



Learning to Confuse: Generating Training Time Adversarial Data with Auto-Encoder

Ji Feng*, Qi-Zhi Cai*, Zhi-Hua Zhou

NeurIPS 2019

Presented by Jiaru Zhang, AISIG

March 8, 2022



SHANGHAI JIAO TONG
UNIVERSITY

① Introduction

② The proposed method

③ Experiments

④ Conclusion



Section 1

Introduction



Introduction

- Problem: Adding imperceptible noises to the training data to confuse classifier in testing.



Clean training examples



Acc: High



Adversarial training samples



Acc: Low

Training

Testing



Section 2

The proposed method

Problem formulation

The learning target of a neural network f_θ with parameter θ is

Target

$$\theta^* = \arg \min_{\theta} \sum_{(x,y) \sim \mathcal{D}} [\mathcal{L}(f_\theta(x), y)] \quad (1)$$

Noise generator: g_ξ

Constraint on noise

$$\forall x, \|g_\xi(x)\|_\infty \leq \epsilon \quad (2)$$

In this work, an encoder-decoder network with activation $\epsilon \cdot (\tanh(\cdot))$ in the last layer is used.



Optimization

- The equality constraint can be relaxed into

$$\theta_i = \theta_{i-1} - \alpha \cdot \nabla_{\theta_{i-1}} \mathcal{L}(f_{\theta_{i-1}}(x + g_{\xi}(x)), y) \quad (4)$$

- The basic idea is to alternatively update f_{θ} on **noisy data** via gradient descent, and g_{ξ} on **clean data** over gradient ascent.
- However, f_{θ} and g_{ξ} won't converge in practice.



Optimization

- Collecting the update trajectories for f_θ
- Update g_ξ based on such trajectories.

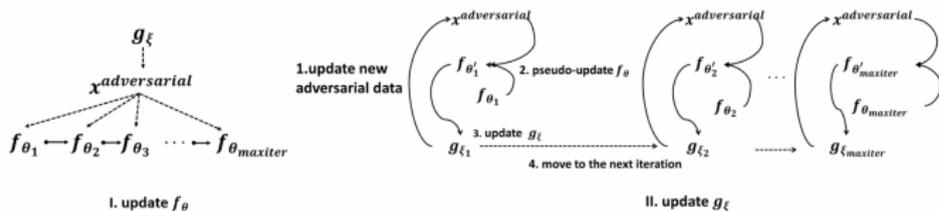


Figure 1: An overview for learning to confuse: Decoupling the alternating update for f_θ and g_ξ

- Implementation trick: save g_ξ instead of f_{θ_i} .

Label specific adversaries

- It can be easily transfer to the label specific conditions.

Label specific adversaries

Replace

$$\max_{\xi} \sum_{(x,y) \sim \mathcal{D}} [\mathcal{L}(f_{\theta^*}(\xi)(x), y)] \quad (5)$$

into

$$\min_{\xi} \sum_{(x,y) \sim \mathcal{D}} [\mathcal{L}(f_{\theta^*}(\xi)(x), \eta(y))], \quad (6)$$

where η is a predefined label transformation function.

Section 3

Experiments

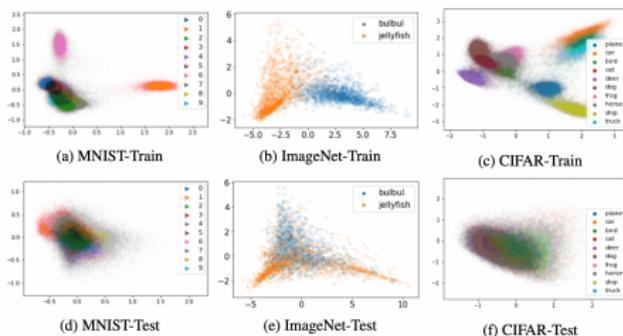


Performance Evaluation

- The test accuracy obviously dropped when trained on the adversarial datasets.

| | MNIST | ImageNet | CIFAR-10 |
|------------------|------------------|------------------|------------------|
| Clean Data | 99.32 ± 0.05 | 88.5 ± 2.32 | 77.28 ± 0.17 |
| Adversarial Data | 0.25 ± 0.04 | 54.2 ± 11.19 | 28.77 ± 2.80 |

- The classifier trained on the adversarial data cannot differentiate the clean samples.



Effect of varying parameters

- There is a sudden drop in performance when the perturbation constraint ϵ exceeds 0.15.
- The proposed method performs better than *random flip*.

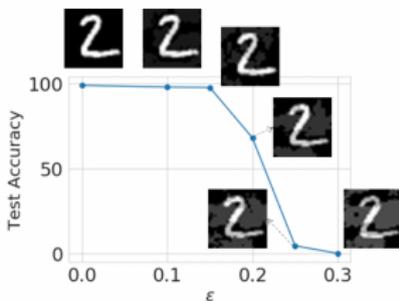


Figure 4: Effect of varying ϵ .

