

Adversarial Machine Learning: *How Secure Is Machine Learning?*

Competence Centers for Excellent Technologies





Rudolf Mayer, <u>rmayer@sba-research.org</u> Tanja Šarčević, <u>tsarcevic@sba-research.org</u>

irtschaft

Bundesministerium Klimaschutz, Umwelt, Energie, Mobilität, Innovation und Technologie **Bundesministerium** Digitalisierung und Wirtschaftsstandort











AI/ML/DL is everywhere

• AI, ML and Deep Learning among most hyped technologies

- A lot of hype, **AND** tremendous advances
 - Surpassing human-level performance on a number of tasks
 - Advances based on a number of new learning concepts
 - Many application areas:



• What about security?







Setting

• Learning paradigms/domains considered



Attacks

• Attack vectors specific to Machine Learning



• How to secure Machine Learning

Adversarial Machine Learning

4

Machine Learning Pipeline

Machine Learning Workflow



- Two steps:
 - **Training (offline)**:

estimate model parameters **O** from **X** and **Y**

• **Prediction**: apply Θ in prediction function $f_{\Theta}: X \to Y$

Specific ML Setting considered

- Classification / categorisation
 - Assign samples to a predefined list of categories
 - o Input
 - Vectors X (n-dimensional, real numbers)
 - Labels $Y = \{0, 1\}$
 - Space separated by prediction function *f* (*decision boundary*)



ML & Security

- Perfect security is difficult (impossible) to achieve
 - Goal: raising the threshold for an attack to be successful
 - → Balancing the cost of protection with the cost of recovering from an attack
- ML systems can be subject to attacks on their *integrity, availability or confidentiality*





7

ATTACKS AGAINST MACHINE LEARNING

Types of Attacks & Attack Vectors



Classification: Confidential

SBA Research

Security & Machine Learning

- Recent topic: Adversarial machine learning
 - Attacks & defences
 - History of approx. 15 years
 - Adversarial examples lately gained a lot of publicity



- Machine Learning historically: rather focused on optimising accuracy / generalisation power
 - Security was not a major topic: assumed training data comes from natural or well-behaved distribution
 - Does not generally hold in security-sensitive settings
 - Adversaries not considered

9

Vulnerabilities and Attacks

- Different attack vectors on the Machine Learning process ullet
 - Training and/or prediction phase 0





Poisoning (Backdoor) attacks







Evasion Attack: Adversarial Examples

- *Fooling* model in the prediction step
 - Minimal perturbation of an input leads to misclassification
 - Often not perceptible to human vision!



- Effective and robust
 - Small perturbation sufficient for successful attack
 - Able to attack other models besides Deep NNs!
 - Often resistant against digital → analog → digital conversion (e.g. scanning a printout)
- Attacks against *integrity* of prediction



Adversarial Examples: Simple Data

- Adversarial input generated using various algorithms
 - Needs to query the model
 - Simple approach: greedy search for decision
 boundary by changing
 pixels (minimising changes)



More advanced: Fast Gradient Signs, Iterative FGS, C&W, ... adv. label 1 9 5 4 3 4 7 8 1 1
FGS
IFGS
IFGS</

Adversarial Examples: More Complex Data

- Adversarial examples for object recognition
 - Perturbations often invisible to human perception
 - Maybe perturbation visible, but not recognized as relevant by human



Adversarial Examples: More Complex Data

- Adversarial examples for object recognition
 - Perturbations often invisible to human perception
 - Maybe perturbation visible, but not recognized as relevant by human





15

Adversarial Examples: Critical Threat



"panda"

57.7% confidence





"gibbon" 99.3 % confidence





89.2% confidence



"Priority" 92.4% confidence



Grosse et al. Adversarial Perturbations Against Deep Neural Networks for Malware Classification: https://arxiv.org/abs/1606.04435

Adversarial Attack: Demo

d∰o ART - IBM Research × +			- a ×	
← → C â art-demo.mybluemix.net		९ 🛧 🞯 🕈 🖪 🚺	🔤 💿 📄 🗯 🛒 🔒 🗄	
Try it out 1. Select an image to targe				*
and the second s				
2. Simulate Attack	Adversarial name type Fact Grapient Method V Original	Visua; Code Madified		ζ
Determine strength	None low too too			
3. Defend attack			N	
Gaussian Noise				
Spatial Smoothing	None and the second sec			
Feature Squeezing	None La ma ar			*

https://art-demo.mybluemix.net/

Defending Adversarial Attacks

Can we defend against these attacks?





