

Adversarial Machine Learning: problem evaluation and implementation

Luca Ghislotti

University of Bergamo
Department of Engineering

Computer Security
AY 20/21

March 28, 2021

Overview

1 The problem

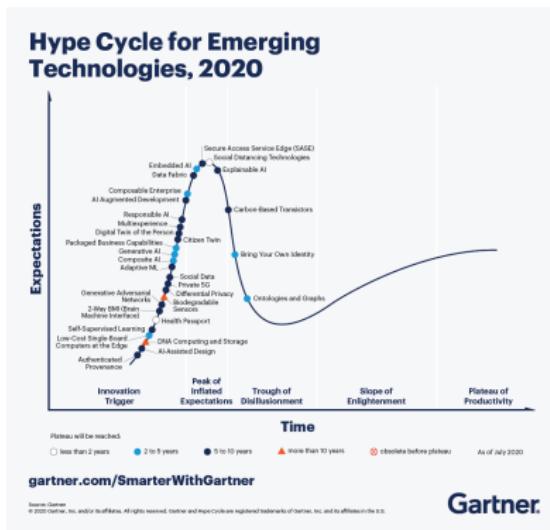
- A spot of context
- What do we mean with Adversarial Machine Learning?
- Why is it such an important issue?
- Different ways to implement an attack

2 Implementation: an example

- Project setup: the tools you'll need
- An example: a Ferrari that turns out to be a purse

3 References

A spot of context



As machine learning techniques have entered computing mainstream, their uses have multiplied. Online advertising and algorithmic trading are now inconceivable without machine learning, and machine learning techniques are increasingly finding their ways into health informatics, fraud detection, computer vision, machine translation, and natural language understanding.

Figure: Gartner hype cycle for emerging technologies

What do we mean with Adversarial Machine Learning?

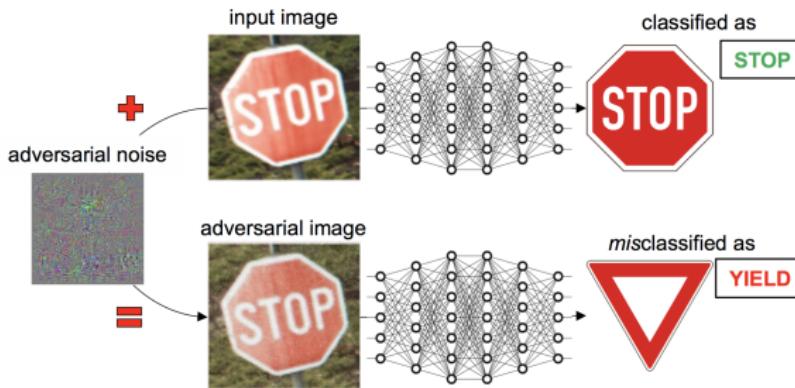


Figure: An example of image manipulation

Many applications of machine learning are adversarial in nature. Some are adversarial because they are safety critical, such as autonomous driving. An adversary in these applications can be a malicious party aimed at causing congestion or accidents, or may even model unusual situations that expose vulnerabilities in the prediction engine.

Why is it such an important issue?

One potential concern with learning algorithms is that they may introduce a security fault into systems that employ them. The key strengths of learning approaches are their adaptability and ability to infer patterns that can be used for predictions and decision making. However, these advantages of machine learning can potentially be subverted by adversarial manipulation of the knowledge and evidence provided to the learner.

In rapidly changing environments, machine learning techniques are considerably advantageous over handcrafted rules and other approaches because they can infer hidden patterns in data, they can adapt quickly to new signals and behaviors, and they can provide statistical soundness to a decision-making process.

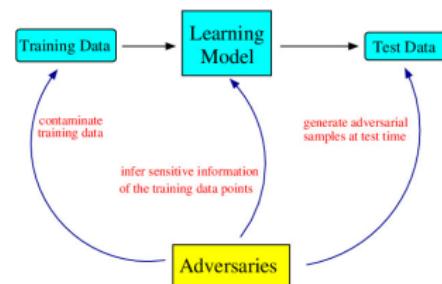


Figure: Different points of interaction between the adversary and the system

Different ways to implement an attack

We classify attacks along three dimensions: *timing*, *information*, and *goals*.

The *timing* dimension

The first crucial consideration in modeling attacks is *when* the attack takes place.

Attacks on models assume that the model has already been learned, and the attacker now either changes its behavior, or makes changes to the observed environment, to cause the model to make erroneous predictions. **Poisoning attacks**, in contrast, take place before models are trained, modifying a part of the data used for training.

