

Adversarial XAI methods in Cybersecurity

Aditya Kuppa and Nhien-An Le-Khac
UCD

Why *

Question 1. The ability to construct a coherent and complete “story” with the facts of a situation is the most important task when making a decision or recommendation.

Agree 93%

Disagree 7%

Question 2. As a forecasting/recommendation task becomes more complex and difficult, I tend to rely more on judgment and less on formal, quantitative analysis.

Agree 64%

Disagree 36%

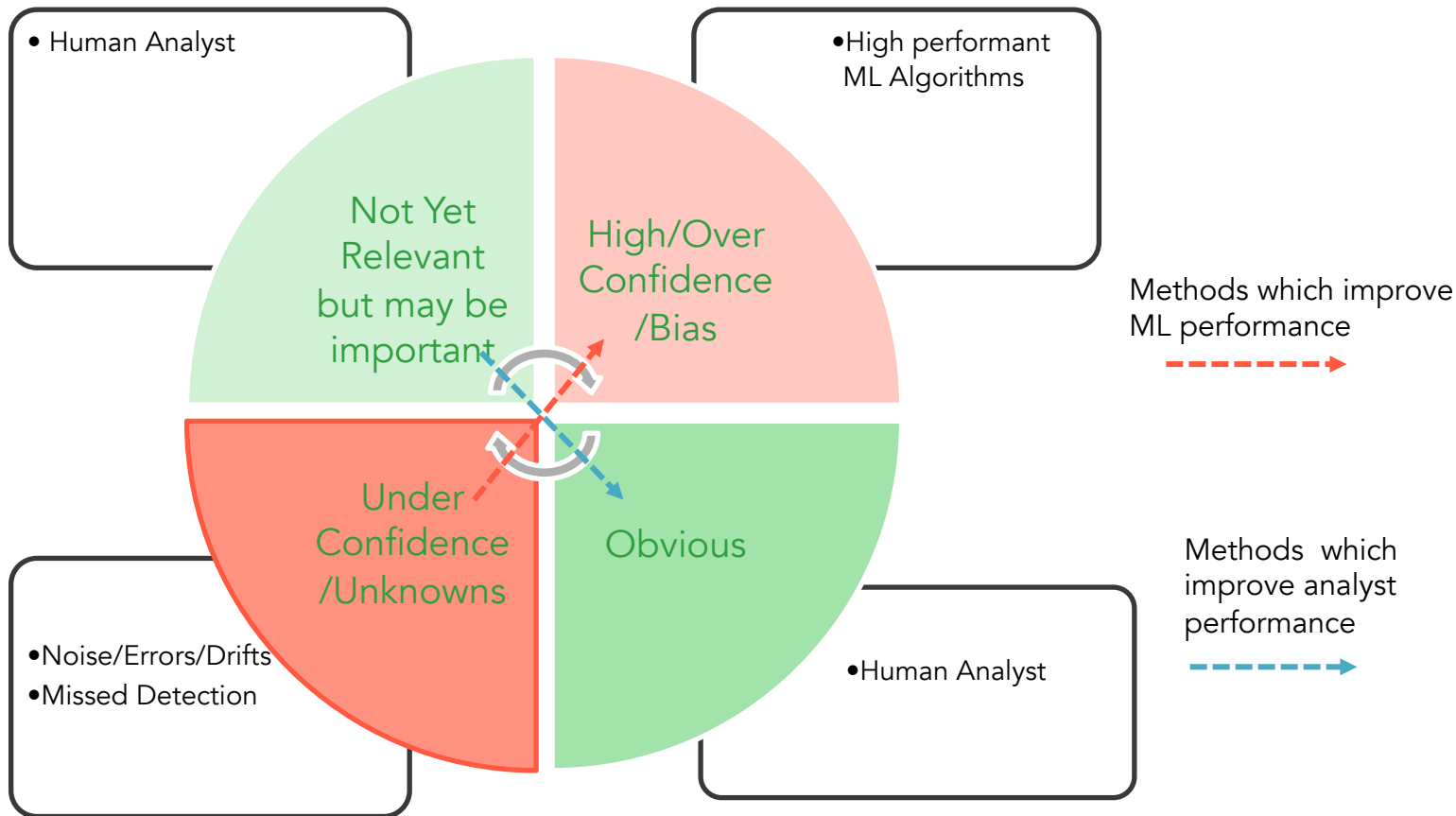
Question 8. As I become more uncertain about my ability to predict outcomes, I give greater weight to negative information about alternatives.

