

## " Machine Learning Based Cyber Attacks Targeting on Controlled Information: A Survey "

2022.10.03

Professor: 박종혁

**Presented by:** 

YOTXAY SANGTHONG

Seoul National University of Science and Technology, Seoul, South Korea.

SeoulTech UCS Lab

1

#### **Table Contents**

Abstract

- I. Introduction
- II. Attack Methodology
  - 2.1 Reconnaissance
  - 2.2 Data Collection
  - 2.3 Feature Engineering
  - 2.4 Attacking the Objective
  - 2.5 Evaluation
- III. ML-based Stealing Attacks and Protections
  - 3.1 Controlled User Activities Information
  - 3.2 Controlled ML Model Related Information

#### Abstract

- This survey presented the recent advances in this new type of attack and corresponding countermeasures.
- The ML-based stealing attack is reviewed in perspectives of three categories of targeted controlled information.
- Including controlled user activities, controlled ML model-related information, and controlled authentication information.

## **I. Introduction**

- Driven by the needs to protect the enormous value within data and the evolution of the emerging data mining techniques, information leakage becomes a growing concern for governments, organizations and individuals.
- This survey introduced the stealing attack in the cyber security area. The information leakage can be defined as the violation of confidentiality of methods/mechanisms /framework which stores information or has access to information.
- For example, authors [26] extracted user's foreground app running in Android in order to exploit it for the phishing attack, while the user activity information was protected by a nonpublic system level permission.

## I. Introduction (Cont.)



Fig. 1. Introduced Stealing Controlled Information Attack Categories.

## I. Introduction (Cont.)

The contributions:

- This paper introduced the ML-based stealing attack, which aims at stealing the controlled/protected information and leads to huge economic loss. ML algorithms are applied in the attack to increase the success rate in various aspects.
- The classification of the ML-based stealing attacks is built based on the targeted controlled information preferentially. Based on this classification, the vulnerabilities in various systems and corresponding attacks are sorted out and revealed.
- The authors surveyed the advances of the ML-based stealing attacks between 2014 and 2019. A methodology applied for the ML-based stealing attack against the controlled information is generalized to five phases reconnaissance, data collection, feature engineering, attacking the objective, and evaluation.
- They discussed the challenges of attacks stealing controlled information and forecast their future directions in terms of how they might affect our digital society.

## II. Attack Methodology



Fig. 2. ML-based stealing attack methodology

#### 2.1 Reconnaissance

- Reconnaissance refers to a preliminary inspection of the stealing attack. The two aims of this inspection include defining adversaries' targets and analyzing the accessible data in order to facilitate the forthcoming attacks.
- The target of adversaries in the published literature is usually the confidential information controlled by systems and online services.
- The attacker needs to exhaustively search all possible entry points of the targeted system, reachable data paths, and readable data.
- When the attacker aims at user's activities, the triggered hardware devices and their corresponding logged information will be investigated.
- For example, the attacker always searches and explores the readable system files, such as interrupt timing data and network resources.

## 2.2 Data Collection

- Active collection refers to the attacker actively interacts with the targeted system for data collection.
- Specifically, an attacker designs some initial queries to interact with the system and subsequently collects the data. The goal of the attacker guides the design of malicious interactions, referring to the analysis results from the reconnaissance phases.
- For example, if an attacker intends to identify which app is launched in a user's mobile, some system files like *proc f s* recording app launching activities.

#### **2.3 Feature Engineering**

- After the datasets are prepared, feature engineering is the subsequent essential phase to generate representative vectors of the data to empower the ML model. The two key points in feature engineering for ML-based attacks consist of dataset cleaning and extracting features.
- An obstacle of feature engineering is cleaning the noises and irrelevant information in the raw data. In general, deduplication and interpolation can be used to reduce the noise from accessible resource.
- To reduce the noise, a Fast Fourier Transform (FFT) filter and an Inverse FFT (IFFT) filter are applied.