Different ways to implement an attack

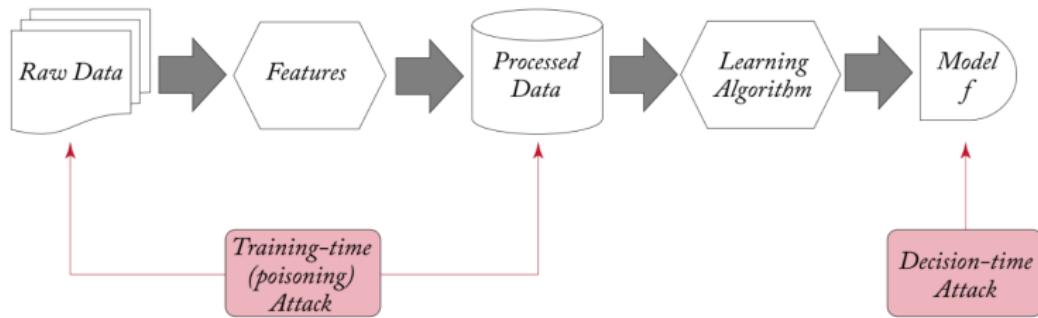


Figure: A schematic representation of the distinction between decision-time attacks (attacks on models) and poisoning attacks (attacks on algorithms)

Different ways to implement an attack

The *information dimension*

The second important issue in modeling attacks is what information the attacker has about the learning model or algorithm, a distinction which is commonly distilled into white-box vs black-box attacks. In particular, **white-box attacks** assume that either the model (in the case of attacks on decisions) or algorithm (in poisoning attacks) is fully known to the adversary, whereas in **black-box attacks** the adversary has limited or even no information at all.

Different ways to implement an attack

The *goals* dimension

Attackers may have different reasons for attacking, such as evading detection or reducing confidence in the algorithm. We differentiate two broad classes of attack goals: targeted attacks and reliability attacks. In a **targeted attack**, the attacker's goal is to cause a mistake on specific instances of a specific nature. A **reliability attack**, in contrast, aims to degrade the perceived reliability of the learning system by maximizing prediction error.

Implementation: an example

Project setup: the tools you'll need

In order to try out and test the project you'll need:

- **Python 3**
- **Tensorflow (2.2 required, 2.4.1 provided)** Free and open-source software library for machine learning. <https://tensorflow.org/>
- **Keras (2.2 required, 2.4.3 provided)** Open-source software library that provides a Python interface for artificial neural networks. <https://keras.io/>
- **Matplotlib (3.3 required, 3.3.4 provided)** Comprehensive library for creating static, animated, and interactive visualizations in Python. <https://matplotlib.org/>
- **Conda/Conda-Navigator GUI (optional)** Package, dependency and environment manager. <https://docs.conda.io/projects/conda/en/latest/index.html>

Everything's available for Windows/MacOS/Linux (Pop!_OS 20.10 Kernel 5.8 used)

An example: a Ferrari that turns out to be a purse

This example wants to be a very simple implementation of a type of attack, called *Fast Gradient Sign Method* or *FGSM*. The fast gradient sign method works by using the gradients of the neural network to create an adversarial example.

For an input image, the method uses the gradients of the loss with respect to the input image to create a new image that maximises the loss. This new image is called the adversarial image.

This can be summarised using the following expression:

$$adv_x = x + \epsilon \cdot sign(\nabla_x \cdot J(\theta, x, y))$$

Where:

- adv_x : adversarial image
- J : loss
- x : original input image
- y : original input label
- ϵ : multiplier to ensure the perturbations are small
- θ : model parameters

