# Adversarial and Secure Learning 2019

Introductory information

Bojan Kolosnjaji, TUM I20
Ching-Yu Kao, Fraunhofer AISEC

# Today's agenda

1) Introduction to the research area

2) Seminar instructions

    a) Deliverables

    b) Grading

    c) FAQ

# Introduction to the research area

# Learning in adversarial environment

- Problem considered in the research community at least since early 2000s

- With the hype over machine learning (deep learning) the problem gains

  importance

- Adversarial perturbations studied in vision, text, malware...
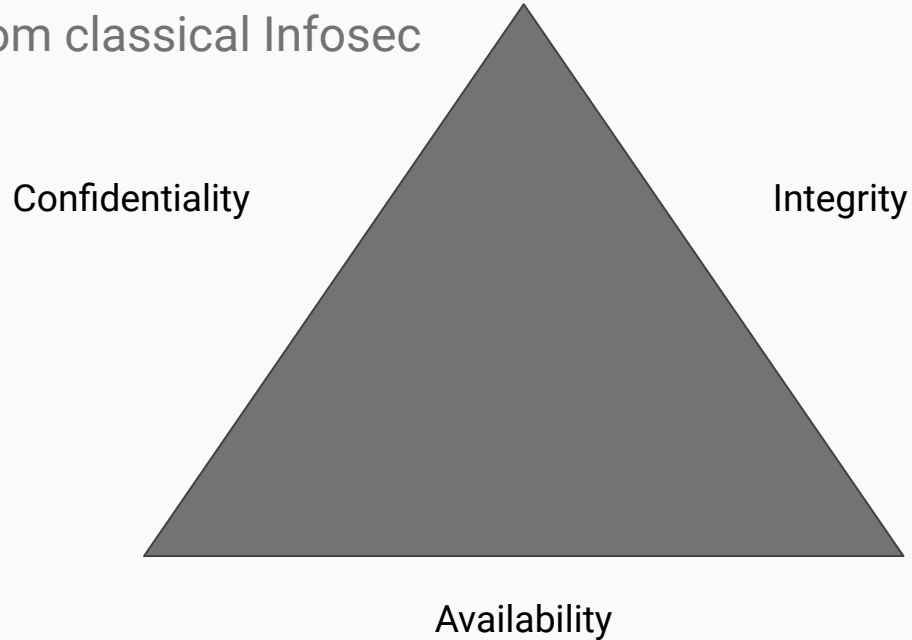
# Adversarial attacks

- Attack goals
  - Destroy system functionality
  - Change the model output
  - Reveal confidential information
- Attacker knowledge
  a. Perfect Knowledge (white-box attack)
  b. Limited Knowledge (gray-box attack)
  c. Zero Knowledge (black-box attack)
- Attacker capability
  a. Make queries to the model
  b. Change training data
  c. Change validation/test data

# Adversarial attacks

- Evasion attacks - modify test input to get a desired classifier decision

- Poisoning attacks - modify training data to change the resulting model (either just make it less accurate or enable particular properties (backdoor)