Agree 86%

Disagree 14%

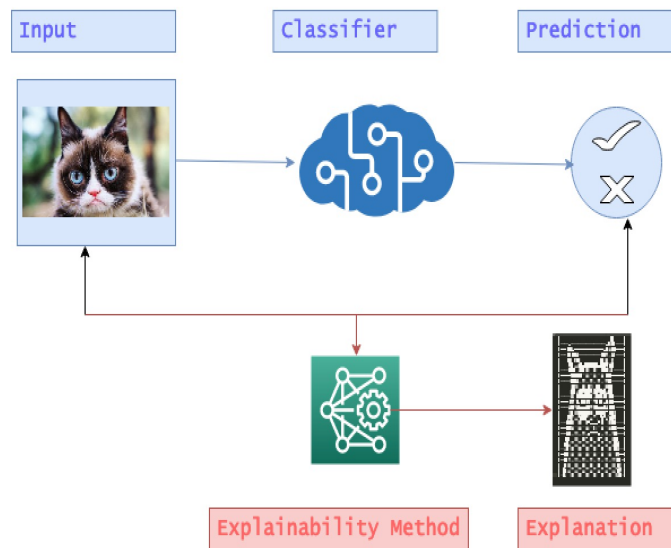
* Behaviour Response of CFA Charter Holders (Olsen, Robert A. "Professional investors as naturalistic decision makers: Evidence and market implications." *The Journal of Psychology and Financial Markets* 3.3 (2002): 161-167.)

Why ?*



Explanations

- The field of explanations of intelligent systems was active in the 1970s for expert systems; 1980's for neural networks; and then to recommendation systems in the 2000s.
- Explainability methods
 - Post-hoc/During/Pre-hoc
 - Scope - Local, Global
 - Dependency - Model, Data and Domain
- Interpretability methods coupled with the human in the loop improves the trust and security in the decision making process of ML systems.



Explanations in Security Domain

- Problems in Security Domain
 - Imbalanced DataSets
 - Attribution of threat and Context is important and hard to infer.
 - Threats are always evolving and there is a need to improve robustness of underlying systems.



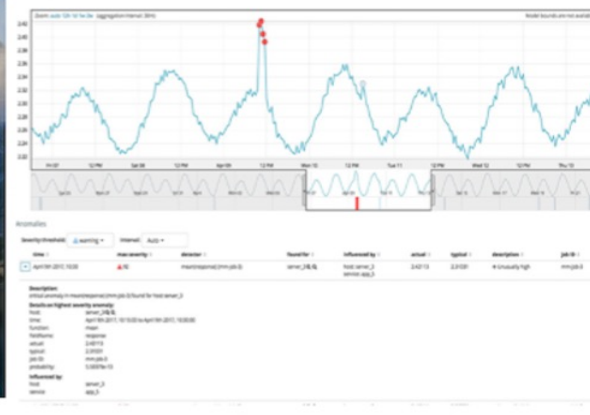
Login Failed



The screenshot shows a security dashboard with a "malicious" badge in the top right. The main content is titled "Behavioral Indicators" and lists several indicators under "Malicious Indicators" and "Suspicious Indicators".