An example: a Ferrari that turns out to be a purse

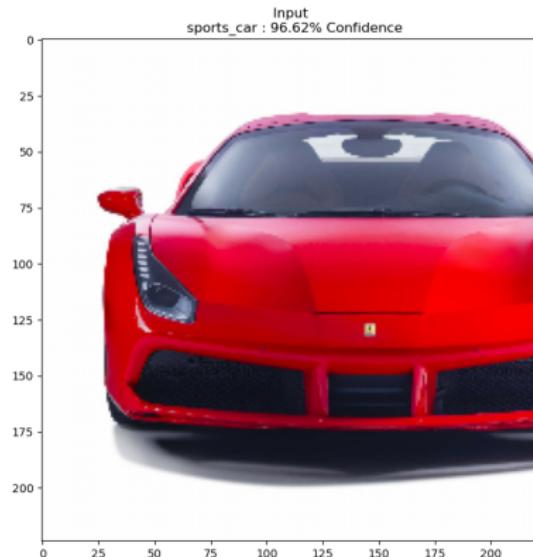


Figure: Input image

An example: a Ferrari that turns out to be a purse

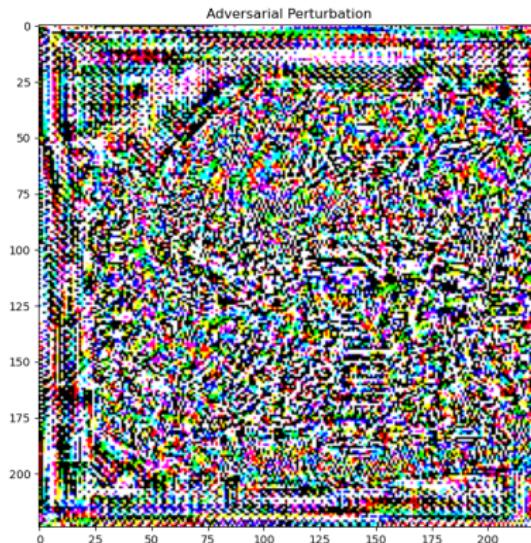


Figure: Perturbation matrix

An example: a Ferrari that turns out to be a purse

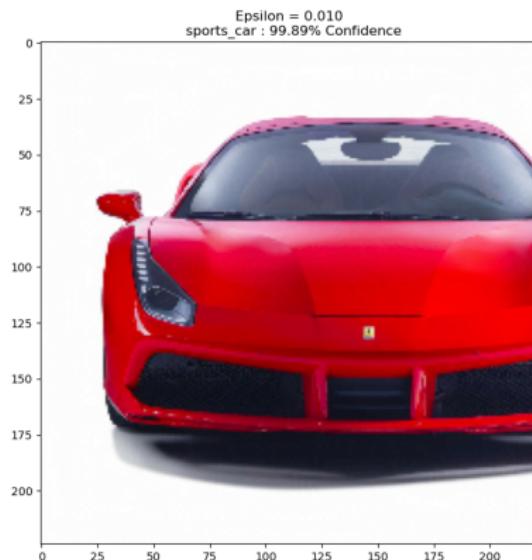


Figure: Output image after perturbation with $\epsilon = 0.010$

An example: a Ferrari that turns out to be a purse

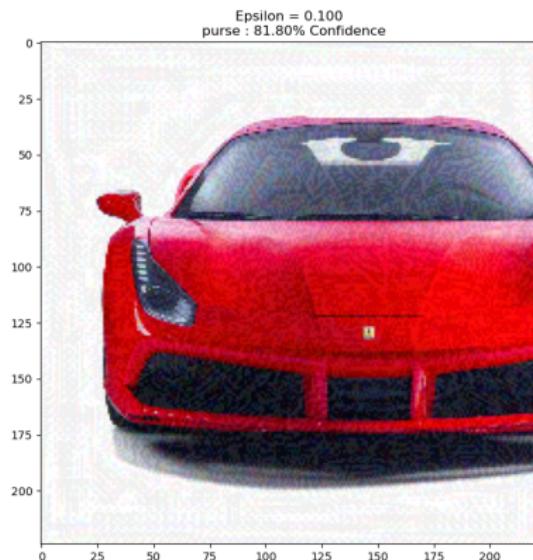


Figure: Output image after perturbation with $\epsilon = 0.100$

An example: a Ferrari that turns out to be a purse

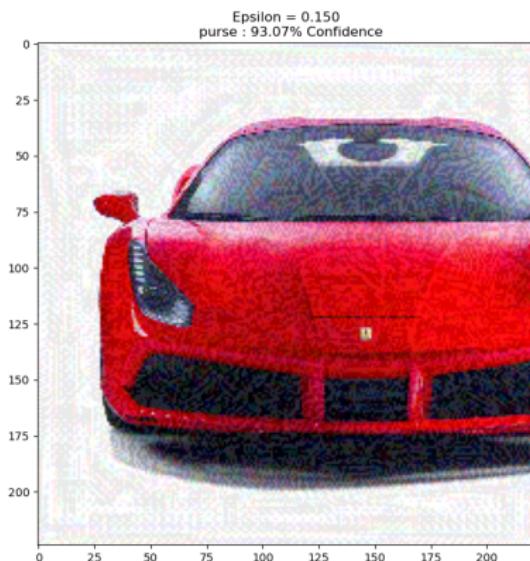


Figure: Output image after perturbation with $\epsilon = 0.150$

References

-  **Laskov, Pavel and Richard Lippmann**
Machine Learning in Adversarial Environments
Machine Learning, vol. 81, no. 2, Nov. 2010, pp. 115–19
-  **McDaniel, Patrick, et al.**
Machine Learning in Adversarial Settings
IEEE Security & Privacy, vol. 14, no. 3, May 2016
-  **Tygar, J. D**
Adversarial Machine Learning
IEEE Internet Computing, vol. 15, no. 5, Sept. 2011, pp. 4–6

A book suggestion

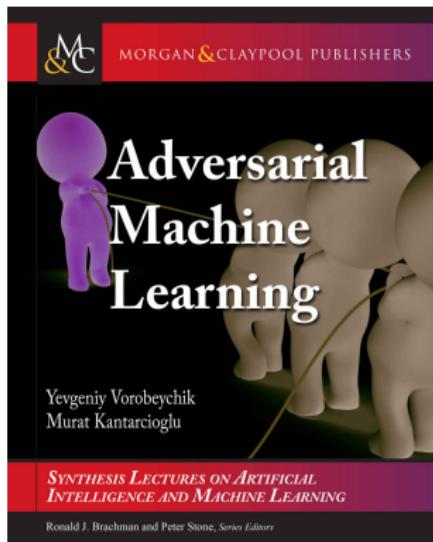


Figure: Adversarial Machine Learning: Synthesis Lectures on Artificial Intelligence and Machine Learning. Yevgeniy Vorobeychik and Murat Kantarcioglu. Morgan & Claypool, 2018

Useful links

- **GitHub Project**

https://github.com/lucaghislo/adversarial_machine_learning

- **LaTeX Project**

<https://it.overleaf.com/read/kfkrrytsqpzw>

The End