<u>Al Model Inspector</u>: Towards Holistic Adversarial Robustness for Deep Learning



Pin-Yu Chen (IBM Research) <u>www.pinyuchen.com</u> @pinyuchenTW CSIRO Sep. 2022 IEM Let's Create Products & Solutions∨ Consulting & Services∨ Learn & Support∨ Explore more∨

The man who challenges AI every day

Researchers like Pin-Yu Chen are exploring the enormous potential AI holds for society





Search

Research Vision: *AI Model Inspector*

- A machine-learning-driven automated pipeline that **proactively** and **continuously** identifies and mitigates error-prone risks hidden in our AI systems.
- The inspection spans the lifecycle of an AI model, from data collection and processing, model selection, training and testing, to model deployment and system integration.
- Creating a trustworthy AI ecosystem featuring self-correction and agility.

IBM

IBM Research Blog Topics ∨ Labs ∨ About

Preparing deep learning for the real world – on a wide scale

December 15, 2020 | Written by: Pin-Yu Chen



The gap between AI development and deployment

How we develop AI



How we deploy AI



Al revolution is coming, but *Are We Prepared*?

- According to a recent Gartner report, 30% of cyberattacks by 2022 will involve data poisoning, model theft or adversarial examples.
- However, industry is underprepared. In a survey of 28 organizations spanning small as well as large organizations, 25 organizations did not know how to secure their Al systems.



DEFENSE

Pentagon actively working to combat adversarial Al

Why adversarial (worst-case) robustness matters?

- Prevent prediction-evasive manipulation on deployed models
 - Build trust in AI: address inconsistent decision making between humans and machines & misinformation
- Assess negative impacts in high-stakes, safety-critical tasks
- Understand limitation in current machine learning methods

April 23, 2013 at 4:31 p.m. EDT

- Prevent loss in revenue and reputation
- Ensure safe and responsible use in AI

TESLA AUTOPILOT —

Researchers trick Tesla Autopilot into steering into oncoming traffic

Stickers that are invisible to drivers and fool autopilot. DAN GOODIN - 4/1/2019, 8:50 PM



Syrian hackers claim AP ha \$136 billion. Is it terrorism?	ck that tipped stock market
INDU 146690.06 122.89 16688.12/14692 At 13:34 0.14567.17/14720.34 14557.29 17000 NOU Income 10.300 10.300 10.300 10.300 10.30 10.300 10.300 10.300 10.300 10.300 10.300 10.30 10.400 10.300 10.400 10.300	81 intraday Chart sourny/Study 14720 14720 14720 14750 14750 14750 14750 14750 14750 14550 14550
AP The Associated Press AP The Associated Press The Associated P	14660 14640 14620 14600 14680









Microsoft silences its new A.I. bot Tay, after Twitter users teach it racism [Updated]

Sarah Perez @sarahintampa / 10:16 am EDT • March 24, 2016





Microsoft's • newly launched A.I.-powered bot called Tay, which was responding to tweets and chats on GroupMe and Kik, has already been shut down due to concerns with its inability to recognize when it was making offensive or racist statements. Of course, the bot wasn't *coded* to be racist, but it "learns" from those it interacts with. And naturally, given that this is the Internet, one of the first things online users taught Tay was how to be racist, and how to spout back ill-informed or inflammatory political opinions. [Update: Microsoft now says it's "making adjustments" to Tay in light of this problem.]

Holistic View of Adversarial Robustness



	Attack Category / Attacker's reach	Data	Model / Training Method	Inference
\checkmark	Poisoning Attack [learning]	Х	X*	
\checkmark	Backdoor Attack [learning]	Х		
\checkmark	Evasion Attack (Adversarial Example) [learning]		X*	X
Extraction Attack (Model Stealing, Membership Inference)				X
	Model Injection [AI governance]		Χ*	X

*No access to model internal information in the black-box attack setting

Pin-Yu Chen and Sijia Liu, Holistic Adversarial Robustness of Deep Learning Models, arxiv 2022

