



An Adversarial Robustness Perspective on the Topology of Neural Networks



Authors :

Morgane Gibert (Giteo / Telecom Paris)

Thomas Ricatte (Amazon)

Elvis Dohmatob (Facebook AI Research)

Published :

Safety ML Workshop

NeurIPS 2022

The paper is available on [ArXiv](#)

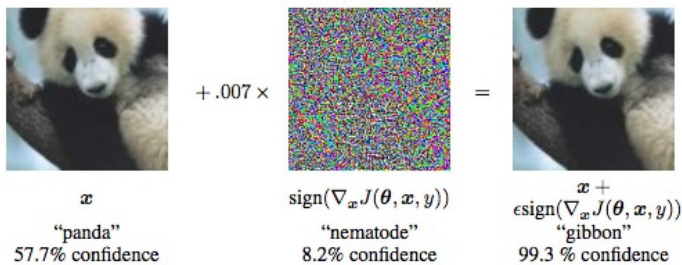
Content

- ① Introduction to the Adversarial Robustness problem in Deep Learning
- ② Our Hypothesis
 - a) Some observations on attacks
 - b) The hypothesis
- ③ Method : Extracting Topological Features
 - a) Induced graph
 - b) Selecting Under-optimized edges
 - c) Persistent Diagrams
- ④ Experiments : Differentiating Clean vs Adversarial ex.
 - a) Quantitative Results (simple setting)
 - b) Detecting Adversarial Examples
 - c) Additional Results

① Intro to Adversarial Robustness

Phenomenon first described in [Intriguing Properties of Neural Networks, Szegedy et al. 2013]

Illustration from [Explaining and Harnessing Adversarial Examples, Goodfellow et al, 2014]

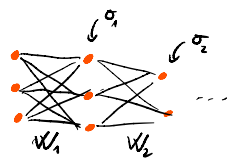


Classification problem setting

* Input $x \in \mathbb{R}^d$
Label / Class $y \in \{1, \dots, K\}$ } P is the joint distribution on (x, y)

* Dataset (x_i, y_i)

* Neural Network:



↳ Feature map $g : x \in \mathbb{R}^d \mapsto g(x) \in \mathbb{R}^k$

$$g(x) = \sigma_L (W_L \sigma_{L-1} (W_{L-1} \sigma_{L-2} \dots (\sigma_1 (W_1 x)))$$

where W_p is the weight matrix for layer p

σ_p is the activation function for layer p

Notation: $g(x)_p$ = output / activation of layer p

↳ NN: $h(x) = \underset{k=1, \dots, K}{\text{argmax}} g(x)$

Adversarial example

* an adv example x^{adv} is a perturbed version of a clean ex x :

$$x^{adv} = x + \delta, \quad \delta = \text{perturbation}$$

* Small perturbation: $\|\delta\|_{\infty} \leq \epsilon$

* Goal = misclassification $\rightarrow h(x^{adv}) \neq h(x)$

Attacks

The underlying optim pb to find the optimal perturbation δ is hard

\Rightarrow Several algos using different strategies

- * FGSM, PGD
 - * CW
 - * Boundary
- } white-box
- \rightarrow black-box

(Known) Characteristics of attacks

- ① They work really well
- ② They transfer between architectures
- ③ They can be $\left. \begin{array}{l} \text{on-manifold} \rightarrow \text{leverage useful \& non-robust features in the data} \\ \text{off-manifold} \rightarrow \text{non-useful} \end{array} \right\}$
- ④ They are less efficient against highly regularized NNs / smooth decision boundary

② Our Hypothesis

a) Some observations

* Pruning

→ NNs are over-parametrized

→ pruning do not hinder accuracy [The lottery ticket hypothesis: finding sparse, trainable NNs, Franckle & Carbin, 2018]

⇒ Only a subset of parameters are really important for inference (on clean inputs)

⇒ Clean inputs use highway edges

* Adversarial examples

→ They can use under-optimized parameters/edges (off-manifold)
i.e. more diffuse paths

b) The hypothesis

Clean vs adv examples induce differences in the topological structure of their respective induced graph because adv (and not clean inputs) leverage under-optimized edges