#### 2.4 Attacking the Objective

- The ML-based stealing attack into two attack modes as illustrated in Fig. 3. The five actions correspond to the first three phases within the MLBSA methodology.
- As stated in the data collection phase, the inputs and their query results are collected as the required accessible dataset, which reveals the target information.
- Based on the target information, the ground truth of the dataset is set up in this phase. With proper feature engineering methods, the training dataset is prepared to attack the objective.
- But the subsequent actions to steal the controlled information using machine learning differ between two attack modes.

2.4 Attacking the Objective (Cont.)



SeoulTech UCS Lab

Fig. 3. The ML-based stealing attack into two attack modes

## 2.4 Attacking the Objective (Cont.)

- For the first attack mode as shown in Fig. 3a, this attack mode is applied in the ML-based stealing attack against the user activity information, the authentication information, and training set information.
- Regarding the testing dataset is collected from a victim's system/service, the testing samples are not labeled while querying the attack model.
- Since the attack model is **built to infer the controlled information from these accessible data**, the output of the model is the targeted controlled information.
- The second attack mode illustrated in Fig. 3b, is mostly applied against the ML modelrelated information. In a black-box setup, stealing the ML model attack aims at calculating the detailed expression of the model's objective function.
- Reconstructing the original model is essentially a reconstruction attack. Using the equation-solving and path-finding methods, the inputs and their query outputs for solving the specific objective function expression is interpreted as the training set.
- Therefore, this attack can be simplify regarded as an ML-based attack. Additionally, based on the attackers' inputs and the query outputs, the training set is synthesized and used to build a substitute model for reconstruction.

#### **2.5 Evaluation**

- Evaluation metrics differ between two attack modes.
- For the first attack mode, the attack model is the attacker's weapon.
- Most metrics commonly used to measure the effectiveness include accuracy, precision, recall, FPR, FNR, and F-measure.
- Accuracy: It is also known as success rate and inference accuracy. Accuracy means the number of correctly inferred samples to the total number of predicted samples. Accuracy is a generic metric evaluating the attack model's effectiveness.

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$$

• **Precision:** It is regarded as one of the standard metrics for attack accuracy. Precision illustrates the percentage of samples correctly predicted as controlled class *A* among all samples classified as *A*.

$$Precision = \frac{TP}{TP + FP}$$

• **Recall:** It is regarded as another standard metric for attack accuracy. Recall is also called sensitivity or True Positive Rate (TPR). It is the probability of the amount of class A correctly predicted as class *A*. Similar to precision, recall also reveals the model's correctness on a specific class. These two metrics are almost always applied together.

$$Recall = \frac{TP}{TP+FN}$$

#### 2.5 Evaluation (Cont.)

• **F-measure:** This metric or **F1-score** is the harmonic mean of recall and precision. F-measure provides a comprehensive analysis of precision and recall .

$$F - measure = \frac{2 \times Recall \times Precision}{Recall + Precision}$$

• **False positive rate (FPR):** This metric denotes the proportion of class *B* samples mistakenly categorized as class *A* sampled.

FPR assesses the model's misclassified samples.  $FPR = \frac{FP}{TN+FP}$ 

• **False negative rate (FNR):** This metric stands for the ratio between class *A* samples mistakenly categorized as class *B* samples. Similar to FPR, FNR assesses the model's misclassified samples from another aspect. FPR and FNR are almost always applied together to measure the model's error rate.  $FNR = \frac{FN}{TP+FN}$ 

• **Battery consumption:** It is also known as **power consumption**. Battery consumption refers to the target mobile's battery while the target system is a mobile system, which indicates the efficiency of the attack model. For the second attack.

## 2.5 Evaluation (Cont.)

• For the second attack mode, ML-based attacks of stealing the ML model are assessed with other metrics.

