

A History of Adversarial Attacks

Prof. Stephen Roberts¹; Benjamin Etheridge¹

AI Security Reading Group, October 2024

¹Machine Learning Research Group
Department of Engineering Science

Table of Contents

Adversarial Examples

BNNs, GPs, Bayesian Opt

Table of Contents

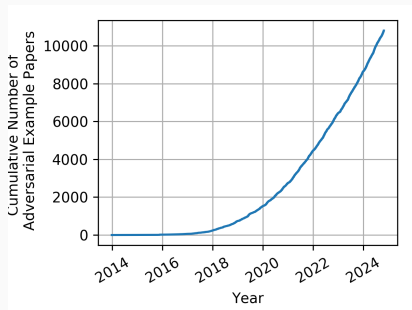
Adversarial Examples

BNNs, GPs, Bayesian Opt

Intro - Why should we care?

Intro - Why should we care?

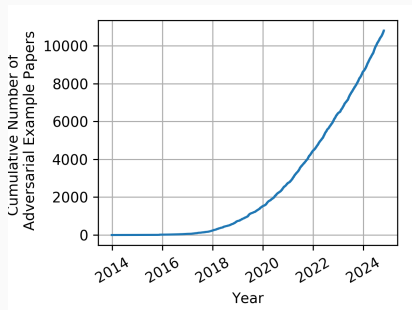
- Extremely active field



<https://nicholas.carlini.com/writing/2019/all-adversarial-example-papers.html>

Intro - Why should we care?

- Extremely active field



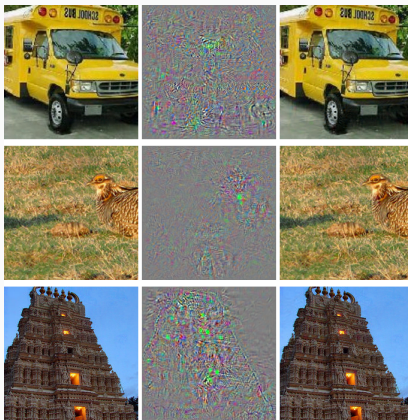
<https://nicholas.carlini.com/writing/2019/all-adversarial-example-papers.html>

- Threat model has changed - no longer purely an academic concern.

Intriguing properties of neural networks

Szegedy et al. - ICLR 2014 - arXiv:1312.6199

- Imperceptible perturbations can result in significant changes in output



Intriguing properties of neural networks

Szegedy et al. - ICLR 2014 - arXiv:1312.6199

- Box constrained L-BFGS Computation, minimising $\|L\|_2$ of perturbation η s.t:
 1. $f(x + \eta) = l$ (misclassification to label l)
 2. $x + \eta \in [0, 1]^m$ (?)

Intriguing properties of neural networks

Szegedy et al. - ICLR 2014 - arXiv:1312.6199

- Box constrained L-BFGS Computation, minimising $\|L\|_2$ of perturbation η s.t:
 1. $f(x + \eta) = l$ (misclassification to label l)
 2. $x + \eta \in [0, 1]^m$ (?)
- Second Order Optimisation - Inefficient

Intriguing properties of neural networks

Szegedy et al. - ICLR 2014 - arXiv:1312.6199

- Box constrained L-BFGS Computation, minimising $\|L\|_2$ of perturbation η s.t:
 1. $f(x + \eta) = l$ (misclassification to label l)
 2. $x + \eta \in [0, 1]^m$ (?)
- Second Order Optimisation - Inefficient
- ICLR 2024 Test Of Time Runner Up

Explaining and Harnessing Adversarial Examples

Goodfellow et al. - ICLR 2015 - [arXiv:arXiv:1412.657](https://arxiv.org/abs/1412.6575)

- Introduces Fast Gradient Sign Method (FGSM)

AUTHOR1 et al. YEAR - arXiv:REF



AUTHOR1 et al. YEAR - arXiv:REF



AUTHOR1 et al. YEAR - arXiv:REF



AUTHOR1 et al. YEAR - arXiv:REF



Table of Contents

Adversarial Examples

BNNs, GPs, Bayesian Opt

Sample frame title

In this slide, some important text will be highlighted because it's important. Please, don't abuse it.

Remark

Sample text

Important theorem

Sample text in red box

Examples

Sample text in green box. The title of the block is “Examples”.

This is a text in first column.

$$E = mc^2$$

- First item
- Second item

This text will be in the second column and on a second thought this is a nice looking layout in some cases.