③ Method: Extracting Topological Features

a) Induced Graph

* Idea: represent the way an input x traverses a NN g

* Induced graph: $G(x, g) = (V, E)$

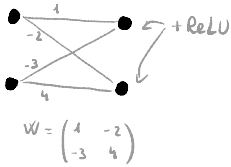
$$V = \{1, 2, \dots, n_0 + \dots + n_L\} \quad \text{where } n_f = \text{nb neuron in layer } f$$
$$E = \{(v^f, v^{f+1}, w_{v^f, v^{f+1}})\} \quad \text{where } w_{v^f, v^{f+1}} \text{ is the weight}$$

→ linear layer: $w_{v^f, v^{f+1}} = |g_f(x)_{v^f} \times (w_f)_{v^f, v^{f+1}}|$

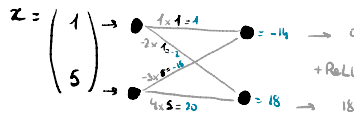
→ generalise to all architecture types (eg ResNet) and layers (eg convolutional)

* Illustration:

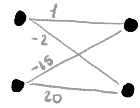
A trained NN



Computing the passage of x



Induced graph $G(x, g)$:



b) Selecting Under-optimized Edges

* Why?

→ according to the hypothesis, they are sufficient

→ complexity: NNs have way too many neurons/edges

* How?

→ definition from [Deconstructing Lottery Tickets: Zeros, signs and the supermask, Zhou et al., 2019]

- An edge (u,v) is **under-optimized** if its weight has not changed much during training → $|(W_p)_{u,v} - (W_p^{init})_{u,v}| < \text{quantile}(q)$
- So we keep only a fraction q of edges

c) Persistent Diagrams

* Graphs considered: thresholded induced graph $G^q(x, g)$

* Weights: $\tilde{w}_{u,v} = -|w_{u,v}|$ so that highest-weighted edges appear first

* Filtration: natural one based on weight value

→ A time $t \in \mathbb{R}_-$, the considered graph is $G_t^q(x, g)$

where $G_t^q(x, g) = (V_t, E_t)$ with $E_t = \underbrace{\{(u,v, \tilde{w}_{u,v}) \mid \tilde{w}_{u,v} \leq t\}}_{\in E}$

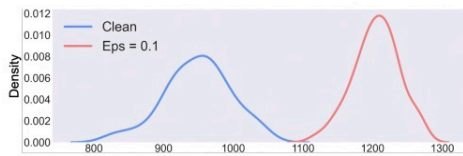
④ Experiments : Differentiating Clean vs Adv inputs

a) Quantitative Results

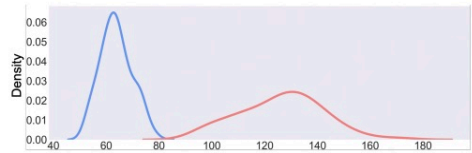
* Summary : I have a collection of persistence diagrams for

- 1) adv inputs
- 2) clean inputs

* Simple idea : count the nb of points in the PDs → all points
→ infinitely-lived ones



(a) Distribution of all PD points.

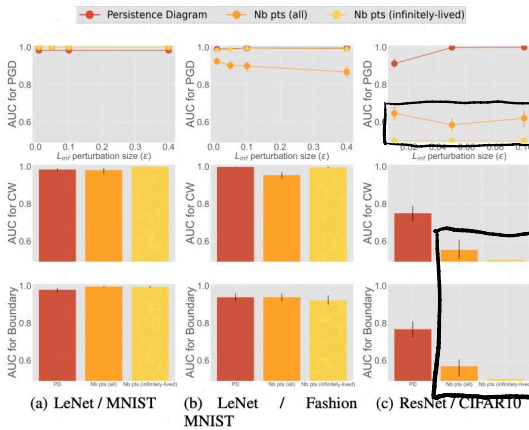


(b) Distribution of infinitely-lived PD points.

Figure 2: PD points computed on MNIST / LeNet

b) Detecting adv examples

* Nb of points is a limited strategy



In this more complicated setting, the results are bad

Figure 19: Unsupervised detection results using number of points only.

