



Introduction

Many adversarial attacks have recently been proposed, which can all successfully find a small perturbation on the input image to fool the DNN based classifier. E.g., the swan can be mis-predicted after adding perturbations from many different attacks listed here.



Inspired by anomaly detection in control flow programs, we propose the hypothesis that only a small set of synapses and neurons are critical to the predicted class (effective path), and adversarial attacks use abnormal path to modify the predicted class as shown, which is verified by evaluation and leads to a novel method for adversarial defense.

Effective Path Extraction

> We sort the activation values in the receptive field of the predicted class by their contributions, and add the first several of them to

 $\min_{\tilde{K}_p^L} |\tilde{K}_p^L|, s.t. \sum_{\tilde{k}_p}$ $\mathcal{W}^L = \{ w_{k,p}^L | k \in \tilde{K}_p^L \}$ $\mathcal{N}^{L-1} = \{ n_k^{L-1} | k \in \mathcal{N} \}$

effective path according to a threshold. This extraction process starts at the last layer and moves backward to the first layer.



Adversarial Defense Through Network Profiling Based Path Extraction Yuxian Qiu¹, Jingwen Leng¹, Cong Guo¹, Quan Chen¹, Chao Li¹, Minyi Guo¹, Yuhao Zhu² ¹Shanghai Jiao Tong University, ²University of Rochester



$$w_k^{L-1} \times w_{k,p}^L \ge \theta \times n_p^L$$

$$\{\tilde{K}_p^L\}$$

 $\times -6$

Effective Path vs Network Pruning

- Both effective path and network pruning extract a sparse subset of a network. However, effective path has some properties to tell it apart from network pruning and enable adversarial detection:
- > Path specialization: Different classes activate not only sparse but also a distinctive set of paths. E.g., different digits in MNIST have low effective path similarity with each other as shown on the right.
- \succ When comparing with the predicted class's effective path, normal examples' effective paths (blue line) achieve higher similarity than adversarial examples' (other lines), which confirms that adversarial attacks use abnormal path. We show each layer's similarity for ResNet-50 below.





Adversarial Example Detection

- detection achieves:
- Higher accuracy for different attacks and defense models (linear/AdaBoostg /gradient boosting/random forest)
- Less training samples. AUC can reach 0.9 for all attacks with less than 1000 training samples.

- > Better generalizability. our work generalizes well to unseen attacks because effective path captures their common behavior.
- > Faster detection. Extracting the partial effective path for only necessary layers will lead to more than 200x speedup to CDRP.
- > Wider application. Effective path can also be used to detect wrongly predicted examples (AUC=0.86) and random noise (AUC=0.91).



LONG BEACH CALIFORNIA June 16-20, 2019

> When comparing with prior work CDRP, effective path based