- Privacy attacks - learn about sensitive data from the model based on queries

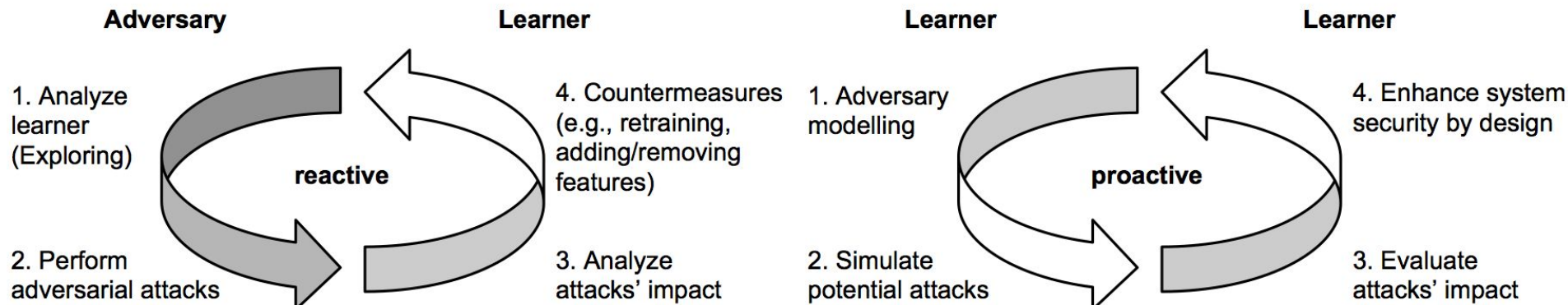# Adversarial attacks

- CIA Triad - from classical Infosec

Confidentiality

Integrity

Availability

# Attack and defense

- Attack optimization -> find minimal perturbation that achieves the goal

  (min g(x) , s.t. d(x,x')<$d_{max}$ (x,x'))
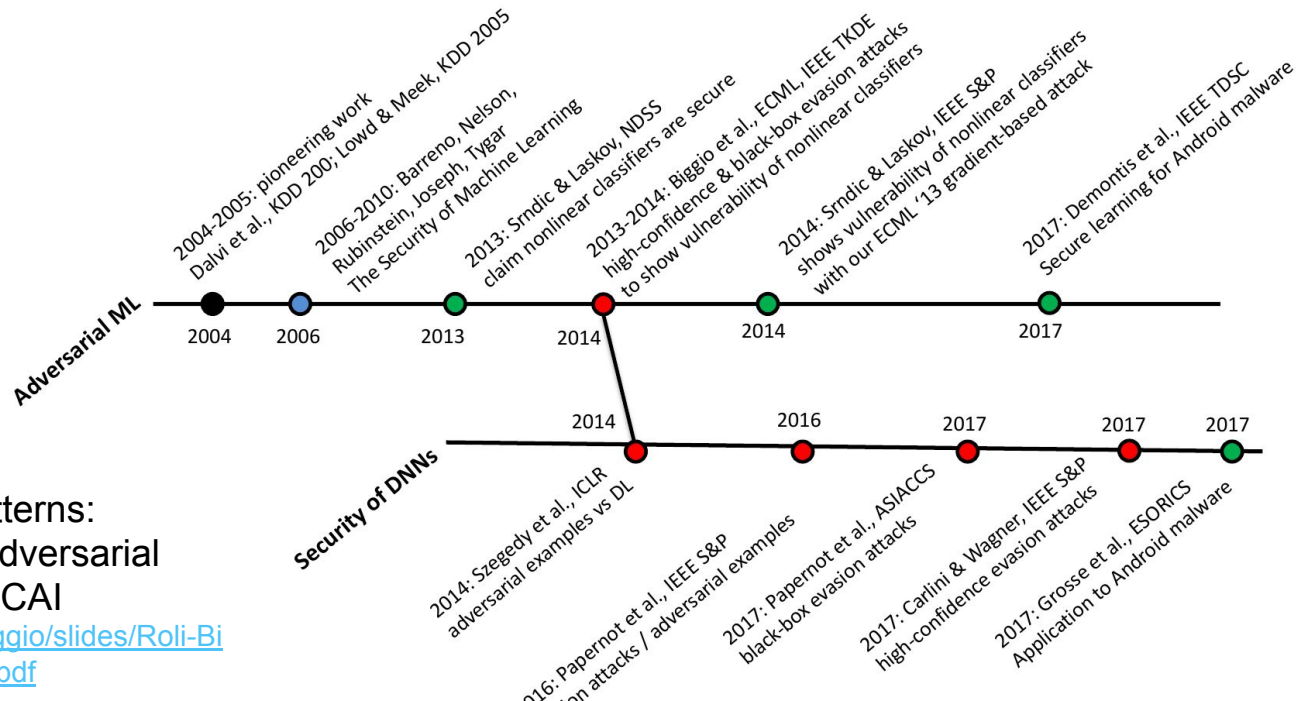
- Defense optimization -> find an approach that either:
  - Rejects the attack points as outliers or is
  - Robust to attack points