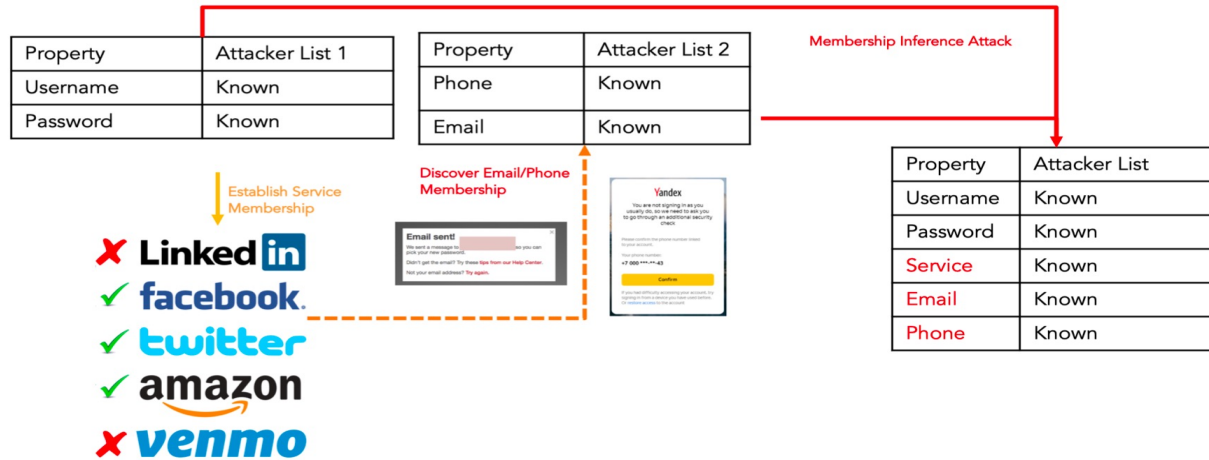
- Malicious Indicators:**
 - Exploit/Shellcode:** Requested execution of the SHELX loader (other part of generating kernel exploits). Details: "Requested execution of the SHELX loader (other part of generating kernel exploits). Source: Based on API Call."
 - Search the SYSTEM ADMIN PAGE (LOCAL SYSTEM EXECUTION):** Details: "Attempt to search for the SYSTEM ADMIN PAGE (LOCAL SYSTEM EXECUTION). Source: Based on Network Target."
- Suspicious Indicators:**
 - Exploit/Shellcode:** Possible heap spraying attempt detected.
 - Installation/Persistence:** Changes memory access rights in a remote process to writebackcode. Details: "Writes a file to the start menu. Source: 'The API call is for the SYSTEMADMINPAGE (LOCAL SYSTEM EXECUTION)'. Source: Based on API Call."
 - Writes data to a remote process.**
 - Network Phishing:** Found potential IP address in binary memory.

The screenshot shows a Yandex login security check screen. It features the Yandex logo and the text: "You are not signing in as you usually do, so we need to ask you to go through an additional security check." Below this, it asks to "Please confirm the phone number linked to your account." and displays "Your phone number: +7 000 ***-**-43". A yellow "Confirm" button is visible. At the bottom, it provides instructions: "If you had difficulty accessing your account, try signing in from a device you have used before. Or restore access to the account."



Goals and Motivation

- Security analysis of XAI methods
 - How can an attacker, given only outputs of explanation method and model predictions, can conduct powerful black-box model extraction, membership inference attacks?
 - How explanation outputs facilitate the generation of adversarial samples and poison/backdoor samples to evade the underlying classifier
- Motivating Example – Credential Stuffing (Membership Inference attack)



Threat Model Assumptions

CHARACTERISTIC	TYPE	MEA [26]	MIA [68]	PA [67]	AE [45]
<i>Knowledge</i>	TRAINING DISTRIBUTION	X	X	✓	X
	FEATURE SET	✓	✓	✓	✓
	FEATURE EXTRACTOR	✓	✓	✓	✓
	FEATURE TRANSFORMERS	✓	✓	✓	✓
	INFERENCE API	✓	✓	✓	✓
	EXPLANATIONS INTERFACE/METHOD	✓	✓	✓	✓
	CONFIDENCE INTERVALS	✓	✓	✓	✓
<i>Goal/Intent</i>	COMPROMISING INTEGRITY (EVASION)	X	X	✓	✓
	COMPROMISING PRIVACY	✓	✓	X	X
<i>Capability</i>	MANIPULATE TRAINING DATA	X	X	✓	X
	MANIPULATE TEST DATA	X	✓	X	✓
<i>Strategy</i>	TRAIN A SURROGATE MODEL FOR PARAMETER EXTRACTION	X	✓	X	X
	TRAIN A SURROGATE MODEL FOR QUERY REDUCTION	✓	X	X	✓
	SATISFY DOMAIN CONSTRAINTS	X	✓	✓	✓
<i>Frequency</i>	ITERATIVE	✓	✓	✓	✓
<i>Perturbation Scope</i>	INSTANCE SPECIFIC	✓	✓	✓	✓
<i>Perturbation Constraints</i>	OPTIMISATION	✓	X	✓	X
	DOMAIN	✓	✓	✓	✓