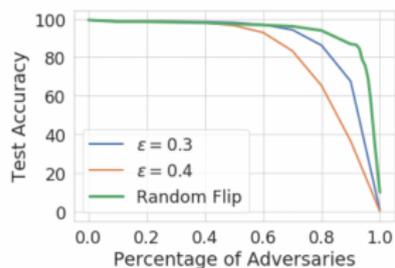


Figure 5: Varying the ratio of adversaries under different ϵ .

Generalization Gap

- A clear generalization gap is observed during the training process.
- It is conjectured that the deep model tends to overfit towards the adversarial noises.

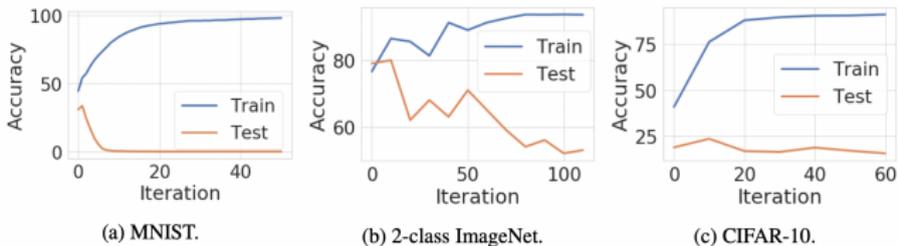


Figure 8: Learning curves for f_θ

Weight Visualizations

- The victim SVM weights went to the opposite direction and tend to overfits on image corners.

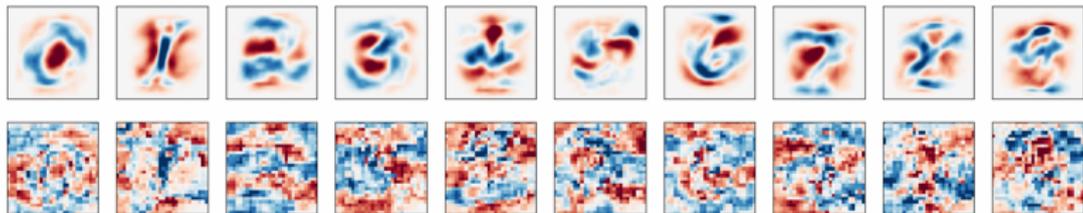


Figure 10: LinearSVM weights visualization for MNIST. Top row: Weights trained on clean training data. Bottom row: Weights trained on adversarial training data.

Section 4

Conclusion

Conclusion

- This paper proposed a general framework for generating training time adversarial data.
- A simple yet effective training scheme to train both networks.
- Experiments on image data confirm the effectiveness.

Related consecutive work

- A concurrent work minimizes the gradients of weights to make models harder to converge in transfer learning ¹.
- “Inversely adversarial noise” generated by PGD has a similar effect and is used to synthesize *Unlearnable Examples* ².
- Gradient manipulation is used to generate poisoned dataset ³.
- Adversarial examples make stronger poisons ⁴.
- Adversarial training serves as a defense with theoretical guarantee ⁵.

¹ Juncheng Shen, Xiaolei Zhu, De Ma. TensorClog: An Imperceptible Poisoning Attack on Deep Neural Network Applications, in IEEE Access, vol. 7, pp. 41498-41506, 2019

² Hanxun Huang, Xingjun Ma, Sarah Monazam Erfani, James Bailey, Yisen Wang. Unlearnable Examples: Making Personal Data Unexploitable. In ICLR, 2021.

³ Liam H Fowl, Ping-yeh Chiang, Micah Goldblum, Jonas Geiping, Arpit Amit Bansal, Wojciech Czaja, Tom Goldstein. Protecting Proprietary Data: Poisoning for Secure Dataset Release. In arxiv preprint, 2103.02683.

⁴ Liam H Fowl, Micah Goldblum, Ping-yeh Chiang, Jonas Geiping, Wojciech Czaja, Tom Goldstein. Adversarial Examples Make Strong Poisons. In NeurIPS, 2021.

⁵ Lue Tao, Lei Feng, Jinfeng Yi, Sheng-Jun Huang, Songcan Chen. Better Safe Than Sorry: Preventing Delusive Adversaries with Adversarial Training. In NeurIPS, 2021.