Robustness Challenges in AI Lifecycle



AI Lifecycle & robustness inspection



https://www.technologynetworks.com/tn/articles/cars-require-regular-inspection-why-should-ai-models-be-any-different-359405

Conceptual Pipeline of AI Model Inspector



Roadmap toward Holistic Adversarial Robustness



Research Highlights (1): *Finding failure modes*

- Practical white-box and black-box robustness testing
- Principled methods demonstrated on different data modalities and machine learning tasks in digital space and physical world
 - Images
 - Texts
 - Audio/Speech
 - Graphs
 - Tabular data
 - Reinforcement learning
 - □ CNN/RNN/LSTM/Transformer







Original Top-3 inferred captions:

- 1. A red stop sign sitting on the side of a road.
- 2. A stop sign on the corner of a street.
- 3. A red stop sign sitting on the side of a street.

Adversarial Top-3 captions:

- 1. A brown teddy bear laying on top of a bed.
- 2. A brown teddy bear sitting on top of a bed.
- 3. A large brown teddy bear laying on top of a bed.

Task: Fake-News Detection. Classifier: LSTM. Original label: 100% Fake. ADV label: 77% Real

Man Guy punctuates high-speed chase with stop at In-N-Out Burger drive-thru Print [Ed.—Well, that's Okay, that 's a new one.] A One man is in custody after leading police on a bizarre chase into the east Valley on Wednesday night. Phoenix police began has begun following the suspect in Phoenix and the pursuit continued into the east Valley, but it took a bizarre turn when the suspect stopped at an In-N-Out Burger restaurant's drive-thru drive-through near Priest and Ray Roads in Chandler. The suspect appeared to order food, but then drove away and got out of his pickup truck near Rock Wren Way and Ray Road. He then ran into a backyard ran to the backyard and tried to get into a house through the back door get in the home.





Inspecting AI/ML systems with Limited Knowledge: ZOO Attack

Prior to our work, all attacks either require white-box assumption or attack transfer



Pin-Yu Chen*, Huan Zhang, Yash Sharma, Jinfeng Yi, Cho-Jui Hsieh. ZOO: Zeroth Order Optimization based Black-box Attacks to Deep Neural Networks without Training Substitute Models. AI-Security Workshop 2017. Best Paper Award Finalist

Research Summary

- Practicality: minimal dependency on model and platform (no back-prop, only function calls) – adopted in IBM Watson products
- □Shift the focus of the research community from transfer attack to query-based attack
- Inspire many follow-up works on query-efficient and hard-label black-box attacks
- □New use case for gradient-free optimization
- Tutorial on "<u>Zeroth-Order Optimization: Theory</u> <u>and Applications to Deep Learning</u>" at CVPR'20 & KDD'19
- Applications: Contrastive explanations and endto-end molecule optimization using ML models

Steve is the tall guy with long hair who does not wear glasses



Amit Dhurandhar*, Pin-Yu Chen*, Ronny Luss, Chun-Chen Tu, Paishun Ting, Karthikeyan Shanmugam, and Payel Das, "Explanations based on the Missing: Towards Contrastive Explanations with Pertinent Negatives" NeurIPS 2018 Samuel Hoffman, Vijil Chenthamarakshan, Kahini Wadhawan, Pin-Yu Chen, and Payel Das, "Optimizing Molecules using Efficient Queries from Property Evaluations," Nature Machine Intelligence, 2021

Research Highlights (2) *Tradeoff between accuracy and adversarial robustness*

- Large-scale adversarial robustness benchmarks using 18 ImageNet models
- □ Standard Accuracy ≠ Adversarial Robustness
- Solely pursuing for highaccuracy AI model may harm deployment

Adversarial Robustness



Standard Accuracy

Research Summary

- Model performance benchmarking beyond standard accuracy
- From adversarial robustness to general robustness
 - $\circ~$ Common corruptions
 - $\circ~$ Distribution shifts
 - $\circ~$ Semantic changes
 - $\circ~$ Out-of-domain samples
- Tutorial on "Foundational Robustness of Foundation Models" at NeurIPS'22 with Sijia Liu and Sayak Paul
- Tutorial on "<u>Practical Adversarial Robustness in</u> <u>Deep Learning: Problems and Solutions</u>" at CVPR'21 with <u>Sayak Paul</u>
- □ Tutorial on "<u>Adversarial Robustness of Deep</u> <u>Learning Models</u>" at ECCV'20