The applied evaluation metrics are defined and listed below:

• **Test error** is the average error based the same test set (*D*) testing at learned model and targeted model [125]. A low test error means *f*<sup>^</sup>matches *f* well.

$$Error_{test}(f, \hat{f}) = \frac{\sum_{x \in D} diff(f(x), \hat{f}(x))}{|D|}$$

• **Uniform error** is an estimation of the portion of full feature space that the learned model is different from the targeted one, when the testing set (*U*) are selected uniformly [125].

$$Error_{uniform}(f, \hat{f}) = \frac{\sum_{x \in U} diff(f(x), \hat{f}(x))}{|U|}$$

• **Extraction accuracy** indicates the performance of model extraction attack based on the test error and the uniform error [125].

$$Accuracy_{extraction} = 1 - Error_{test}(f, \hat{f}) = 1 - Error_{uniform}(f, \hat{f})$$

#### 2.5 Evaluation (Cont.)

- **Relative estimation error (EE)** measures the effectiveness of model extraction attack using its learned hyperparameters ( $^{\lambda}$ ) contrasting to the original hyperparameters ( $^{\lambda}$ ) [131].  $Error_{EE} = \frac{|\hat{\lambda} - \lambda|}{\lambda}$
- **Relative mean square error (MSE)** measures how well the model extraction attack reconstructs the regression models via comparing the mean square error after learning hyperparameters using cross-validation techniques [131].

$$Error_{MSE} = \frac{|MSE_{\hat{\lambda}} - MSE_{\lambda}|}{MSE_{\lambda}}$$

• **Relative accuracy error (AccE)** measures how well the model extraction attack reconstructs the classification models via comparing accuracy error after learning hyperparameters using cross-validation techniques [131].

$$Error_{AccE} = \frac{|AccE_{\hat{\lambda}} - AccE_{\lambda}|}{AccE_{\lambda}}$$

## 2.5 Evaluation (Cont.)

Table 2. Outline of Reviewed Papers (info: information)

Reference	Year	Targeted Info	Accessible Data	Goals
[96]	2016	Unlock pattern;	Hardware internut data	Unlock pattern & foreground app inference attacks via analyzing
[20]	2010	Foreground app	Hardware Interrupt data	interrupt time collected from interrupt log file.
[110]	2018	Visited websites;	Interrupt data; Network	Search and attack the kernel records leaking user's specific events
[117]	2010	Foreground app	&Memory process record	(i.e. app starts, website launch, keyboard gesture).
[151]	2018	Visited websites;	Memory data; Network	Several side-channel inference attack on iOS mobile device
[131]	2010	Foreground app; Map	source; File system data	Several side challier interence attack on 105 mobile device.
[136]	2015	Visited websites;	Kernel data-structure	Protect by injecting noise into the value of kernel data
[150]	2015	Input keystrokes	fields	structure values to secure procfs.
[50]	2016	Manufacturing	Acoustic sensor data;	An attack capture acoustic & magnetic sensor data to steal a
[20]	2010	activities	Magnetic sensor data	manufacturing process specification or a design.
[117]	2017	User activities info	Sensor data	Contextual model detect malicious behavior of sensors like leaking.
[125]	2016	Parameters of	Input features &	Model extraction attacks leverage confidence info with
[123]	2010	an ML model	Query outputs	predictions against MLaaS APIs in black-box setting.
[101]	2017	Internal info of	Input features &	Build a local model to substitute the target model and use
		an ML model	Query outputs	it craft adversarial examples in black-box setting.
[131]	2018	Hyperparameters	Input features &	Hyperparameters stealing attack via observing minima objective
[151]	2010	of an ML model	Query outputs	function against MLaaS in black-box setting.
[98]	2018	Hyperparameters	Input features &	Build a metamodel to predict hyperparameters with a given
[20]	2010	of an ML model	Query outputs	classifier in black-box setting to generate adversarial examples.
[34]	2015	Training data for	Input features & Query	Model inversion attacks used confidence info leaking
[54]	2015	an ML model	outputs & model structure	training samples with predictions against MLaaS in two settings.
[40]	2017	Training data for	Input features & Query	Online Attack using GAN against collaborative deep learning
[47]	2017	an ML model	outputs & model structure	model leaking user's training sample.
[116]	2017	Training data for	Input features &	Membership inference attacks use shadow training technique to
[110]	2017	an ML model	Query outputs	leak the specific record's membership of original training set.
[110]	2010	Training data for	Input features &	Enlarge the scope of membership inference attacks by releasing
[110]	2019	an ML model	Query outputs	some key assumptions.
[38]	2018	The property of	Input features & Query	Infer global properties of the training data unintended to be shared
[20]		training set	outputs & model structure	in white-box setting.