Defences against Attacks on ML

- Defence against adversary is often an arms race
- Adversary is often "in the drivers seat"
 - $_{\circ}$ Decides which data to present to model \ge
 - Training data hard to verify / sanitise
 - Often direct access to model / parameters / service \vec{s}
- Often a trade off: security vs. model effectiveness or user experience ("cost")
- Operational vs. integrated defences (model robustness)



Defending Adversarial Attacks: Model Robustness

- Training a classifier robust to adversarial attacks
 - By pro-actively generating adversarial inputs
 - Letting the classifier learn these inputs \rightarrow "Harden" classifier
 - Impacts clean sample performance
 - Mostly effective against anticipate attack algorithm (e.g. FGSM, C&W, ..)
- Cleansing data inputs
 - Blurring or other image manipulation approaches
 - Passing it through an auto-encoder

→Embedded patterns might be removed



ATTACKS AGAINST MACHINE LEARNING

Poisoning & Backdoor



Classification: Confidential

SBA Research

Poisoning and Backdoors

- Attacks manipulating the learning model
 - Manipulation using some inputs, creating "poisoned" training data
 - **Generally for one class** (1-50% of those samples)





- Attacker requires access to training data or model
 - → *Supply chain attack*
 - E.g. when training in the cloud, using a pre-trained model in transfer learning, ...
- Attacks against integrity of model





Backdoored Neural Networks (BadNet)

Benign Network



BadNet

Behave **identically** on **clean** inputs



Clean Input



Gu et al. BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain. ML and Security 2017



Backdoored Neural Networks (BadNet)

Benign Network



BadNet



Backdoored Input



BadNets misbehave on backdoored inputs....

Backdoors: Simple Data

• Backdoor in the form of a pixel (or pixel pattern) on MNIST dataset



Very effective, without affecting classification of clean examples too
 much



Backdoors: Realistic Threat

- Poisoning of traffic-sign recognition
 - Often targets state-of-the-art **Convolutional NNs**
 - Backdoor symbol is noticeable, but not suspicious







Gu et al. BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain. ML and Security 2017



Backdoored Neural Networks

• Why do backdoors work?

Models in general have too much memory capacity!



- Comparing clean versus backdoored activations:
 - Some neurons active only on backdoor inputs
 - "Backdoor neurons"

Defending Poisoning/backdoor Attacks

Can we defend against these attacks?





Backdoors: Pruning Defence



- Defender prunes not-activated neurons
 - Identified using **validation** data (if available!)



Pruning Defence: Face Recognition





Pruning Defence: Traffic Sign



ATTACKS AGAINST MACHINE LEARNING

Data and Model confidentiality



Classification: Confidential

SBA Research

Confidentiality of Training Data

- Membership Inference Attack
 - Identify whether a sample was used to train a specific ML model
- Model Inversion Attack, e.g. against Face recognition model
 - Can an adversary use a model to recover images of training members?
 - **Reconstructs** input data for specific class (person)
 - Not perfect, yet scary 80% of faces recognized by humans







TARGET DATA

An adversary wants to know if some

Data Record was in the training set

of a target model

TARGET

MODEL

ATTACK MODEL

2. predVec

4. in / out

1. Data Record

3. predVec+Label

Model Extraction/Stealing

- Adversary wants to learn close approximation of model in as few queries as possible
 - Target: f'(x) = f(x) on ≥99.9% of inputs



- Efficient attacks can:
 - Undermine pay-for-prediction *(AI-as-a-Service)* model
 - Facility privacy attacks
 - Enable evasion attacks

Defending Model Stealing Attacks

Can we defend against these attacks?



Reactive & Proactive Defence: Model Watermarking

- Owner marks a model using their "signature"
 - Model verification!
- Signature (watermark) as:
 - Modification in model parameters (white-box access)
 - Set of adversarial inputs (black-box access)



Conclusions

Conclusions

- Machine Learning needs to consider security & privacy
 - Can get easily fooled & exploited
- Attacks can compromise:
 - Confidentiality (e.g. model inversion)
 - Integrity & Availability
- Supply chain needs to be considered
 - As-a-service, transfer learning from existing models, ...

• Adversaries are everywhere!



37

Questions?

MLDM



- Rudolf Mayer, rmayer@sba-research.org \bullet
- Tanja Šarčević, tsarcevic@sba-research.org •
- https://www.sba-research.org/research/mldm/ 38 \bullet