# Defense strategies

- Kerckhoffs principle - do not rely on obscurity
- Assess security against various levels of adversary's knowledge and capability
- Arms race

# Adversarial Attack -discoveries



Roli, Biggio, "Wild Patterns:
Half-day Tutorial on Adversarial
Machine Learning", IJCAI
http://www.diee.unica.it/~biggio/slides/Roli-Biggio-IJCAI-ECAI18-tutorial.pdf
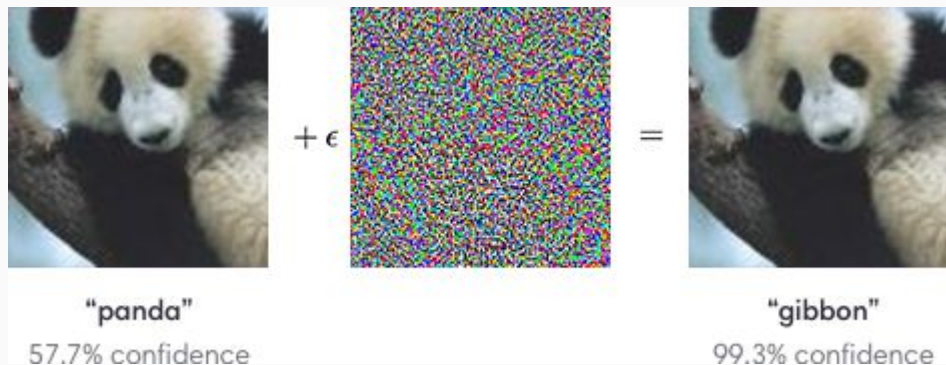
# Adversarial Learning

- Attacked methods
  - Linear regression
  - **SVM**
  - **Neural Networks**
  - Random Forests
  - ...

# Adversarial Learning

- Data under modification
  - **Computer Vision**
  - **Natural Language Processing**
  - Speech
  - **Malware**
  - ...

# Vision papers best known...



"panda"
57.7% confidence

"gibbon"
99.3% confidence

Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." *arXiv preprint arXiv:1412.6572* (2014).
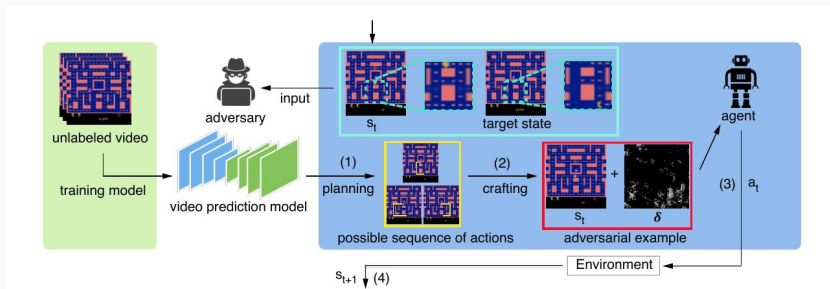


Papernot, Nicolas, et al. "Practical black-box attacks against deep learning systems using adversarial examples." *arXiv preprint* (2016).

# But also other domains...



Mei, Shike, and Xiaojin Zhu. "The security of latent dirichlet allocation." *Artificial Intelligence and Statistics*. 2015.



Lin, Yen-Chen, et al. "Tactics of adversarial attack on deep reinforcement learning agents." *arXiv preprint arXiv:1703.06748*(2017).

# Machine Learning in getting adoption in AV industry

Microsoft | TechNet

Search

## Microsoft Malware Protection Center
Threat Research & Response Blog

# Windows Defender: Rise of the machine (learning)

Rate this article ★★★★★

November 16, 2015   By msft-mmpc

f 0    y 0    in 0    💬 5

Windows Defender harnesses the power of machine learning, contributing to making Windows 10 Microsoft's most secure client operating system and providing increased protection against security threats facing consumers and commercial enterprises today.

