from (a) and (d). Each column in (b), (c), (e), (f) indexes the targeted class for attack.

III. EXPERIMENTS

A. Experiments Setup and Parameters Setting

We compare our attack with C&W attack. We would like to show that our attack can achieve high image quality as C&W attack. We implement untargeted attacks and targeted attacks. In MNIST, we set $\alpha = 0.5$ with different epsilons. For CIFAR-10 dataset, we set $\epsilon = -8/256$ and $\alpha = 0.5$. For C&W attack, we set κ is equal to zero in both datasets.

B. Targeted attack

We showed the adversarial examples generated by our attack and C&W attack in Fig1 and carried out targeted attack. our attack can still maintain high image quality. We also used FID scores to assess the image quality showed in Table I.

TABLE I

FID scores comparison of adversarial examples and training data.

	MNIST	CIFAR-10
Our attack	124.854	246.86
C&W attack	111.745	262.177

C. Untargeted attack

We showed the results of untargeted attack in Fig.2 and Fig.3. Fig.2 showed that most perturbations of our attack are on the main part of digits instead of backgrounds when we set $\epsilon = 1$ in MNIST.

IV. FUTURE WORK AND CONCLUSION

In this paper, we propose to compute MI between input samples and adversarial examples as an information theoretic similarity measure. The

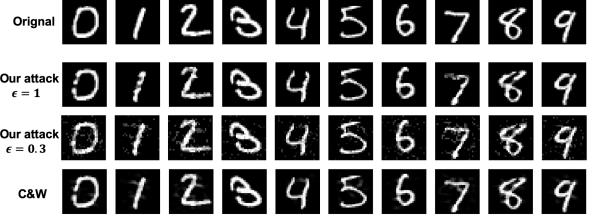


Fig. 2. Visual comparison of successful adversarial examples in MNIST under different untargeted attacks with different epsilons.



Fig. 3. Visual comparison of successful adversarial examples in CIFAR-10 under different untargeted attacks.

results demonstrated that MI can help make adversarial examples look more similar to normal images. Our framework can be extended to computing MI between the middle layer output of the classifier with normal images and adversarial examples as inputs. Our future work will explore the effectiveness of MI-based adversarial examples against known defense methods.

REFERENCES

- [1] D. Su, H. Zhang, H. Chen, J. Yi, P.-Y. Chen, and Y. Gao, “Is robustness the cost of accuracy?—a comprehensive study on the robustness of 18 deep image classification models,” *ECCV, arXiv preprint arXiv:1808.01688*, 2018.
- [2] N. Carlini and D. Wagner, “Towards evaluating the robustness of neural networks,” in *IEEE Symposium on Security and Privacy (SP)*, 2017, pp. 39–57.
- [3] M. I. Belghazi, A. Baratin, S. Rajeswar, S. Ozair, Y. Bengio, A. Courville, and R. D. Hjelm, “Mine: mutual information neural estimation,” *arXiv preprint arXiv:1801.04062*, 2018.