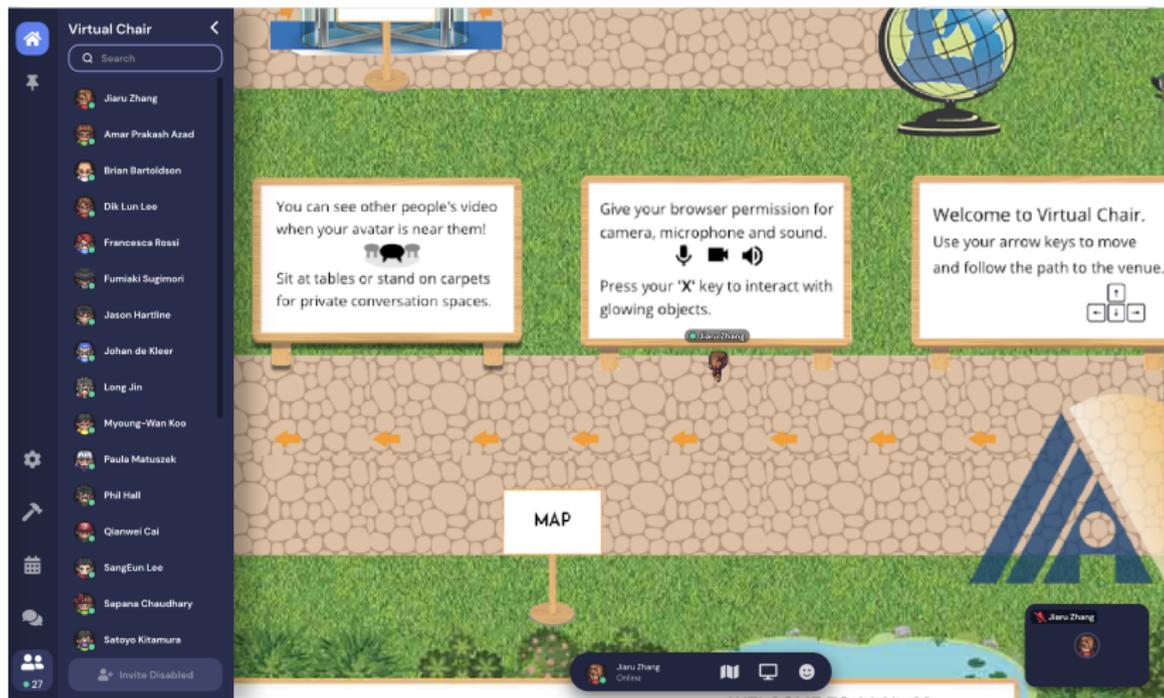


Section 5

Others



Experience in AAAI



Experience in AAI

 **AAAI-2022**

[Home](#) [Venue](#) [Schedule](#) [Papers](#) [Plenary](#) [Events](#) ▾

Select ▾ Select ▾ Show Favorites



Each dot represents a paper. They are arranged by a measure of similarity.

If you hover over a dot, you see the related paper.

If you click on a dot, you go to the related paper page.

You can search for papers by author, keyword, or title

Drag a rectangle to summarize an area of the plot.

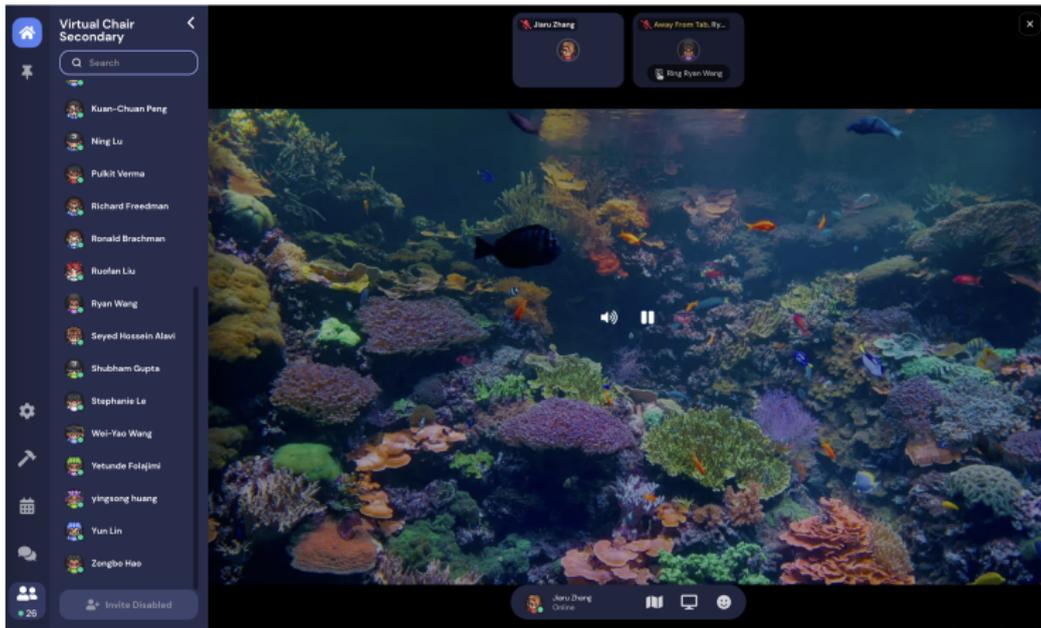
IAO TONG

Experience in AAAI



SHANGHAI JIAO TONG UNIVERSITY

Experience in AAAI



Experience in AAAI