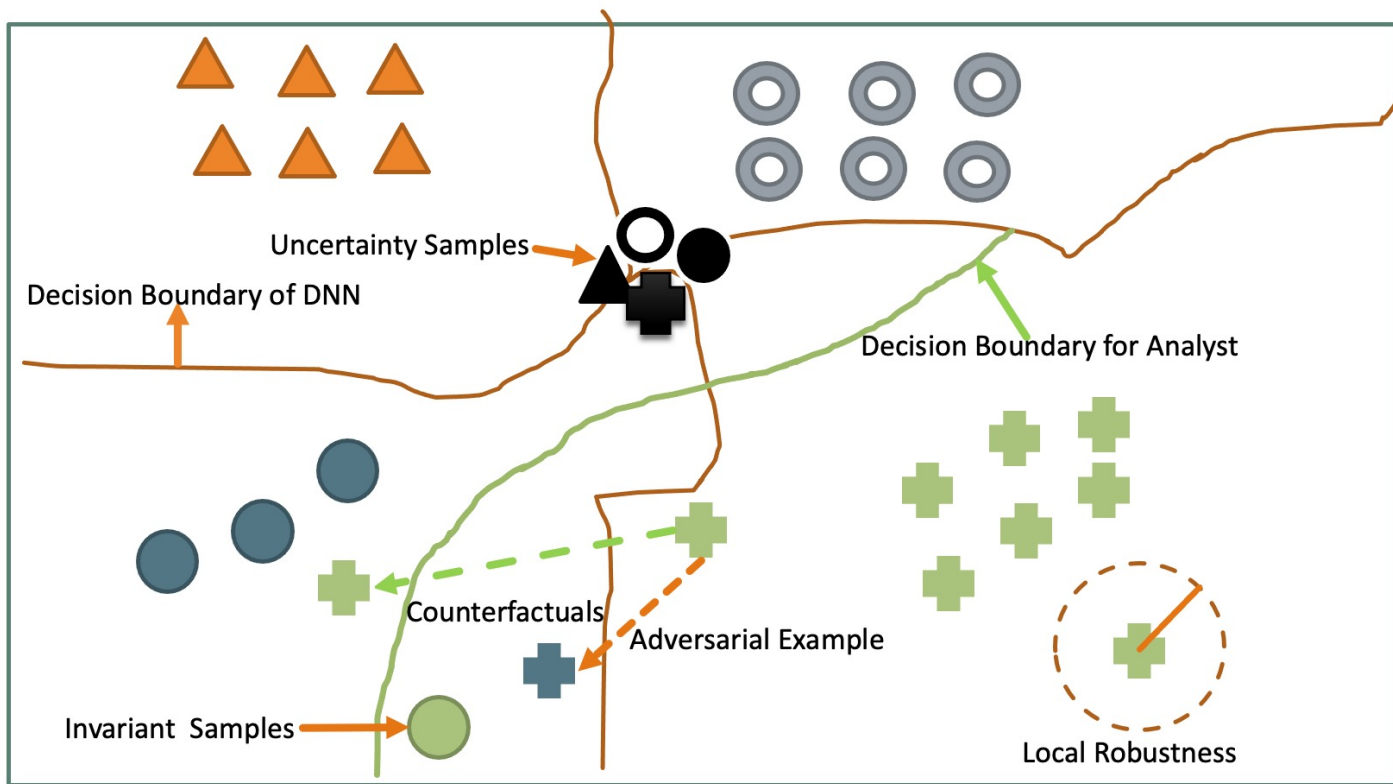
Counterfactuals (why/why-not)

