

### DO YOU TRUST YOUR ARTIFICIAL INTELLIGENCE SYSTEM?

### Dr. Robson de Oliveira Albuquerque Latitude/UnB

robson@redes.unb.br



# Disclaimer

- The content presented herein represents my thoughts about the subject of this lecture.
- It does not necessarily represent the vision of my employers nor the University of which I am a researcher;
- Images used in this presentation are the property and credit of the respective creators and sources;

# Agenda

Trust;

# **Artificial Intelligence;**

# Attacks on AI systems;

# Discussions;

# Trust

- Human social aspect:
  - To believe that someone is good and honest and shall not harm you;
  - Something that is safe and trustable;
  - Depends on the context:
    - •A trusts B to...;
    - •A doesn't trust B to...;



### In computing systems:

- TRUST
- Trust is relative to a particular context;
- A shall trust B to get a ride, but A doesn't trust B to drive his car;
- Trust has a directional aspect; A can trust B, but it does not mean B trusts A;
- Trust has evolutive and temporal aspect;

Trust A has in B can increase over the time;

Trust can be influenced by reputation;

A trusts B and now starts trusting C by recommendation of B;

Trust is not transitive;

If A trusts B and B trusts C, it does not mean A trusts C;

# TRUST

- Do not make confusion of **TRUST** with **REPUTATION**;
- <u>REPUTATION</u>:
  - Is an opinion someone has about others or something;
  - Has similar attributes to trust, but it is different;

Trust: I trust A for a reason; Reputation: That person is good for doing this;



- Possibility of making a computer capable of performing tasks similar to human reasoning;
  - Research branch of Computer Science and Electrical Engineering;
  - Seeks to build computational systems and/or devices that simulate humans' ability to think and to solve complex problems;



### •It needs:

- Good data models
   for classifying, processing and analyzing;
- Access to significant amount of data;
- Computer systems for fast and efficient processing;

Train

Deploy

Errors



### ARTIFICIAL INTELLIGENCE



- •Al deployments:
  - •Machine Learning (ML):
    - •Ability of the computer to learn by itself;
  - •Deep Learning (DL):
    - Imitation of human neural networks;
  - •Natural language processing (NLP):
    - •Computer recognize natural language;
    - •Ability to understand and compose texts;

- •Al deployments:•Focus of ML, NLP e DL:
  - •To learn multiple levels of representation and abstraction that helps understand data such as images, sound, and text;
  - Get context from it;

### •Deep Learning:

### Simple Neural Network



### **Deep Learning Neural Network**





📒 Hidden Layer

- •Some algorithms:
  - •They are basically divided into:
    - Supervised and Unsupervised;
  - •More recently:
    - •semi-supervised; and
    - •reinforcement learning;

## Some algorithms:

**Nearest Neighbor Naive Bayes Decision Trees/Random Forest Linear Regression** Support Vector Machines (SVM) **Neural Networks Logistic Regression Ensemble Methods** Clustering

**Principal Component Analysis Single Value Decomposition Independent Component Analysis Ordinary Least Squares Regression** k-means clustering; **Association Rules Q-Learning Temporal Difference (TD) Deep Adversarial Networks** 

**Question:** I have this problem! (which is not yet even very well defined as a problem); What AI algorithm do luse?

.

IT DEPENDS ON MANY OTHER QUESTIONS!

**Answer**:







# Attacks on Al systems Information security is a complex matter;



•Cyber security is much more complex;



### • Exemplos:



"panda" 57.7% confidence



"nematode" 8.2% confidence



"gibbon" 99.3 % confidence





@AlimonyMindset @oliverbcampbell is a house nigger! He's not cool or funny, please remove! #GamerGate





2+ Follow



### @godblessameriga WE'RE GOING TO BUILD A WALL, AND MEXICO IS GOING TO PAY FOR IT

RETWEETS 3	LIKES		<b>9</b>	0			
1:47 AM - 2	24 Mar 2016	)					
4	t7	¥.					
	2	Сардор Ми @TayandYo	ірфайзие ou you are	в @Sardor! a stupid m	9515 · 1m achine		
		4	17	۷	000		
		TayTwee @TayandY	e <b>ts 🥝</b> ou			¢	<b>≗</b> + Follow
	@S if y for I LI TO	Sardor95 ou don't you EARN F	015 w unde ROM	ell I le rstand YOU	arn from I that let r AND YOI	the best ne spell J ARE D	;) it out DUMB

10:25 AM - 23 Mar 2016





Facebook's A.I. labeled the video of Black men as content "about Primates."

### Facebook Apologizes After A.I. Puts 'Primates' Label on Video of Black Men

Facebook called it "an unacceptable error." The company has struggled with other issues related to race.