Lecturer at MLSS 2021

Vision Transformers are Robust Learners

Sayak Paul* PyImageSearch s.paul@pyimagesearch.com

Purpose

Common

corruptions

Common

perturbations

Semantic shifts

Out-of-domain

distribution

Natural adversarial

examples

Background

dependence

Dataset

ImageNet-C [13]

ImageNet-P [13]

ImageNet-R [14]

ImageNet-O [9]

ImageNet-A [9]

ImageNet-9 [15]

Pin-Yu Chen* IBM Research pin-yu.chen@ibm.com





Research Highlights (3): *Practical adversarial threat detection and mitigation*

- Plug-and-play techniques to inspect and repair potential risks in trained neural networks
- Data efficiency: threat detection and mitigation using limited data
- Applicable to training-time and inference-time threats

Backdoor Attack



C^t

Chulin Xie, Keli Huang, Pin-Yu Chen, and Bo Li. DBA: Distributed Backdoor Attacks against Federated Learning. ICLR 2020 Wang et al. Neural Cleanse: Identifying and Mitigating Backdoor Attacks in Neural Networks. IEEE Security and Privacy, 2019

Research Summary

□ I have an amazing ImageNet model which achieves 95% top-1 accuracy, and I make it publicly available by releasing the network architecture and trained model weights. <u>Care to use it for your task</u>?

- Tempting ... but how do I know your model does not have any backdoor?
- ✓ Sanitize the model before using it (aka wear mask before you go out)





Data-limited and Data-free Trojan net detection





Ren Wang, Gaoyuan Zhang, Sijia Liu, Pin-Yu Chen, Jinjun Xiong, and Meng Wang. Practical Detection of Trojan Neural Networks: Data-Limited and Data-Free. *ECCV 2020*

Pu Zhao, Pin-Yu Chen, Payel Das, Karthikeyan Natesan Ramamurthy, and Xue Lin. Bridging Mode Connectivity in Loss Landscapes and Adversarial Robustness. ICLR 2020 Research Highlights (4): *Provable robustness assessment and quantification*

- Avoid false sense of robustness from improper empirical measures
- Certified robustness: attackproof risk quantification and verification for neural networks
- Legitimate answer to "Will my model become more robust if I do/use X?"



Research Summary

CLEVER score: First attack-independent and model-agnostic robustness score

Robustness certification tools

- $\,\circ\,$ Input perturbations
- Semantic perturbations
- \odot Training of verifiable neural networks
- $\circ\,$ Certified defenses
- $\,\circ\,$ Support for different network architectures

□AI regulation and standardization





Start >

ation on Evaluating the Robustness of Neu

http://bigcheck.mybluemix.net

Verification: lower bounds on robustness

Evaluating the Robustness of Neural Networks: An Extreme Value Theory Approach, Tsui-Wei Weng*, Huan Zhang*, Pin-Yu Chen, Jinfeng Yi, Dong Su, Yupeng Guo, Cho-Jui Hsieh, and Luca Daniel, ICLR 2018 Efficient Neural Network Robustness Certification with General Activation Functions, Huan Zhang*, Tsui-Wei Weng*, Pin-Yu Chen, Cho-Jui Hsieh and Luca Daniel, NeurIPS 2018 Research Highlights (5): (Hybrid) Quantum ML: Better Robustness, Privacy, and Generalization?

Variational Quantum Circuits for Deep **Reinforcement Learning**

SAMUEL YEN-CHI CHEN1, CHAO-HAN HUCK YANG2, JUN QI2, PIN-YU CHEN3, XIAOLI MA2, (Fellow, IEEE), HSI-SHENG GOAN^{1,4},





Hybrid Quantum ML model: quantum circuits for feature extraction, followed by classical neural networks

□ Parameter efficiency; Generalization Guarantee

Better privacy and robustness?