# Even deep learning...

## Symantec Adds Deep Learning to Anti-Malware Tools to Detect Zero-Days
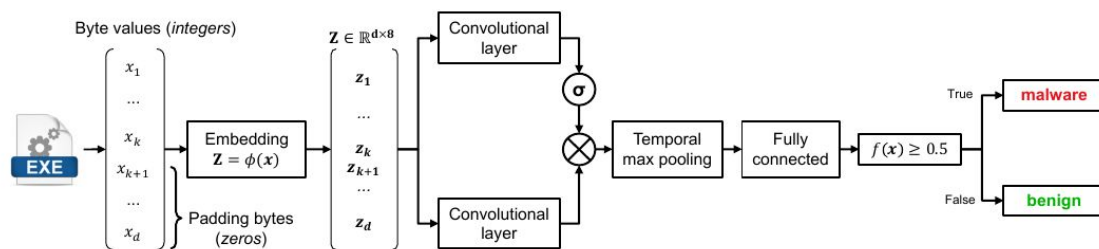
### Robotics

## Antivirus That Mimics the Brain Could Catch More Malware

**Computer malware can often evade antivirus security software if the** author changes a few lines of code or designs the program to automatically mutate before each new infection.

Artificial neural networks, trained to recognize the characteristics of malicious code by looking at millions of examples of malware and non-malware files, could perhaps offer a far better way to catch such nefarious code. An approach known as deep learning, which involves

# Some current work at Chair I20

- Generating adversarial binaries to evade neural network malware detectors based on raw code



Algorithm 1 Adversarial Malware Binaries

**Input:** $\boldsymbol{x}_0$, the input malware (with $k$ informative bytes, and $d-k$ padding bytes); $q$, the maximum number of padding bytes that can be injected (such that $k + q \le d$); $T$, the maximum number of attack iterations.

**Output:** $\boldsymbol{x}'$: the adversarial malware example.

1: Set $\boldsymbol{x} = \boldsymbol{x}_0$.
2: Randomly set the first $q$ padding bytes in $\boldsymbol{x}$.
3: Initialize the iteration counter $t = 0$.
4: **repeat**
5:     Increase the iteration counter $t \leftarrow t + 1$.
6:     **for** $p = 1, \ldots, q$ **do**
7:         Set $j = p + k$ to index the padding bytes.
8:         Compute the gradient $\boldsymbol{w}_j = -\nabla_\phi(x_j)$.
9:         Set $\boldsymbol{n}_j = \boldsymbol{w}_j / \|\boldsymbol{w}_j\|_2$.
10:        **for** $i = 0, \ldots, 255$ **do**
11:           Compute $s_i = \boldsymbol{n}_j^\top (\boldsymbol{m}_i - \boldsymbol{z}_j)$.
12:           Compute $d_i = \|\boldsymbol{m}_i - (\boldsymbol{z}_j + s_i \cdot \boldsymbol{n}_j)\|_2$.
13:        **end for**
14:        Set $x_j$ to $\arg \min_{i:s_i>0} d_i$.
15:     **end for**
16: **until** $f(\boldsymbol{x}) < 0.5$ or $t \ge T$
17: **return** $\boldsymbol{x}'$

# Seminar instructions

# Our goal

- Get an **overview** of the academic state-of-the-art

- Extend your **knowledge** in machine learning, learn to look at it from the

  **security** standpoint

- Get a feeling on **how to evaluate risk**

# What to deliver?

- 1 presentation
  - 45 minutes + 15 minutes of discussion
  - Present the topic, use all papers

- 1 report
  - 14 Pages LNCS (look up LNCS template in Latex)

Presentation (start on 14.05.)