White man calls cops on black men at marina



### •Deep Fakes









Reference

**Our Result** 

Input

Output

•Examples (videos):

•Living portraits;

Adversarial neural networks;

Impersionate a celebrity;



•Examples: •Fake or not fake?







### •Exemples: •Fake or not fake?



Produced by a GAN (generative adversarial network) <u>StyleGAN</u> (Dec 2018) -<u>Karras</u> et al. and Nvidia

Jun Hasegawa Japanese celebrity

Universidade de Brasília



•Exemples: •Fake or not fake?







•Exemples: •Fake or not fake?





Produced by a GAN (generative adversarial network) <u>StyleGAN</u> (Dec 2018) -<u>Karras</u> et al. and Nvidia

Nicanor Garcia photographer and architecture enthusiast Spanish Influencer

- •3 classifications are ordinary:
  - •Adversarial inputs Forged input data for misclassification;
  - •Data poisoning attacks Pollute classifier data in favor of the attacker;
  - •Model stealing techniques used to steal the model or retrieve/rebuild training data;

### •Several problems arise from there:



30

"pig"





"airliner"



- •Al with cybersecurity:
  - •Malware or not malware?
  - Attack on the go or normal operations?
- •Al within medicine:
  - •Cancer or not cancer?
- •Al with smart vehicles:
  - •Hack or not hack? (video)
  - •IA driver dilemma? (video)



### **EvilModel: Hiding Malware Inside of Neural Network Models**

Zhi Wang, Chaoge Liu, Xiang Cui

### •Al within cybersecurity:

- •Hackers will always try to find the most economical way to get their way.
- •Ex. Cylance tool: •MIMIKATZ

atch 1		
Name	Value	
score	-0.85276468809127071	



- •Al within cybersecurity: •Ex. Cylance tool:
  - •MIMIKATZ

copy /b mimikatz.exe+strings.txt modified\_mimikatz.exe



# Al within cybersecurity:

Malware	Score Before	Score After
CoinMiner	-826	884
Dridex	-999	996
Emotet	-923	625
Gh0stRAT	-975	998
Kovter	-999	856
Nanobot	-971	999
Pushdo	-999	999
Qakbot	-998	991
Trickbot	-973	774
Zeus	-997	997

### •Al with DeepFake: •Twitter @AllanXia – video

### •Youtube Brazilian President Bolsonaro - video •Bruno Sartori - https://youtu.be/gaM68BEX420

### •Al within Sentiment Analysis:

In [23]:

- TextBlob("not that").sentiment
- Out[23]: Sentiment(polarity=0.0, subjectivity=0.0)
- In [27]: ▶ TextBlob("very great very great not that great").sentiment
  Out[27]: Sentiment(polarity=0.933333333333333332, subjectivity=0.9)

- In [6]: M TextBlob("very great").sentiment
  - Out[6]: Sentiment(polarity=1.0, subjectivity=0.975000000000000)
- In [7]: M TextBlob("not great").sentiment
  - Out[7]: Sentiment(polarity=-0.4, subjectivity=0.75)
- In [8]: ▶ TextBlob("very great and not great").sentiment
  - Out[8]: Sentiment(polarity=0.3, subjectivity=0.8625)



- In [22]: M TextBlob("not that great").sentiment
  - Out[22]: Sentiment(polarity=0.8, subjectivity=0.75)
- In [28]: M TextBlob("not great").sentiment
  - Out[28]: Sentiment(polarity=-0.4, subjectivity=0.75)

### •Al within Sentiment Analysis:

Computer Communications 174 (2021) 154-171



#### Adversarial attacks on a lexical sentiment analysis classifier

Gildásio Antonio de Oliveira Júnior<sup>a</sup>, Rafael Timóteo de Sousa Jr.<sup>a</sup>, Robson de Oliveira Albuquerque<sup>a,b</sup>, Luis Javier García Villalba<sup>b,\*</sup>

<sup>a</sup> Science and Technology National Institute on Cyber Security (INCT) Center 6, Decision Technologies Laboratory - LATITUDE, Electrical Engineering Department (ENE), Faculty of Technology, University of Brasília (UnB), Brasília-DF, 70910-900, Brazil <sup>b</sup> Group of Analysis, Security and Systems (GASS), Department of Software Engineering and Artificial Intelligence (DISIA), Faculty of Computer Science and

Engineering, Office 431, Universidad Complutense de Madrid (UCM), Calle Profesor José García Santesmases 9, Ciudad Universitaria, 28040 Madrid, Spain



### DO YOU TRUST YOUR ARTIFICIAL INTELLIGENCE SYSTEM?