## 2.5 Evaluation (Cont.)

#### Table 2. Outline of Reviewed Papers

		-	•	-	
[91]	2019	The property of training set	Input features & Query outputs & model structure	Membership inference attacks against collaborative deep learning model leaking others' unintended feature.	
[95]	2018	Training data for an ML model	Input features & Query outputs	Protect against black-box membership inference attack using an adversarial training algorithm.	
[100]	2017	Training data for an ML model	Input features & Query outputs	Protect training set of model from leakage with teacher and student models using PATE.	
[70]	2017	Training data for an ML model	N/A	Protect training dataset in stored from leakage before training.	
[79]	2015	Input PINs; User input texts	Acoustic sensor data; Accelerometer data	Attack infers users' inputs on keyboards via accelerometer data within user's smartwatch.	
[122]	2016	Input PINs; User input texts	Audio sensor data	Attack infers a user's typed inputs from surreptitious video recordings of a tablet's backside motion.	
[42]	2018	Cryptographic keys	TLB Cache data	TLBleed attack TLBs to leak secret keys about victim's memory activities via reversing engineer and ML strategies.	
[152]	2016	Secret keys	CPU Cache data	Mitigate access-driven side-channel attacks with CacheBar managing memory pages cacheability.	
[132]	2016	Password info	PII & leaked password & site info	Attack with seven mathematical guessing models for seven password guessing scenario using different personal info.	
[128]	2014	Password info	Corpus & Site leaked list	Password guessing attack by analyzing its semantic patterns.	
[90]	2016	Password info	Corpus library	Mitigate against password guessing attack by modeling password guessability in password creation stage.	

#### **3.1 Controlled User Activities Information**

- It is essential for security specialists to protect user activities information. Not only because the private activities are valuable to adversaries.
- Also the adversary can exploit some specific activities to perform malicious attacks such as the phishing attack.



Fig. 4. The ML-based stealing attack against user activities information.

#### **3.1.1 Stealing controlled user activities from kernel data.**

- The dataset collected from the kernel about system process information is too noisy and coarse-grained to disclose any intelligible and valuable information.
- However, through analyzing plenty of such data, the adversary could deduce some confidential information about the victim's activities with the help of ML algorithms.
- *Stealing User Activities with Timing Analysis:* The security implications of the kernel information through integrating some specific hardware components into Android smartphones.
- The targeted user activities were unlock patterns and foreground apps. Moreover, users' browsing behavior was targeted by the attacker.

Reference	Dataset for Experiment	Description	Feature Engineering	ML-based Attack Method
[26]	Interrupt data for unlock pattern and for apps	Collect from procfsDeduplication; Interpolation; Interrupt Increment Computation; Gram Segmentation; DTW		HMM with Viterbi algorithm; <i>k</i> -NN classifier with DTW
[119]	Time series for apps, website,keyboard guests	Collect from <i>procfs</i>	Automatically extract with <i>tsfresh</i> ; DTW	Viterbi algorithm with DTW; SVM classifier with DTW
[151]	1200 x 6 time series of data about app; 1000 website traces	120 apps(App Store+iOS ) +10 trace x 6 time series; 10 traces for each website	Manually defined; SAX, BoP representation	SVM classifier; <i>k</i> -NN classifier with DTW
[136]	Consecutively reading data; Resident size field data	Collect from <i>procfs</i>	N/A; Construct a histogram binned into seven equal-weight bins	SVM classifiers

#### Table 3. Stealing Controlled User Activities using Kernel Data

**3.1.1 Stealing controlled user activities from kernel data (Cont.).** 

- *Stealing User Activities with iOS Side-channel Attack:* In iOS systems, one popular sidechannel attack vector of Linux system about the process information is inaccessible.
- **Protection using Privacy Mechanism:** An attack exploiting the kernel process information via decreasing the data's resolution was defended. A differential privacy (DP) mechanism was utilized to prevent the attackers from gaining any useful storage information.