- Counterfactual data instances of the input have
 - Similar feature values as input
 - Different model predictions from that of input
 - Lay closer to the decision boundary of an input class

$$x_{cf} = \underset{x_{cf1}, \dots, x_{cfk}}{\operatorname{argmin}} \mathcal{L}(\mathcal{T}(x_{cf}), \mathcal{T}(x)) + \operatorname{Dist}(x_{cf} - x)$$

Method	\mathcal{L}	Dist	CF per x	Optimisation Method
Latent CF [38]	Latent Vector Loss	ℓ_1	1	Gradient Descent
DICE [39]	Hinge-loss	ℓ_1 and Median Absolute Deviation(MAD)	k	Gradient Descent
Permute Attack [40]	-	ℓ_2	1	Genetic Algorithm

Class Decision Boundary (AE vs CF)



Attack Method

- Given a black-box access to a target model T prediction interface $T(x) = y$, x_{cf} counterfactual, $E(x) = x_{cf}$ explanation interface, D_{aux} auxiliary dataset and S a surrogate model
 - Attacker aims to compromise the confidentiality and integrity of the underlying ML system
- Explanation-based Poisoning Attack
 - Identify and Perturb robust features, which are consistently same across their counterfactual class.
- Explanation-based Adversarial Sample Generation
 - Adapt Counterfactual method which works in feature space to sample space.
- Explanation-based Membership Inference Attack
 - 1-Class Nearest neighbor classifier for each class is trained on counterfactuals to establish membership
- Explanation-based model extraction
 - Knowledge distillation technique to transfer knowledge from the target model to the surrogate model

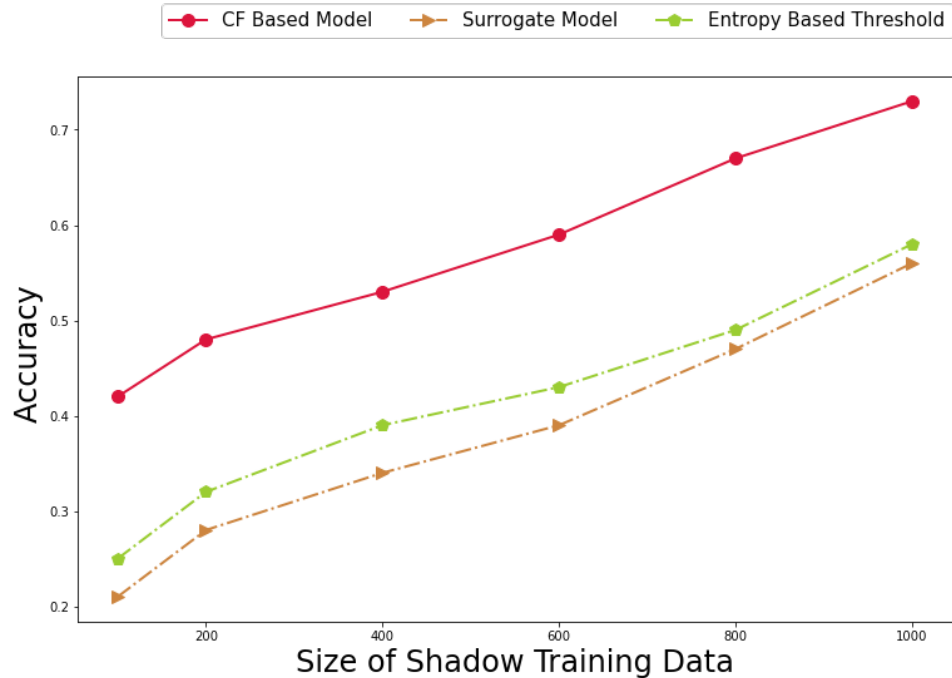
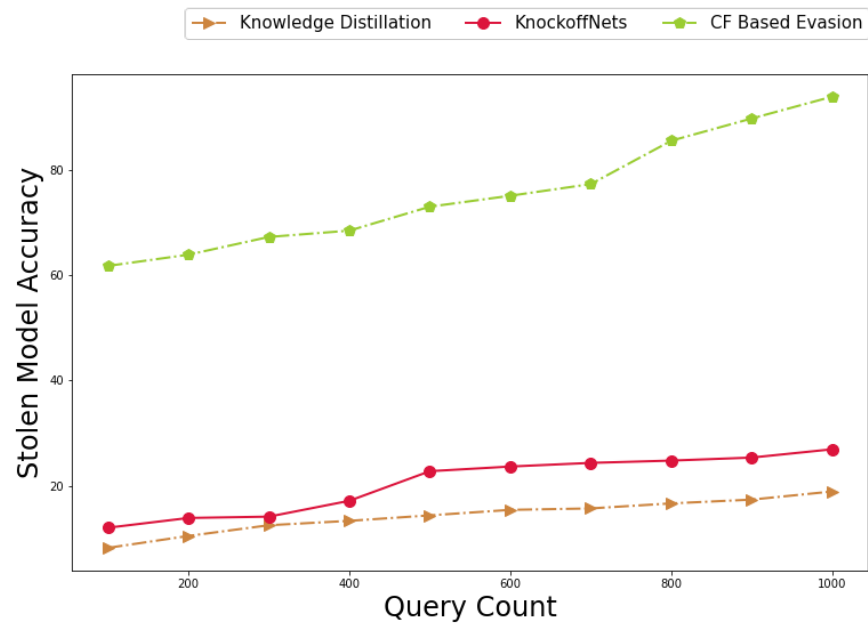
DataSets