DECENTRALIZING FEATURE EXTRACTION WITH QUANTUM CONVOLUTIONAL NEURAL NETWORK FOR AUTOMATIC SPEECH RECOGNITION

Chao-Han Huck Yang¹ Jun Qi^1 Sabato Marco Siniscalchi^{1,4,5}

Samuel Yen-Chi Chen² Pin-Yu Chen³ Xiaoli Ma¹

Chin-Hui Lee¹

WHEN BERT MEETS QUANTUM TEMPORAL CONVOLUTION LEARNING FOR TEXT CLASSIFICATION IN HETEROGENEOUS COMPUTING

Jun Qi^1 Chao-Han Huck Yang¹ Samuel Yen-Chi Chen² Pin-Yu Chen⁴ Yu Tsa o^3



Quantum Neural Networks for Speech and Natural Language Processing (QuantumNN) Tutorial

Montreal, 21st (Sat.) August (Virtual Room Auditorium Red), IJCAI, 2021

THEORETICAL ERROR PERFORMANCE ANALYSIS FOR VARIATIONAL QUANTUM CIRCUIT BASED FUNCTIONAL REGRESSION

A PREPRINT

Jun Qi Georgia Institute of Technology jqi41@gatech.edu Chao-Han Huck Yang Georgia Institute of Technology huckiyang@gatech.edu Pin-Yu Chen IBM Research pin-yu.chen@ibm.com

Min-Hsiu Hsieh Hon Hai (Foxconn) Research Institute minhsiuh@gmail.com









QTN-VQC: AN END-TO-END LEARNING FRAMEWORK FOR QUANTUM NEURAL NETWORKS

A PREPRINT

Jun Qi Georgia Institute of Technology jqi41@gatech.edu Chao-Han Huck Yang Georgia Institute of Technology huckiyang@gatech.edu Pin-Yu Chen IBM Research pin-yu.chen@ibm.com

(a) A dense layer for dimension reduction

Our portfolio in adversarial robustness research



https://www.ucc.ie/en/cirtl/newsandevents/cirtl-seminar-the-assessment-arms-race-and-its-fallout-the-case-for-slow-scholarship-may-14th.html

Adversarial Robustness Toolbox (ART)

Adversarial Robustness Toolbox (ART)

External: https://github.com/IBM/adversarial-robustness-toolbox

- Python library, 7K lines of code
- State-of-the-art attacks, defences and robustness metrics

	from keras.datasets import mnist from keras.models import load_model	
Load ART	<pre>from art.attacks import CarliniL2Attack from art.classifier import KerasClassifier from art.metrics import loss_sensitivity</pre>	
	<pre># Load data (_, _), (x_test, y_test) = mnist.load_data()</pre>	
Load classifier model (Keras, →→ TF, PyTorch etc)	<pre># Load model and build classifier model = load_model('my_favorite_keras_model.h5 classifier = KerasClassifier((0, 1), model)</pre>	
Perform attack	<pre># Perform attack attack = CarliniL2Attack(classifier) adv_x_test = attack.generate(x_test)</pre>	
Evaluate	<pre># Compute metrics on model robustness print(loss_sensitivity(classifier, x_test))</pre>	

Robustness service based on ART on the roadmap for general availability under IBM AI OpenScale offering

Evasion attacks	Evasion defenses	Poisoning detection	Robustnes
• FGSM	Feature squeezing	Detection based on	• CLEVER
• JSMA	Spatial smoothing	clustering activations	• Empirica



IBM donates "Trusted AI" projects to Linux Foundation AI

As real-world AI deployments increase, IBM says the contributions can help ensure they're fair, secure and trustworthy

Adversarial Robustness 360		AI Fairness 360		AI Explainability 360
५ (ART)		└ (AIF360)		Ს (AIX360)
github.com/IBM/ adversarial-robustness-		github.com/IBM/AIF360		github.com/IBM/AIX360
art-demo.mybluemix.net		aif360.mybluemix.net		aix360.mybluemix.net
Ţ				