#### 3.1.2 Stealing controlled user activities using sensor data.

- The stealing attack using sensor data should be studied seriously by the defenders, not only from the application of effective ML mechanisms, but also from the popularity of sensing enabled applications.
- The sensor information can reveal the controlled information indirectly as demonstrated in this stealing attack, such as acoustic and magnetic data.

#### 3.1.2 Stealing controlled user activities using sensor data (Cont.)

Table 4. Stealing Controlled User Activities using Sensor Data

Reference	Dataset for Experiment	Description	Feature Engineering	ML-based Attack Method
[50]	Audio signature dataset	Recorded with a phone put within 4 inches of the printer	STFT, noise normalization	A regression model
[117]	Sensor dataset	Sensor data collected benign and malicious activities	N/A	Markov Chain, NB, LMT, (alternative algorithms e.g. PART)

- *Stealing Machine's Activities with Sensor-based Attack:* A side-channel attack was proposed manufacturing equipment exploiting sensor data collected by mobile phones, which revealed its design and the manufacturing process.
- The attacker managed to reconstruct the targeted equipment. As a result of reconnaissance, the security threat of the manufacturing sector was indicated.
- After the dataset was gathered, the ML-based attack was completed by feature engineering, attacking with model training, and evaluation.

#### **3.1.3 Summary.**

- ML-based attacks steal user activities information from operating systems. According to the data sources, there are two kinds of attacks using kernel data and using sensor data.
- Kernel data reveals some system-level behaviors of the target system, while sensor data reflects the system's reactions on its specific functionality used by users [26]. The kernel data is analyzed by the adversary from a time dimension.
- *Countermeasures:* Regarding the protection mechanism, differential privacy is an important method for the attacks stealing user activities information.
- The in-depth research in protecting against user activities information can explore the differential privacy appliance or a management system design for kernel files and sensor data.
- Noise injection and access restriction are two effective protections, and the detection can alert the stealing attack.

#### **3.2 Controlled ML Model Related Information**

- ML model related information consists of the model description, training data information, testing data information, and testing results.
- The model description and training data information are controlled, otherwise, it is easy for an attacker to interpret the victim's query result.
- The generalized attack in this category is illustrated in Figure 5. In this category, ML-based attacks aim at stealing the training samples or the ML model.
- Stealing the controlled training sample attacks use an ML model to determine whether the input sample is contained in the target training set.



Fig. 5. The ML-based stealing attack against ML model related information.

#### **3.2.1 Stealing controlled ML model description.**

- It is important to protect the confidentiality of ML models online. If fraud detection are based on ML models then understanding the model means that adversaries can evade detection.
- A specific ML model is defined by two important elements including ML algorithm's parameters and hyperparameters. Since the model is controlled, its parameters and hyperparameters should be deemed confidential by nature.

Reference	Dataset for Evaluation	Description	Targeted ML Model	Attack Methods	
	Circles, Moons, Blobs,	Synthetic, 5,000 with 2 features,			
	5-Class [125]; Synthetic,1000 with 20 feature				
	Steak Survey [126], 331 records with 40 features,			Fountion solving	
	GSS Survey [118], 16,127 records with 101 features,		Logistic Regression;		
[105]	Adult (Income/race) [126],	48,842 records with 108/105 features,	Decision Tree;	attack; Path-finding attack	
[125]	Iris [126],	150 records with 4 features,	SVM;		
	Digits [107],	1,797 records with 64 features,	Three-layer NN		
	Breast Cancer [126],	683 records with 10 features,			
	Mushrooms [126],	8,124 records with 112 features,			
	Diabetes [126]	768 records with 8 features			
	MNIST [69]	70 000 handwritten digit images	DNN; SVM; <i>k</i> -NN;	Jacobian-based Dataset	
[101]	GTSRB [121]	49 000 traffic signs images	Decision Tree;		
	010100 [121]	47,000 traine signs images	Logistic Regression	Hughientation	
	Diabetes [126],	442 records with 10 features,			
	GeoOrig [126],	1,059 records with 68 features,	Regression algorithms:	Equation solving	
[131]	UJIIndoor [126];	19,937 records with 529 features;	Logistic regression		
	Iris [126],	100 records with 4 features;	algorithms: SVM: NN		
	Madelon [126],	4,400 records with 500 features;	algorithms, 5 v W, 1414		
	Bank [126]	45,210 records with 16 features			
[98]	MNIST [69]	70,000 handwritten digit images	NNs	Metamodel methods	