- MEA
 - CICIDS17 Network Traffic dataset which contains a wide range of attack types like SSH brute force, Botnet, DoS, DDoS, web, and infiltration
- AE and Poisoning Attack
 - 30120 malware from virus share and for benign samples we scrapped 20334 clean files from free ware sites
- MIA
 - Leaked Password Dataset -- The dataset consists of 1.4 billion email password pairs with 1.1 unique emails and 463 million unique passwords. This dataset is aggregated password leaks from different incidents.

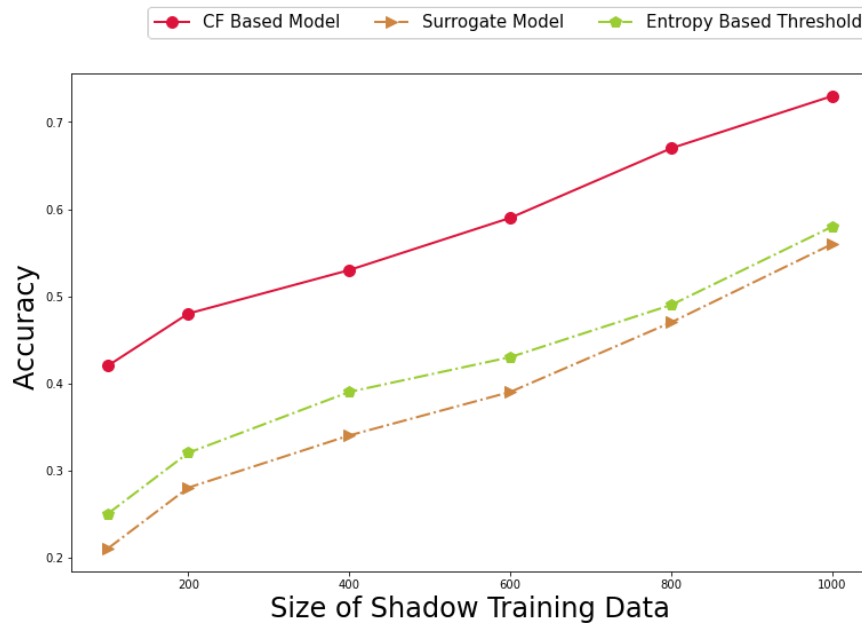
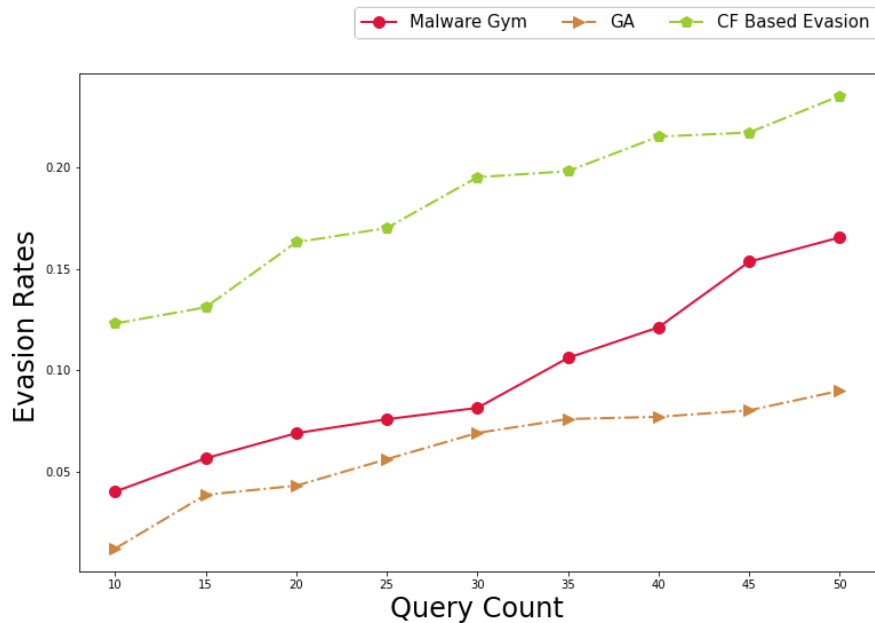
Results

