# Adversarial Attacks and Defenses in Images, Graphs and Text: A Review

- 2020, 108 citations

## Vulnerable networks

- CNN
- FC DNN
- RNN
- GCN (graph convolutional networks) - used in fraud detection
    - only necessary to change couple of edges

## Counter-measures

- Gradient masking
- Robust optimization
- Adversary detection

Deep neural-nets reason differntly -> understanding adversarial attacks should help understand this difference

# Definitions and notations

## Threat model

### Adversary's goal

- Poisoning attack - change the behavior of DNN by modifying/inserting few train examples
    - public honeypot - collection of training data for malware detectors
- evasion attack - craft fake examples classifier cannot recognize
    - targeted
    - untargeted

### Adversary's knowledge

- White-box attack - widely studied, easily analyzed mathematically
- Black-box attack - practical
- Semi-white (gray) box attack - train generative model in white-box setting, then use in black-box scenario (TREMBA)

### Victim models

- conventional machine learning models - SVM, Naive-Bayes
- DNN - not well understood how they work, studying security necessary
    - FC
    - CNN - sparse version of FC
    - GCN
    - RNN

## Security evaluation

### Robustness

- **Minimal pertubation** - The smallest pertubation to fool the network
    - $\delta_{min} = \arg\min_{\delta} \|\delta\|$ s.t. $F(x + \delta) \neq y$
- **Robustness** - The norm of minimal pertubation on particular example
    - $r(x, F) = \|\delta_{min}\|$
- **Global robustness** - Expectation of robustness over the whole dataset

- $\rho(F) = \underset{x \sim \mathcal{D}}{\mathbb{E}} r(x, F)$

**Adversarial risk (loss)**

- **Most-adversarial example** - Given classifier $F$, datapoint $x$ and $\epsilon$ ball, $x_{adv}$ is the adversarial example with the largest loss
    - $x_{adv}$ is the point, where the classifier is the most likely to be fooled
    - $x_{adv} = \underset{x'}{\arg\max} \mathcal{L}(x\ F)$ s.t. $\|x' - x\| \le \epsilon$
- **Adversarial loss** - Loss value of the most-adversarial example
    - $\mathcal{L}_{adv}(x) = \mathcal{L}(x_{adv}) = \underset{\|x'-x\|<\epsilon}{\max} \mathcal{L}(\theta, x', y)$
- **Global adversarial loss (adversarial risk)** - The expectation of adversarial loss over the data distribution $\mathcal{D}$
    - $\mathcal{R}_{adv}(F) = \underset{x \sim \mathcal{D}}{\mathbb{E}} \mathcal{L}_{adv}(x) = \underset{x \sim \mathcal{D}}{\mathbb{E}} \underset{\|x'-x\|<\epsilon}{\max} \mathcal{L}(\theta, x', y)$

**Adversarial risk vs. risk**

- The concept of *Adversarial risk* is similar to the definition of classifier risk (empirical risk)
    - $\mathcal{R}(F) = \underset{x \sim \mathcal{D}}{\mathbb{E}} \mathcal{L}(\theta, x, y)$
    - Global adversarial risk (loss) is in a sense empirical risk but on the most adversarial examples, low empirical risk doesn't have to mean low adversarial risk

# Generating adversarial examples

Studying adversarial examples in the image domain essential, because: - perceptual similarity between fake and original is intuitive (unlike in other domains - graphs, audio) - imaga data have simple structure

## Studied datasets

- MNIST
- CIFAR10
- ImageNet

# White-box attacks

Given classifier $C$ (model $F$), victim sample $(x, y)$, synthesize fake image $x'$, that is perceptually similar to original $x$, but fools the classifier $C$ - find $x'$ satisfying $\|x' - x\| \le \epsilon$, such that $C(x') = t \ne y$ - $\|\cdot\|$ usually $l_p$ norm

## Biggio's attack

- adversarial examples on MNIST targeting SVMs and 3-layer FC DNNs
- inspired studies on safety of deep learning

## Szegedy's limited-memory BFGS (L-BFGS)

- first to attack image classifiers
- find minimally distorted adversarial example $x'$ by solving:
    - $\min \|x - x'\|_2^2$, s.t. $C(x') = t$ and $x' \in [0, 1]^m$
- Szegedy et al.