Table 5. Stealing Controlled ML Model Description

#### 3.2.2 Stealing controlled ML model's training data.

- Another type of controlled information about MLaaS product is the training data. Training data is not only useful to construct the model using ML algorithms provided by an MLaaS platform,
- Also sensitive as the records can be private information [34, 35]. Hence, the confidentiality of the model's training data should be protected.
- *Model Inversion Attack & Defense:* The model inversion attack was developed via conducting the commercial MLaaS APIs and leveraging confidence information with predictions.
- However, the attack aimed to be applicable across both white-box setting and black-box setting. For the white-box setting, an adversarial client had a prior knowledge about the description of the model as the APIs allowed.

#### 3.2.2 Stealing controlled ML model's training data (Cont.)

Table 6. Stealing Controlled ML Model's Training Data.

Reference	Dataset for Experiment	Description	Feature Engineering	ML-based Attack Method	
[24]	FiveThirtyEight survey,	553 records with 332 features,	NI/A	Decision Tree,	
[24]	GSS marital happiness survey	16,127 records with 101 features	IN/A	Regression model	
[40]	MNIST [69],	70,000 handwritten digit images,	Features learned	Convolutional Neural	
[49]	AT&T [111]	400 personal face images	with DNN	Network (CNN) with GAN	
	CIFAR10 [65],	6,000 images in 10 classes,			
	CIFAR100 [65],	60,000 images in 100 classes,			
	Purchases [52],	10,000 records with 600 features,	Regarded shadow model	NN	
[116]	Foursquare [140],	1,600 records with 446 features,	resulted as features and		
	Texas hospital stays [47],	10,000 records with 6170 features,	label records as in/out		
	MNIST [25],	10,000 handwritten digit images,			
	Adult (income) [126]	10,000 records with 14 attribute			
	Include 6 sets in [116],	Same as above cell,	Regarded shadow model	Random Forest,	
[110]	News [53],	20,000 newsgroup documents in 20 classes,	resulted as features and	Logistic Regression,	
	Face [68]	13,000 faces from 1,680 individuals	label records as in/out	Multilayer perceptron	
	Adult (income) [126],	299,285 records with 41 features,			
[29]	MNIST [69],	70,000 handwritten digit images,	Neuron sorting,	NN	
[30]	CelebFaces Attributes [80],	more than 200K celebrity images,	e than 200K celebrity images, Set-based representation		
	Hardware Performance Counters	36,000 records with 22 features			
	Face [68],	13,233 faces from 5,749 individuals,			
	FaceScrub [96], 76,541 faces from 530 individuals,			Logistic regression,	
[91]	PIPA [149],	60,000 photos of 2,000 individuals,	N/A	gradient boosting,	
	Yelp-health, Yelp-author [141],	17,938 reviews, 16,207 reviews,		Random Forests	
	FourSquare [140], CSI corpus [129]	15,548 users in 10 locations, 1,412 reviews			
	CIFAR100 [65],	60,000 images in 100 classes,	Regarded shadow model		
[95]	Purchase100 [52],	197,324 records with 600 features,	resulted as features and	NN	
	Texas100 [47]	67,330 records with 6,170 features	label records as in/out		