ATTACK TYPE	D_{aux}	EXPLANATION METHOD	\mathcal{T}	ORIGINAL ACCURACY	EVASION ACCURACY
ADVERSARIAL ATTACK	MALWARE	PERMUTE	AV1	93.5%	65.23%
			AV2	94.7%	41.89%
POISONING ATTACK	MALWARE	PERMUTE	GBM	POISONING PERCENT	ACCURACY DROP
				0.5%	62.4%
				1%	76.23%
			NN	2%	87.24%
				0.5%	30.9%
				1%	50.89%
			2%	65.31%	
			3%	79.48%	
MEMBERSHIP INFERENCE	LEAKED PASSWORDS	LATENT-CF	AUTOENCODER	METHOD MODEL	ACCURACY/QUERIES
				ENTROPY	49.46/1000
				CF	54.17/1000
					73.17 /1000
MODEL EXTRACTION	CICIDS	DICE	AUTOENCODER	MODEL	ACCURACY
				\mathcal{T}	98.02
				KN	78.91
				KD	53.89
				CF	93.54

MIA and MEA Comparison



AE and Poisoning Attack Comparison



Defense Discussion

- CF methods share a large set of similarities with adversarial examples concepts
 - Adapting methods from adversarial defense literature
- Noise based Defense Intuition
 - Defender has no control over the attacker's *full* training data but only a portion of it.
 - ML model aims to learn the mapping function from the feature space to the label space from the training samples.
- Defender can transform the counterfactual samples so the learned model (surrogate) has a strong correlation between the labels and noise of the feature space instead of only features.
- Adding noise to samples
 - Individual CF sample and All CF samples of same class

$$T_s = \begin{cases} \text{None} & \text{No transformation} \\ \phi, & \text{Random noise } [-1, 1] \\ \delta_{x_i}, & \text{Adv. noise } [-\epsilon, \epsilon] \\ \delta_{y_i} & \text{Adv. noise } [-\epsilon, \epsilon] \end{cases}$$

T_s	Accuracy
None	95.6
ϕ	87.4
δ_{x_i}	37.4
δ_{y_i}	28.22
$\delta_{x_i} + \phi$	32.22
$\delta_{y_i} + \phi$	18.22

Limitations and Future work

- MEA
 - Methods which optimize on multiple properties of CF's improve the stolen model accuracy
 - We only tested non-differential models
- MIA
 - Methods which do not employ latent space to search for CF need large number of queries.
 - Learning password rules and investigate how CF attack can speed up the password cracking methods
- AE and Poisoning
 - The functionality preserving transformation functions applied on the binary are biased towards static features.
 - Our results may not be valid when AV engines use both static and dynamic analysis to make a decision.
 - CF methods can help attackers to find quicker ways to find adversarial/poisoned samples, instead of solving a hard-to converge black-box optimization problem in input space.