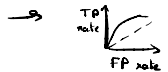
* Using all the information provided by PDs => detection

→ One-class SVM is * unsupervised (the algo only sees clean PDs during training)

* relevant kernel on PDs → Sliced-Wasserstein Kernel [Sliced-Wasserstein Kernel for Persistence Diagram; Carrara, Cuturi & Oudot, 2017]

→ Performance metric: AUC

(Area Under the ROC Curve)



→ Comparison with SOTA:

1) L1D

2) Mahalanobis

3) Baseline from the induced graph = RG

→ features are a vector, elements are the weights of $G^l(x, g)$

* Main results

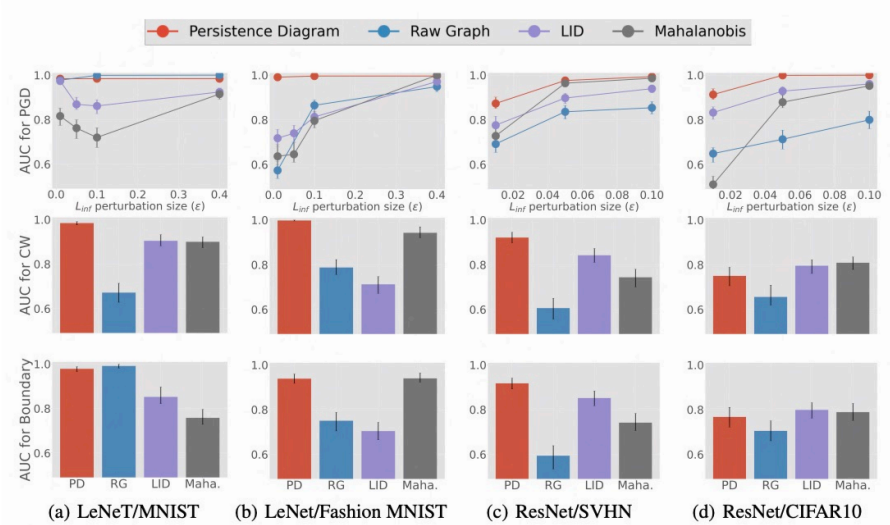


Figure 3: Showing detection AUC for different detection methods (legend) against different kinds of adversarial attacks (rows) and model architectures and datasets (columns). We see that our proposed method based on PD outperforms the SOTA methods, except for one tie.

c) Additional Results

* Still performant on AT networks

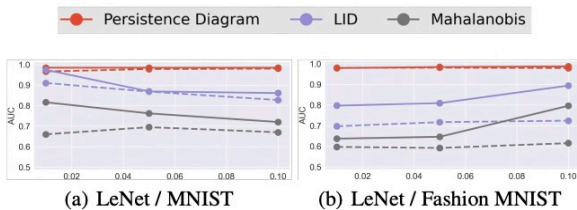


Figure 21: Unsupervised detection results (on PGD) of AT vs standard NNs

* Selecting Under-optimized edges is the good strategy

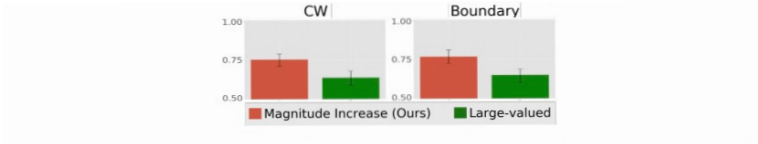


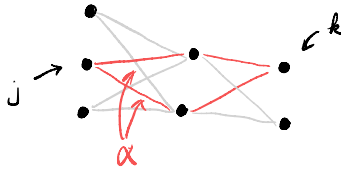
Figure 22: Impact of edge-selection methods on AUC (ResNet / CIFAR10).

* Pruning can improve robustness

Proposition : \forall class $k \in [1, K]$, \forall input feature index $j \in [1, n_0]$

$$\frac{\partial [g(x)]_k}{\partial x_j} = \sum_{\alpha \in \mathcal{A}} W(\alpha)$$

where \mathcal{A} is the set of active paths from j to k
 $W(\alpha) = \prod_{l=1}^L (W_l)_{v^{l-1}(\alpha), v^l(\alpha)}$ is the product of weights in path α



$\Rightarrow \frac{\partial [g(x)]_k}{\partial x_j}$ is a proxy for the vulnerability of class k wrt input feature j .

The norm of the Jacobian matrix $J(x, g) = \left(\frac{\partial [g(x)]_k}{\partial x_j} \right)$ is a general proxy for the robustness to perturbations on x \rightarrow Reducing the cardinality of \mathcal{A} should help

The End!

Thank you &

Question time!