3.2.2 Stealing controlled ML model's training data (Cont.)

- *Stealing the Training Data of Deep Model with GAN:* An attack against the privacy-preserving collaborative deep learning was designed to leak the participants' training data which might be confidential.
- A distributed, federated, or decentralized deep learning algorithm can process each users' training set by sharing the subset of parameters obfuscated with differential privacy.
- However, the training dataset leakage problem had not been solved by using the collaborative deep learning model.
- *Membership Inference Attack:* Learning a specific data record which was the membership of the training set of the targeted MLaaS model.
- Since the commercial ML model only allowed black-box access provided by Google and Amazon, not only the training data but also the training data's underlying distribution were controlled.
- Therefore, an attack model was trained which could recognize such differences and determine whether the input data was the member of targeted training set or not. The attack is intended to recognize the model's behavior testing with target training sample.

#### 3.2.2 Stealing controlled ML model's training data (Cont.)

- **Property Inference Attack:** Different from learning a specific training record, the property inference attack targets at the properties of training data that the model producer unintended to share.
- The target model was defined as a white-box Fully Connected Neural Networks (FCNNs) with the aim to infer some global properties such as a higher proportion.
- During the feature engineering phase, the meta-training set by applying set-based representation instead of using a flattened vector of all parameters [5].

#### 3.2.2 Stealing controlled ML model's training data (Cont.)

Attack Type	Attack Targets		Attack Surfaces		Attacker's Capabilities	
Attack Type	Model Info	Training Set Info	Training Phase	Inference Phase	Black-box Access	White-box Access
Model extraction attack [125]	YES	no	no	YES	YES	no
Model extraction attack [101]	YES	no	no	YES	YES	no
Hyperparameter stealing attack [131]	YES	no	no	YES	YES	no
Hyperparameter stealing attack [98]	YES	no	no	YES	YES	no
Black-box inversion attack [34]	no	YES	no	YES	YES	no
White-box inversion attack [34]	no	YES	no	YES	no	YES
GAN attack [49]	no	YES	YES	no	no	YES
Membership inference attack [116]	no	YES	no	YES	YES	no
Membership inference attack [110]	no	YES	no	YES	YES	no
Property inference attack [38]	no	YES	no	YES	no	YES
Property inference attack [91]	no	YES	YES	no	no	YES

Table 7. Categories of Stealing ML related information attacks from three perspectives.

- As for attack targets, two types of information may be stolen model internal information and training set information. From attack surfaces, attacks may occur during either model's training phase or inference phase.
- Considering the attacker's capability, the ML model usually allows either the black-box access or the white-box access. The first category is used for this subsection's organization.

#### 3.2.2 Stealing controlled ML model's training data (Cont.)

Table 8. Attack's prior knowledge under black-box access and white-box access.

Model's Information	Black-box Access	White-box Access
Predicted Label	YES	YES
Predicted Confidence	YES	YES
Parameters	NO	YES
Hyperparameters	NO	YES

- The black-box access allows the users to query the model and obtain prediction outputs which include the predicted label and confidence value.
- The white-box access allows the users to access any information of its model which includes predicted label, predicted confidence, parameters, and hyperparameters.

#### **3.2.3 Summary.**

- ML-based stealing attacks against model related information target at either model descriptions or model's training data.
- Attackers steal model's training data mostly at inference phase, except the GAN attack [49] and the property inference attack [91] which happen at training phase of collaborative learning.
- *Countermeasures:* Concerning the ML pipeline, the protection methods will be applied in data preprocessing phase, training phase, and inference phase, respectively.
- Differential privacy is the most common countermeasure to defend against the stealing attack, however, it alone cannot prevent the GAN attack. Privacy, regularization, dropout, and rounding techniques are popular protections at the training and inference phases.



Fig. 6. The ML-based stealing attack against authentication information

- Keystroke information and secret keys. After reconnoitering and querying, attackers targeting at keystroke information and secret keys interact with the target system to collect data, which refers to the active collection.
- The attack involved active collection shares a similar workflow as Fig. 4 depicted.

들어주서서 감사합니다! Sanote

# Thank You

# For your Attention!

연락처: yotxaysangthong@seoultech.ac.kr +82 10-8999-3151