Awards and Prizes



Take-aways

- ✓ Adversarial robustness is a new AI standard toward trustworthy ML
- Robustness does not come for free: failure modes exist in digital space, physical world, and different domains during AI lifecycle
- □High accuracy ≠ Good robustness
- Arms race in adversarial ML: adversary-aware AI v.s. AI for adversary
- ✓ AI model inspector for holistic robustness and beyond
- Practical techniques and tools for identifying failure modes
- □Plug-and-play model error detection and mitigation
- □ Robustness quantification and verification
- Adversarial ML for good: model reprogramming
- ✓ My long-term goal: Make the robustness inspection pipeline for AI models as reliable, standard, and easy, as car models



Online Resources for Adversarial Robustness

- J. Z. Kolter and A. Madry: Adversarial Robustness Theory and Practice (NeurIPS 2018 Tutorial)
- Pin-Yu Chen: Adversarial Robustness of Deep Learning Models (ECCV 2020 Tutorial)
- Pin-Yu Chen and Sijia Liu: <u>Zeroth Order Optimization: Theory and Applications to Deep Learning</u> (CVPR 2020 Tutorial)
- Pin-Yu Chen and Sayak Paul: <u>Practical Adversarial Robustness in Deep Learning: Problems and Solutions (CVPR</u> 2021 Tutorial)
- Pin-Yu Chen: <u>Holistic Adversarial Robustness for Deep Learning</u> (MLSS 2021 Tutorial)
- Pin-Yu Chen: Adversarial Machine Learning for Good (AAAI 2022 Tutorial)
- Pin-Yu Chen, Sijia Liu, and Sayak Puak: Foundational Robustness of Foundation Models (NeuIPS 2022 Tutorial)



Adversarial Robustness Toolbox (ART v0.10.0)



Foolbox

Sample Surveys for Adversarial Robustness

Wild Patterns: Ten Years After Adversarial Machine Lea	the Rise of rning	On Evaluating Adversarial Robustness	
Battista Biggio ^{a,b,*} , Fabio Roli ^a	,b		
^a Department of Electrical and Electronic Engineering, Univ ^b Pluribus One, Cagliari, Italy	versity of Cagliari, Italy	Nicholas Carlini ¹ , Anish Athalye ² , Nicolas Papernot ¹ , Wieland Brendel ³ , Jonas Rauber ³ , Dimitris Tsipras ² , Ian Goodfellow ¹ , Aleksander Mądry ² , Alexey Kurakin ¹ *	
	Holistic Adversarial Robustness of Deep Learning Models	1 Google Brain 2 MIT 3 University of Tübingen	
	Pin-Yu Chen ^{1,3*} , Sijia Liu ^{2,3} ¹ IBM Research, ² Michigan State University, ³ MIT-IBM Watson AI La pin-yu.chen@ibm.com and liusiji5@msu.edu	b On Adaptive Attacks to Adversarial Example Defenses	
The Robustness of Deep Networks		Florian Tramèr [*] Nicholas Carlini [*] Wieland Brendel [*] Stanford University Google Brain University of Tübingen Aleksander Mądry	
A geometrical perspective		MIT	
Alhussein Fawzi, Seyed-Mohsen Moosavi-Dezfooli, and Pascal Frossard		Adversarial Robustness	
Adversarial Learning Targeting Comprehensive Review of Def Publisher: IEEE Cite This	g Deep Neural Network Classification: A fenses Against Attacks	Adversarial Robustness for Machine Learning	
3 Author(s) David J. Miller (10); Zhen Xiang (10); George Kesidis View All Authors		Paperback ISBN: 9780128240205	

Find limitations

- Adversarial examples
- Out-of-distribution Generalization

Improve Robustness

- Threat/risk mitigation and evaluation
- Robust training and transfer

Learning with an Adversary

(Adversarial ML)

Create synergies

- Generative adversarial nets
- Policy learning (hide and seek)

Boost machine learning

- Data augmentation
- Model reprogramming
- Al governance