# Presentation

Needs to be:

- Correct
- Complete
- Comprehensible

# Presentation - Correct

- Present information from the paper correctly

- Don't speculate without a reason or proof

- Don't claim something you cannot explain well

# Presentation - Complete

- Explain **all** key points of the papers

- Be careful about **time constraints** and distribution

- **Convey information** without leaving out important insight

# Presentation - Comprehensible

- Speak loud and clear

- Think about the audience - fellow students

- Motivate the audience for discussion

- Don't fight your audience, answer all questions friendly

# Presentation - Concise Text

- A PPT is just a presentation aid. It should not be a paper in its own right and your bullet points should be, if possible, less than a line long. Specifically, keep bullet points short by making them clear and concise. Do not be afraid from using incomplete sentences or phrases. In reality, this is the preferred method because it helps to highlight the points you are making during your talk.
- This is because having lots of text on your slides makes it difficult to understand the point you are trying to make. Furthermore, your audience will end up reading the text and ignoring you.

# Presentation - Emphasis

- ONLY USE CAPITALIZATION WHEN NEEDED
- Use color sparingly
- **Only bold** key **words and phrases**
- Light text on a dark background is bad

# Presentation - Pictures

# Presentation - Structure

- Introduction to the topic

- Present all papers together

  - Introduction

  - Main Points

  - Back up arguments

  - Conclusions (key takeaways)

# Presentation - Audience

- Read papers, or at least abstracts, prior to each presentation day

- Listen carefully, write down questions

- Ask questions, comment

- Active participation is appreciated!

# Presentation - Grading

- Presentation skills
  - General organization, use of slides
  - Language, slide text and graphics
  - Pace, use of time
- Subject-related competence
  - Subject knowledge
  - Staying on topic
  - Identifying interesting/important points

Report (deadline 09.06.2019.)

# Report

- 14 Pages LNCS
- Summarize key points of all papers together - not an easy task
- Use a typical paper structure:

Abstract -> Introduction -> Methodology -> Results -> Discussion -> Conclusion

# Report - Abstract

- Summarize the paper

  - Introduction to the problem

  - How was the problem solved? Methodology

  - Short insight in the results

  - What is the impact of the papers?

# Report - Introduction

- Describe the context

- What is the preexisting work?

- What does the preexisting work lack?

- How do the papers close the gap?

# Report - Methodology

- Describe the mechanisms used to tackle the existing problem

- Lead the reader through the problem solving procedures

- Give arguments for the choice of methods

# Report - Results

- Give an overview of the important results

- Add tables, graphs... if you have space

- Shortly comment on the figures

- Avoid phrases like: It is obvious from this graph that ...

# Report - Discussion

- What do the results actually tell us?

- Compare the results with related work

- What are the limitations of the paper results?

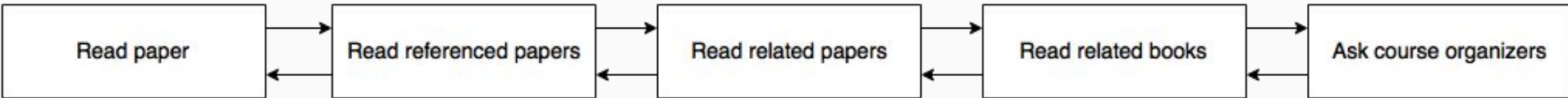- How can the limitations be addressed?

# Report - Conclusion

- Summary of the main findings in 3-4 sentences

- What are the most interesting results?

- What is the impact of the papers?

# Report - Grading

- Paper organization
- Language and grammar
- Subject knowledge
- Ability to summarize
- Proper bibliography and citations (!)

# How to do your research

- Seminar - (kind of) simulation of scientific research
- Try to be independent, but also ask questions

| Read paper | → | Read referenced papers | → | Read related papers | → | Read related books | → | Ask course organizers |

# FAQ

- Allowed to miss a presentation day? Yes, if you have a very good reason.

  - Examples of good reason: health issues, schedule clashes at the Uni

  - Examples of bad reason: HiWi work, homework, football training, bad mood

- Can I set a meeting if I have problems with my papers?

  - Yes, but try to do as much as you can yourself.